

**EPTCS 395**

Proceedings of the  
**Fifth International Workshop on  
Formal Methods for Autonomous Systems**

**Leiden, The Netherlands, 15th and 16th of November 2023**

Edited by: Marie Farrell, Matt Luckcuck, Mario Gleirscher and Maike  
Schwammberger

Published: 15th November 2023  
DOI: 10.4204/EPTCS.395  
ISSN: 2075-2180  
Open Publishing Association

## Table of Contents

Table of Contents .....	i
Preface .....	iii
<i>Matt Luckcuck, Marie Farrell, Mario Gleirscher and Maike Schwammberger</i>	
<b>Invited Presentation:</b> Formal Methods within the TAS Governance Node .....	vi
<i>Alice Miller</i>	
<b>Invited Presentation:</b> SMT: Something you Must Try .....	vii
<i>Erika Ábrahám</i>	
What to tell when? – Information Provision as a Game .....	1
<i>Astrid Rakow, Mehrnoush Hajnorouzi and Akhila Bairy</i>	
Trust Modelling and Verification Using Event-B .....	10
<i>Asieh Salehi Fathabadi and Vahid Yazdanpanah</i>	
Comparing Differentiable Logics for Learning Systems: A Research Preview .....	17
<i>Thomas Flinkow, Barak A. Pearlmutter and Rosemary Monahan</i>	
Extending Neural Network Verification to a Larger Family of Piece-wise Linear Activation Functions .....	30
<i>László Antal, Hana Masara and Erika Ábrahám</i>	
Certified Control for Train Sign Classification .....	69
<i>Jan Roßbach and Michael Leuschel</i>	
Model Checking for Closed-Loop Robot Reactive Planning .....	77
<i>Christopher Chandler, Bernd Porr, Alice Miller and Giulia Lafratta</i>	
Online Reachability Analysis and Space Convexification for Autonomous Racing .....	95
<i>Sergiy Bogomolov, Taylor T. Johnson, Diego Manzananas Lopez, Patrick Musau and Paulius Stankaitis</i>	
Automatic Generation of Scenarios for System-level Simulation-based Verification of Autonomous Driving Systems .....	113
<i>Srajan Goyal, Alberto Griggio, Jacob Kimblad and Stefano Tonetta</i>	
Enforcing Timing Properties in Motorway Traffic .....	130
<i>Christopher Bishopink</i>	
Correct-by-Construction Control for Stochastic and Uncertain Dynamical Models via Formal Abstractions .....	144
<i>Thom Badings, Nils Jansen, Licio Romao and Alessandro Abate</i>	

Towards Formal Fault Injection for Safety Assessment of Automated Systems.....	153
<i>Ashfaq Farooqui and Behrooz Sangchoolie</i>	
Formal Verification of Long Short-Term Memory based Audio Classifiers: A Star based Approach .	162
<i>Neelanjana Pal and Taylor T Johnson</i>	
3vLTL: A Tool to Generate Automata for Three-valued LTL.....	180
<i>Francesco Belardinelli, Angelo Ferrando and Vadim Malvone</i>	
Towards Proved Formal Specification and Verification of STL Operators as Synchronous Observers	188
<i>Céline Bellanger, Pierre-Loïc Garoche, Matthieu Martel and Célia Picard</i>	
Runtime Verification of Learning Properties for Reinforcement Learning Algorithms .....	205
<i>Tommaso Mannucci and Julio de Oliveira Filho</i>	

## Preface

This EPTCS volume contains the papers presented at the Fifth International Workshop on Formal Methods for Autonomous Systems (FMAS 2023), which was held on the 15th and 16th of November 2023. FMAS 2023 was co-located with 18th International Conference on Integrated Formal Methods (iFM'23), organised by Leiden Institute of Advanced Computer Science of Leiden University. The workshop itself was held at Scheltema Leiden, a renovated 19<sup>th</sup> Century blanket factory alongside the canal.

The goal of the FMAS workshop series is to bring together leading researchers who are using formal methods to tackle the unique challenges that autonomous systems present, so that they can publish and discuss their work with a growing community of researchers. Autonomous systems are highly complex and present unique challenges for the application of formal methods. Autonomous systems act without human intervention, and are often embedded in a robotic system, so that they can interact with the real world. As such, they exhibit the properties of safety-critical, cyber-physical, hybrid, and real-time systems. We are interested in work that uses formal methods to specify, model, or verify autonomous and/or robotic systems; in whole or in part. We are also interested in successful industrial applications and potential directions for this emerging application of formal methods.

We continued to hold FMAS as a hybrid event this year. The workshops in 2020 and 2021 have been fully online events because of the restrictions required to deal with the COVID-19 pandemic, and FMAS 2022 continued to facilitate online participation. We feel that a hybrid event, while often challenging to organise, provides accessibility to people not able to travel for the workshop. FMAS 2023 had both presentations and attendees who were remote, and we hope that it was a useful option for those people who made use of it.

FMAS 2023 continued to use the submission categories introduced last year: *vision papers* and *research previews*, both of which were types of short paper; and *experience reports* and *regular papers*, both of which were types of long paper. In total, FMAS 2023 received 25 submissions. We received 11 regular papers, 3 experience reports, 6 research previews, and 5 vision papers. The researchers who submitted papers to FMAS 2023 were from institutions in: Australia, Canada, Colombia, France, Germany, Ireland, Italy, the Netherlands, Sweden, the United Kingdom, and the United States of America. Increasing our number of submissions for the third year in a row is an encouraging sign that FMAS has established itself as a reputable publication venue for research on the formal modelling and verification of autonomous systems. After each paper was reviewed by three members of our Programme Committee, we accepted a total of 15 papers: 8 long papers and 7 short papers.

FMAS 2023 hosted two invited speakers. Prof. Alice Miller, from the University of Glasgow (UK), gave a talk titled "*Formal Methods within the TAS Governance Node*"; which focussed on several strands of work using formal methods within the Trustworthy Autonomous Systems Hub's Governance Node. Prof. Erika Ábrahám, from RWTH Aachen University (Germany), gave a talk titled "*SMT: Something you Must Try*"; which discussed using Satisfiability Modulo Theories (SMT) to solve real-world problems. Prof. Ábrahám's talk was held in a joint session with iFM.

This is the fifth year that we have held an FMAS workshop; five years seems like a number that is round enough to feel significant and warrant some reflection. The idea for FMAS came from conversations between two of the current organisers, Drs Marie Farrell and Matt Luckcuck, and Prof. Michael Fisher during a research project on using robotics and AI in the UK nuclear industry. FMAS 2019 was held as a satellite workshop at FM 2019 in Porto, Portugal. We were unsure if this would be a one-off event or an ongoing series, and there were only five papers presented at that first workshop. After the uncertainty of running workshops during the pandemic, FMAS has rebounded strongly with several years of increasing numbers of submissions. This year, FMAS received five more submissions than last

year, giving us our highest-ever number number of submissions. This is a very pleasing pay-off to five years of hard work, and reflects the progress of this burgeoning community in tackling the challenges that autonomous systems pose for formal modelling and verification techniques.

For four years, FMAS has been organised almost entirely by Drs Matt Luckcuck and Marie Farrell in equal partnership, but it was becoming clear that this would be unsustainable as the workshop grows bigger. To help spread the workload, and to inject some fresh perspective, two colleagues joined the FMAS organising committee this year. Matt and Marie were joined by Jun.-Prof. Dr Maike Schwammburger and Dr Mario Gleirscher, who have been supporting FMAS 2023 in preparation for a bigger restructuring from next year's workshop onward. We are already planning for what FMAS will look like for the next five years and we hope that the improvements that are coming will make FMAS an event that is even more useful and enjoyable for our community.

We would like to thank our brilliant Programme Committee and sub-reviewers for their helpful reviews and discussions behind the scenes. Many of the reviewers for FMAS 2023 have been part of our Programme Committee since the first FMAS workshop, we are pleased to have their continuing support. Whether this is their first time or their fifth, we are proud to have a community of reviewers who are so enthusiastic and supportive of our workshop and the work it receives. We thank them for volunteering their time and effort because without them we could not produce our programme of presentations. We also thank our invited speakers, Prof. Alice Miller and Prof. Erika Ábrahám, for their time; the authors who submitted papers; our EPTCS editor, Martin Wirsing, for overseeing the preparation of the proceedings; the organisers of iFM – Marcello M. Bonsangue, Paula Herber, and Anton Wijs – for supporting our workshop; FME for its sponsorship of our student travel grants; and all of the attendees (both virtual and in-person) of FMAS 2023. We hope to see you all at FMAS 2024.

**Matt Luckcuck, Marie Farrell, Mario Gleirscher, and Maike Schwammburger**

November 2023

## Program Committee

Oana Andrei	University of Glasgow (UK)
Akhila Bairy	Carl von Ossietzky University of Oldenburg (Germany)
Christopher Bischofink	Carl von Ossietzky University of Oldenburg (Germany)
Rafael C. Cardoso	University of Aberdeen (UK)
Louise A. Dennis	University of Manchester (UK)
Marie Farrell	University of Manchester (UK)
Fatma Faruq	ETAS – Empowering Tomorrow’s Automotive Software (UK)
Angelo Ferrando	University of Genova (Italy)
Michael Fisher	University of Manchester (UK)
Mario Gleirscher	University of Bremen (Germany)
Mallory S. Graydon	NASA Langley Research Center (USA)
Ichiro Hasuo	National Institute of Informatics (Japan)
Taylor T. Johnson	Vanderbilt University (USA)
Verena Klös	Technical University of Dresden (Germany)
Matt Luckcuck	University of Nottingham (UK)
Raluca Lefticaru	University of Bradford (UK)
Lina Marsso	University of Toronto (Canada)
Anastasia Mavridou	NASA Ames Research Center (USA)
Claudio Menghi	University of Bergamo (Italy)
Alice Miller	University of Glasgow (UK)
Alvaro Miyazawa	University of York (UK)
Rosemary Monahan	Maynooth University (Ireland)
Yvonne Murray	University of Agder (Norway)
Dominique Méry	Université de Lorraine (France)
Natasha Neogi	NASA Langley Research Center (USA)
Colin Paterson	University of York (UK)
Baptiste Pelletier	ONERA – The French Aerospace Lab (France)
Andrea Pferscher	Graz University of Technology (Austria)
Maike Schwammberger	Karlsruhe Institute of Technology (Germany)
James Stovold	Lancaster University Leipzig (Germany)
Silvia Lizeth Tapia Tarifa	University of Oslo (Norway)
Elena Troubitsyna	KTH Royal Institute of Technology (Sweden)
Gricel Vázquez	University of York (UK)
Hao Wu	Maynooth University (Ireland)
Mengwei Xu	University of Newcastle (UK)

## Subreviewers

Qais Hamarneh	Karlsruhe Institute of Technology (Germany)
Thomas Flinkow	Maynooth University (Ireland)

# Invited Talk: Formal Methods within the TAS Governance Node

Alice Miller

University of Glasgow, Glasgow, UK

The TAS Governance Node part of the £33M Trustworthy Autonomous Systems Programme funded by the UKRI Strategic Priorities Fund. The aim of the node is to explore how to make autonomous systems aware of — and responsive to — changing regulations. Led by the University of Edinburgh, it brings together researchers from the universities of Edinburgh, Glasgow, Nottingham, Heriot-Watt, Sussex, and Kings College London; as well as multiple industrial partners. In this talk I will highlight some of the activities within the node, focussing on those that use Formal Methods. These include:

- an automatic theory repair system for a legal responsibility framework for autonomous vehicles,
- accident anticipation through reasoned simulation,
- robot planning using in-situ model checking; and,
- (formal aspects of) a node-wide automotive case study investigation.



# Invited Talk: SMT: Something you Must Try

Erika Ábrahám

RWTH Aachen University, Aachen, Germany

SMT (Satisfiability Modulo Theories) solving is a technology for the fully automated solution of logical formulas. Due to their impressive efficiency, SMT solvers are nowadays frequently used in a wide variety of applications. These tools are general purpose and as off-the-shelf solvers, their usage is truly integrated. A typical application encodes real-world problems as logical formulas, whose solutions can be decoded to solutions of the real-world problem. In this talk we give some insights into the mechanisms of SMT solving, discuss some areas of application, and present a novel application from the domain of simulation.



# What to tell when? – Information Provision as a Game

Astrid Rakow

Inst. of Systems Engineering for Future Mobility,  
German Aerospace Center (DLR) e.V.  
astrid.rakow@dlr.de

Mehrnoush Hajnorouzi

mehrnoush.hajnorouzi@dlr.de

Akhila Bairy

Dept. of Computing Science,  
Carl von Ossietzky University  
akhila.bairy@uni-oldenburg.de

*Constantly informing systems (CIS)*, which provide us with information over a long period of time, face a particular challenge in providing useful information: Not only does a person’s information base change but also their mood or abilities may change. An information provision strategy should hence take these aspects into account. In this paper, we describe our vision for supporting the design of *CIS*. We envisage using psychological models of the human mind and emotions in a game of information provision. The analysis of these games will give comparative insights into design variants.

## 1 Introduction

Complex technical systems are deeply integrated into our daily lives. They provide us with various services ranging from entertainment to safety-critical services. We use many of these systems over long periods of time and they even adapt to our needs. Such systems and their users exchange information: On one hand, the system acquires user information to adapt. On the other hand, users often need information from the system e.g. some explanation of the system’s behaviour. When such an interaction between the system and its user extends over a longer period of time, an important question therefore arises: what content should be communicated to the user and when? A few brief examples will illustrate the challenges of this question: (1) An autonomous vehicle *AV* initiates emergency braking due to a deer suddenly crossing the road. (2) The *AV* slows down on the way to work due to a dumpy road. (3) The *AV* slows down due to a carnival procession, the user is accompanied by his excited kids. In these scenarios, the *AV* faces the challenge of determining what content helps the user in the current context. It should factor in that the time span is very short in (1), that the driver has prior knowledge in (2), and the right level of detail is especially important, if the driver is occupied otherwise as in (3).

In this paper, we argue that the design of a *constantly informing system (CIS)*<sup>1</sup> can profit from a holistic game theoretic analysis that places a strong emphasis on the human aspect. We propose to analyze the challenge of information provision (*What content should be communicated to the user when?*) by methods of economics to determine which game is to be played (Sect. 2). In the resulting game, we consider the human as a system of limited cognitive resources (Sect. 4). In order to do so, we propose to use models of the human mind and emotions from psychology, which we briefly survey in Sect. 3. We outline how we envision integrating these models into games between a *CIS* and its user (cf. Fig. 1). Their analysis results will allow designers comparing variants of a *CIS* under development.

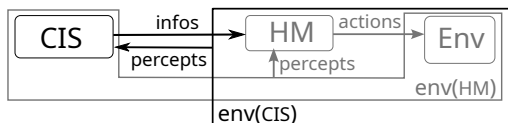


Figure 1: *CIS* observes its environment and provides information to a human model (*HM*). *HM* perceives this information and observes its surrounding *Env* to choose its actions. The goal of *CIS* is to enable *HM* achieve its goals.

<sup>1</sup>Self-explaining systems are also *CIS*. A *CIS* is not necessarily self-explaining.

## 2 The Challenge of Information Provision Seen as a Game

Game theory, developed by von Neumann and Morgenstern [21] in the mid-20th century, studies interactive decision-making where the outcome for each player depends on the actions of the other players [20]. It has been used in various domains such as economics, political science or computer science. In the following, we discuss how the challenge of information provision can be considered a game. By reinterpreting the economic analysis method PARTS [6] for information provision, we outline the broader scale of this challenge and then narrow the focus down to the next steps of our research.

PARTS defines a list of guiding questions to analyze a situation as a game: the **players** (*who is active? who has goals?*), **added value** (*what does player X gain by Y's participation?*), **rules** (*what rules do constrain the players? who can change the rules?*), **tactics** (*how can the players' perception of the game be changed?*) and **scope** (*what are the boundaries of the game?*). PARTS is only a first step towards the definition of a game. The concrete game depends on the characteristics of the application domain, the information available to build a model of the application domain, and the goal of the analysis.

We want to analyze how the “*usefulness of a CIS*” can be optimized. *Useful* information is relevant information<sup>2</sup> that the customer can *process and use* to accomplish their goals, i.e. the right content at the right time. It is part of our future work to concretize this notion. Table 1 summarizes how we interpret the **players'** roles in an information provision setting. While the **added values** that the players will bring to the game vary depending on the system considered, we can identify two main types of rules as shown in Table 2. **Tactics** refer to actions that a player can take to change how other players perceive the game, i.e. to change their belief in who the players are, what their added values and what the rules are. The scope of *our* game is a *CIS* that is a subsystem of an *AV* and serves a single human (the driver) who is captured via a human model *HM*. For a start, we use a game of this very limited scope to investigate how to formally include human capabilities in the design of a *CIS*.

Customers	all recipients of information (humans or technical systems)
Suppliers	suppliers of information (sensors and other <i>CIS</i> ) suppliers of computing services (hardware resources)
Competitors	systems that render the <i>CIS</i> 's information redundant (other <i>CIS</i> , human perception) processes that use the same computing services or human resources ( <i>e.g. music system</i> )
Complementors	other <i>CIS</i> whose information enhances the <i>CIS</i> ' information ( <i>e.g. traffic warnings plus detour</i> )

Table 1: Interpretation of player roles in the “information provision” game

guiding rules	a player is obliged to follow these rules, but can break them at will ( <i>e.g. ethical, legal rules</i> )
constraint rules	a player cannot break these rules; the rules can only be changed by a third party (resulting from the technical implementation)

Table 2: Types of rules in the “information provision” game

As the above exposé shows, information provision is an interdisciplinary challenge where humans and technical systems interact – many *CIS* are embedded into a human cyber-physical system of systems. Legal and ethical aspects, as well as human capabilities, influence what a *CIS* should provide as useful information. Whether information is useful strongly depends on whether the information is relevant *and* on whether humans can make use of it.

<sup>2</sup>Here, relevant information for the *CIS*'s customer is the minimal information that they need to accomplish their goal [26].

While the above was concerned with analyzing the application domain, we will now discuss what game  $G$  in the “catalog” of well-studied games and respective analysis methods is appropriate for our analysis goal. We will take a closer look at the characteristics of the application domain and the goal of the analysis – keeping in mind what aspects are most important and which aspects might be omitted.

In a nutshell, the goal of  $G$ 's analysis is to determine which content our  $CIS$  should provide to a  $HM$  in what context. This analysis is meant to be done at design time in order to derive requirements on the information provision and to compare different design alternatives. Our  $CIS$  controls what information content it provides when to  $HM$ . We discuss what kinds of games are appropriate using the following guiding questions based on the taxonomy of games as developed in [19]: *Is CIS acting strategically? What is the relation between HM and CIS? Who knows what?*

*Is CIS acting strategically?* We will study a  $CIS$ , that *chooses rationally* the content that seems the most useful. We hence will develop a classical game rather than an evolutionary game.

*What is the relation between HM and CIS? Should we use cooperative or non-cooperative game theory?* In cooperative game theory ( $CGT$ ), the players cooperate with each other and form coalitions to obtain maximum reward/payoff.

In non-cooperative game theory ( $NCGT$ ), the players make decisions independently and each player aims to maximise their own payoff. We assume that  $HM$  has a prioritized list of goals.  $CIS$  has only the goal to be helpful within its constraints thus acting subordinately<sup>3</sup>.

In such a subordinate-dominant relationship, both  $CGT$  as well as  $NCGT$  could be a fit. In order to apply  $NCGT$  we envision to model  $HM$  as always greedily accepting all information that it can in its current mode; *E.g. if HM is tired, less information per time is accepted*. We thereby encode a kind of cooperation into  $HM$ 's modes, so that it is not playing against  $CIS$  by maliciously rejecting all offered information. To derive these modes and the dynamics of mode changes we plan to use psychological models of the human mind and emotions. We will give a brief overview of these in Sect. 3. In Sect. 4 we will give technical details on how we envision using such psychological models to derive the  $HM$  for our game.

*Who knows what?* Both  $CIS$  and  $HM$  receive information from the environment. In particular,  $HM$  receives information from our  $CIS$  within its car but also via its five senses that are competing with  $CIS$ .  $CIS$  also does not exactly know how  $HM$  reacts, as it has no precise means to derive the state and reactions of  $HM$ . From  $CIS$ 's point of view  $HM$ 's reactions expose some degree of randomness.  $HM$  is a system of bounded resources and will usually not (have the resources to) think much about the beliefs of any support systems in its car. We hence will define an epistemic game, where  $CIS$  has a belief on  $HM$ 's belief but  $HM$  does not build beliefs on  $CIS$  beliefs.

$CIS$  knows what the possible inputs of  $HM$  are, but assumes that these channels are unreliable. That is, we assume that  $CIS$  can roughly predict what information  $HM$  is able to receive via its senses.  $CIS$  moreover knows what actions  $HM$  can choose (what a driver can do).

To us, the central challenge in the information provision game is capturing the human factors. We envision to derive instantiations of  $HM$  from models of the human mind and emotions. These instantiations are then “plugged” into our game. We then check whether our  $CIS$  can inform  $HM$ , i.e. whether a strategy exists, such that  $HM$  achieves its goals.

---

<sup>3</sup>It is out of this paper's scope to discuss the case when goals of  $HM$  contradict the constraints of  $CIS$ .

### 3 Models of Human Mind

For an automated system to provide useful information, it needs a human behaviour model to determine the impact of the content and the optimal time for providing the content. Well-crafted explanations can reassure users and increase alertness [12]. The impact of percepts and its timing –or more generally the dynamics of the human mind– is captured by cognitive architectures. These are well-established frameworks developed in the field of cognitive psychology [22, 2]. They simulate how knowledge is processed, stored, and utilized in response to external stimuli, perception, and knowledge [25]. There are many cognitive architectures with different foci, like ICON FLUX[23], CASCaS [32], PRIM[28], but none of them captures the human holistically.

In addition to modelling core cognitive functions (such as perception, learning, memory, etc.), a few cognitive architectures can also capture the interplay of emotions and behaviour. These architectures generally adopt two primary approaches to model emotion-induced behaviours: altering architecture structures (e.g. modifying memory modules), or modifying the processing parameters (e.g. adjusting retrieval delay). Emotions act as conduits to human behaviour, significantly impacting decisions and influencing the daily choices of individuals [4]. A change in the person-environment relationship [15] may trigger emotional shifts, which may change this relationship in future.

Over time, different theories of emotion have evolved to explain the functioning of human emotions [24, 14, 11]. The cognitive appraisal theory is one such theory. According to this theory, emotions are influenced by individuals' cognitive evaluations and interpretations of situations rather than by the situations themselves. Many computational models are built on this theory [17, 5, 16]. Implementation of the appraisal process within the context of cognitive architectures is an ongoing effort. MAMID [10] is one example. It processes incoming stimuli through various modules. The affect appraiser module derives emotional states. The resulting affective state influences the cognition processes by adjusting the rules related to goal selection, action determination, module speed and capacity, and mental construct prioritization.

The appraisal process provides an assessment of the current situation, influenced by various parameters, such as individual history, personality, and current affective state. These parameters can be described by modelling the effects of transient states (emotions) and personality traits (characteristics).

To depict how the current situation affects human behavior, [33] incorporated a situation assessor module into their proposed SAMPLE<sup>4</sup> pilot model architecture. This module translates situation descriptions into types that trigger specific responses. The module employs belief nets to hold required knowledge, with nodes denoting particular features and events, and links illustrating causal and correlational connections between them. Furthermore, weighted links denote transition probabilities.

Appraisal theory is also utilised in the field of affective computing to enhance systems' abilities to perceive and generate human emotions. Affective computing seeks to make computers more human-like by enabling them to recognise, interpret, predict, and respond to human emotions. While few models effectively integrate emotion generation and its impact on cognitive processes, EmoCog [16] stands out. This work adeptly fuses diverse models of computational emotions and cognitive architectures into a unified model. In cognitive architectures, emotions can be divided into two parts: emotion effects and emotion generation [16]. Emotion effects consider how emotions might influence cognitive processing, while emotion generation explores how cognitive processes contribute to the generation of emotion. EmoCog, developed by Lin et al., also has primary and secondary appraisals similar to WASABI, developed by Becker-Asano and Wachsmuth [5]. However, EmoCog's secondary appraisal is more versatile,

---

<sup>4</sup>Situation Awareness Model for Pilot-in-the-Loop Evaluation

explaining a broader range of emotions than WASABI. The model includes an appraisal module that adjusts arousal and valence values in short-term memory nodes.

A psychological model of a human can also be achieved by their attention assessment. Wickens et al. developed SEEV<sup>5</sup> model to quantify and predict a person’s attention level across various areas of interest in a given situation [30]. Initially intended for pilots’ attention prediction, [9] and [31] later adapted SEEV to gauge driver attention in road traffic. Leveraging SEEV, Bairy and Fränze crafted a model that predicts the optimal time to deliver explanations to drivers based on their attention level [3].

## 4 Psychological Models in A Game

In a nutshell, we want to define a light-weight game as sketched in Fig. 1. *CIS* plays against its environment env(*CIS*), which contains the human captured via *HM*. As discussed in Sect. 2 we encode into *HM* how a human reacts so that we do not need to consider the human as a separate player. *CIS*’s actions are providing information to *HM*. The *CIS* model can reflect that information can only be provided if *CIS* has the information in its information base and if it has the required computing resources and time to retrieve it. Strategy synthesis will then tell a designer whether *CIS* can provide information such that *HM* achieves its goals<sup>6</sup>. A designer can plug in different variants of *CIS* and compare how well they perform in terms of the goals that can be achieved.

We aim for a *HM* model that encodes how a human reacts to the provided information or other environmental perceptions. We assume that *HM* reacts internally, i.e. emotionally, by consuming cognitive resources or by changing its information base. This internal change might cause a mode change of *HM*. In the new mode *HM* chooses a different strategy to reach its goals. We thus think of *HM* as sketched in Fig. 2. We want to apply strategy synthesis to examine the information provision game of *CIS* (cf. Fig. 1) that uses such a *HM* as part of its environment model. There, a dominant strategy would be the best strategy of *CIS* to provide the *HM* with the necessary information for *HM* to achieve its goals.

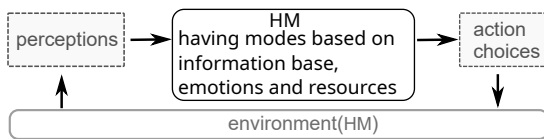


Figure 2: *HM* perceives its environment and chooses actions. During a scenario *HM* changes its mode, reflecting a change in its cognitive resources, emotions, and information.

A central idea of our approach is to consider *HM* as a resource-bounded system that reacts mode-specifically to the provided information content. How *HM* reacts will be derived from psychological models that are instantiated for the considered scenario. We consider the central challenges of making these models usable in our game to be: *attaining a formal representation of psychological models, reducing their complexity, combining different psychological theories, and interpreting the results.*

As outlined in Sect. 3, there is a broad foundation of psychological theories that are shedding light on different aspects of how humans react to provided information. The level of precision of these theories varies heavily as well as their level of formalization, ranging from informal over computational to formal. As Reisenzein puts it “Theoretical fragmentation and a comparative lack of precision are [...] characteristic of psychology in general” [27, p.248].

We focus on executable cognitive architectures as a starting point for our work, as they have an appropriate level of precision. But to use such architectures within a formal game, we need a formal

<sup>5</sup>Saliency, Effort, Expectancy, and Value

<sup>6</sup>This list of goals varies with the considered scenario but we assume that the goal of being safe has highest priority while the goal of feeling well has less priority.

representation of them. The approach of Langenfeld et al. [13] for deriving timed automata for an ACT-R instance promises a detailed study of the human reactions as captured in ACT-R. Since our precursor work [3, 8] showed that the resulting models easily get too complex for our needs, we want to derive simplified human models  $HM_i$  from cognitive architectures by automata learning techniques [1]. Automata learning techniques have been applied in many areas, such as speech recognition, software development, and computational biology [29]. In passive learning, the learning algorithm observes the inputs and the corresponding outputs of the system that is to be learned. Active learning, on the other hand, adopts an interactive approach where the learner also actively asks strategic queries and observes the response. We want to use modes as a means for guiding how the human models get simplified. Our idea is that a mode roughly specifies how  $HM$  reacts, i.e. what kind of strategy  $HM$  will apply. We think this might be done in a similar way to [7]. There Gosh and Verbrugge captured the kinds of reasoning strategies of player in a marble drop game as logical formulas, translated them into cognitive models and executed them in the PRIMs architecture [28]. Instead of player types as in [7] we think of player modes along the line “A tired driver reacts slowly. A driver gets tired when they are bored over a period of time.” or “A driver that already has to keep seven facts in mind, will not easily remember five more facts.” [18] We then want to apply automata learning to the cognitive architecture to derive a simplified model  $HM$ . There the choice of observable propositions will be an important influence on whether meaningful modes can be derived. Since the different architectures have different foci, we plan to use several architectures.

Given we would have  $n$  appropriately light-weight models  $HM_i$  –derived from different cognitive architectures or manually designed–, we build  $n$  games  $G_i$ ,  $1 \leq i \leq n$  (cf. Fig. 3). The performance

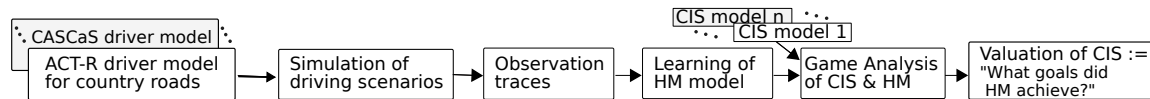


Figure 3: From cognitive architectures to a comparison of *CIS* variants.

scores of the same *CIS*  $C_j$  in the different games  $G_i$  –measured e.g. in terms of  $HM$  reaching its goals– will lead to a  $n$ -tuple  $t_j$  of scores for  $C_j$ . The comparison of two *CIS* variations,  $C_j$  with performance  $t_j$  and  $C_k$  with performance  $t_k$ , will help the designer to get early on insight into how well the *CIS* is tailored to human needs. But the comparison of  $C_j$  and  $C_k$  has to be done with care, as the meaning of the  $HM_i$  models has to be taken into account. This comparison gets especially difficult if the scores relate to interdependent aspects. What if  $HM_h$  suffers from a high workload while  $HM_i$  is in a good mood? Due to the coarse granularity of models that we want to consider, we expect that such cases will often occur, so that plausibility checks and if necessary refining analysis has to be done. To this end, we think of spotlight approaches where the learner gets more fine-grained observations where necessary. So, if this approach is in use, a collection of appropriate  $HM_i$  models could be evolved over time.

## 5 Conclusion

In this paper, we outlined our vision of how game theory can be used to support the development of constantly informing systems (*CIS*), taking into account the human factor. Our vision is to use models of the human mind and emotions from psychology to derive coarse formal human models,  $HM_i$  that can be used in a game between the *CIS* and the human captured by  $HM_i$ . The main challenges of our vision are attaining a formal representation of psychological models, reducing their complexity, combining different psychological theories, and interpreting the results. Nevertheless, we hope to obtain at least comparative results that can guide the early design of *CIS* in tailoring to the human needs.



## A Acknowledgements

This research has been supported by the German Research Council (DFG) in the PIRE Projects SD-SSCPS and ISCE-ACPS under grant no. DA 206/11-1, by RTG SEAS, and by the German Federal Ministry of Education and Research (BMBF) within the project "ASIMOV-D" under grant agreement No. 01IS21022G, based on a decision of the German Bundestag.

## B Abbreviations

AV	Autonomous Vehicle
CIS	Constantly Informing System
CGT	Cooperative Game Theory
G	Game
HM	Human Model
NCGT	Non-Cooperative Game Theory

## References

- [1] Bernhard K. Aichernig, Wojciech Mostowski, Mohammad Reza Mousavi, Martin Tappler & Masoumeh Taromirad (2018): *Model Learning and Model-Based Testing*, pp. 74–100. Springer International Publishing, Cham, doi:10.1007/978-3-319-96562-8\_3.
- [2] John R. Anderson & Christian Lebiere (1998): *The Atomic Components of Thought*. Lawrence Erlbaum associates, doi:10.4324/9781315805696.
- [3] Akhila Bairy & Martin Fränzle (2023): *Optimal Explanation Generation Using Attention Distribution Model. Human Interaction and Emerging Technologies (IHET-AI 2023): Artificial Intelligence and Future Applications 70(70)*, doi:10.54941/ahfe1002928.
- [4] Antoine Bechara, Hannah Damasio & Antonio Damasio (2000): *Emotion, Decision Making and the Orbitofrontal Cortex. Cerebral cortex (New York, N.Y. : 1991)* 10, pp. 295–307, doi:10.1093/cercor/10.3.295.
- [5] Christian Becker-Asano & Ipke Wachsmuth (2010): *Affective computing with primary and secondary emotions in a virtual human. Autonomous Agents and Multi-Agent Systems* 20, pp. 32–49, doi:10.1007/s10458-009-9094-9.
- [6] Adam M. Brandenburger & Barry J. Nalebuff (1995): *The Right Game: Use Game Theory to Shape Strategy. (Cover story). Harvard Business Review* 73(4), pp. 57 – 71.
- [7] Sujata Ghosh & Rineke Verbrugge (2018): *Studying strategies and types of players: experiments, logics and cognitive models. Synthese* 195(10), pp. 4265–4307, doi:10.1007/s11229-017-1338-7.
- [8] Mehrnoush Hajnorouzi & Martin Fränzle: *Model-based Control for Human-Centered Systems*. Unpublished.
- [9] William J Horrey, Christopher D Wickens & Kyle P Consalus (2006): *Modeling Drivers' Visual Attention Allocation While Interacting With In-Vehicle Technologies. Journal of experimental psychology. Applied* 12(2), pp. 67–78, doi:10.1037/1076-898X.12.2.67.
- [10] Eva Hudlicka (2002): *This Time with Feeling: Integrated Model of Trait and State Effects on Cognition and Behavior. Applied Artificial Intelligence* 16, pp. 1–31, doi:10.1080/08339510290030417.
- [11] Carroll E Izard (1992): *Basic emotions, relations among emotions, and emotion-cognition relations.* doi:10.1037/0033-295x.99.3.561.
- [12] Jeamin Koo, Dongjun Shin, Martin Steinert & Larry Leifer (2016): *Understanding driver responses to voice alerts of autonomous car operations. International Journal of Vehicle Design* 70, p. 377, doi:10.1504/IJVD.2016.076740.

- [13] Vincent Langenfeld, Bernd Westphal & Andreas Podelski (2019): *On Formal Verification of ACT-R Architectures and Models*. In: *CogSci*, pp. 618–624.
- [14] Richard S Lazarus (1982): *Thoughts on the relations between emotion and cognition*. *American psychologist* 37(9), p. 1019, doi:10.1037/0003-066X.37.9.1019.
- [15] Richard S Lazarus (1991): *Emotion and adaptation*. Oxford University Press. Available at <https://books.google.de/books?id=tTdIlwpxtWsC>.
- [16] Jerry Lin, Marc Spraragen, Jim Blythe & Michael Zyda (2011): *EmoCog: Computational Integration of Emotion and Cognitive Architecture*. In: *The Florida AI Research Society*. Available at <http://aaai.org/ocs/index.php/FLAIRS/FLAIRS11/paper/view/2625>.
- [17] Stacy Marsella & Jonathan Gratch (2009): *EMA: A process model of appraisal dynamics*. *Cognitive Systems Research* 10, pp. 70–90, doi:10.1016/j.cogsys.2008.03.005.
- [18] George A Miller (1956): *The magical number seven, plus or minus two: Some limits on our capacity for processing information*. *Psychological review* 63(2), p. 81, doi:10.1037/h0043158.
- [19] J.A. Moura & D.H. Hutchison (2019): *Game Theory for Multi-Access Edge Computing: Survey, Use Cases, and Future Trends*. *IEEE Communications Surveys and Tutorials* 21(1), pp. 260–288, doi:10.1109/COMST.2018.2863030.
- [20] Roger B. Myerson (1991): *Game Theory: Analysis of Conflict*. Harvard University Press. Available at <http://www.jstor.org/stable/j.ctvjsf522>.
- [21] John von Neumann & Oskar Morgenstern (2007): *Theory of Games and Economic Behavior (60th Anniversary Edition)*. Princeton University Press, doi:10.1515/9781400829460. Available at <http://www.jstor.org/stable/j.ctt1r2gkx>.
- [22] Allen Newell (1990): *Unified Theories of Cognition*. Harvard University Press, USA.
- [23] Eric Nivel (2007): *Ikon flux 2.0*. Available at <https://api.semanticscholar.org/CorpusID:60599731>.
- [24] Keith Oatley & P. N. Johnson-laird (1987): *Towards a Cognitive Theory of Emotions*. *Cognition and Emotion* 1(1), pp. 29–50, doi:10.1080/02699938708408362.
- [25] David Peebles & Adrian Banks (2010): *Modelling Dynamic Decision Making with the ACT-R Cognitive Architecture*. *Journal of Artificial General Intelligence* 2, pp. 52–68, doi:10.2478/v10229-011-0009-1.
- [26] Astrid Rakow (2023): *Framing Relevance for Safety-Critical Autonomous Systems*, doi:10.48550/arXiv.2307.14355. arXiv:2307.14355. Technical Report.
- [27] Rainer Reisenzein, Eva Hudlicka, Mehdi Dastani, Jonathan Gratch, Koen Hindriks, Emiliano Lorini & John-Jules Meyer (2013): *Computational Modeling of Emotion: Toward Improving the Inter- and Intradisciplinary Exchange*. *IEEE Trans. Affect. Comput.* 4(3), p. 246–266, doi:10.1109/T-AFFC.2013.14.
- [28] Niels Taatgen (2013): *The Nature and Transfer of Cognitive Skills*. *Psychological review* 120, pp. 439–471, doi:10.1037/a0033138.
- [29] Frits Vaandrager, Bharat Garhewal, Jurriaan Rot & Thorsten Wißmann (2022): *A New Approach for Active Automata Learning Based on Apartness*. In: *Tools and Algorithms for the Construction and Analysis of Systems*, Springer International Publishing, pp. 223–243, doi:10.48550/arXiv.2107.05419.
- [30] Christopher Wickens, John Helleberg, Juliana Goh, Xidong Xu & William Horrey (2001): *Pilot Task Management: Testing an Attentional Expected Value Model of Visual Scanning*. Savoy, IL, UIUC Institute of Aviation Technical Report.
- [31] Bertram Wortelen (2014): *Das Adaptive-Information-Expectancy-Modell zur Aufmerksamkeitssimulation eines kognitiven Fahrermodells*. Ph.D. thesis, Carl von Ossietzky Universität, Oldenburg, Germany. Available at <https://oops.uni-oldenburg.de/id/eprint/1970>.
- [32] Bertram Wortelen, Andreas Lüdtkke & Martin Baumann (2013): *Integrated Simulation of Attention Distribution and Driving Behavior*.

- [33] Greg L Zacharias, Adam X Miao, Christine Illgen, Jake M Yara & GM Siouris (1996): *SAMPLE: Situation awareness model for pilot in-the-loop evaluation*. In: *Proceedings of the 1st Annual Conference on Situation Awareness in the Tactical Air Environment*, Citeseer.

# Trust Modelling and Verification Using Event-B

Asieh Salehi Fathabadi

University of Southampton, United Kingdom

a-salehi-fathabadi@soton.ac.uk

Vahid Yazdanpanah

University of Southampton, United Kingdom

v.yazdanpanah@soton.ac.uk

Trust is a crucial component in collaborative multiagent systems (MAS) involving humans and autonomous AI agents. Rather than assuming trust based on past system behaviours, it is important to formally verify trust by modelling the current state and capabilities of agents. We argue for verifying actual trust relations based on agents' abilities to deliver intended outcomes in specific contexts. To enable reasoning about different notions of trust, we propose using the refinement-based formal method Event-B. Refinement allows progressively introducing new aspects of trust - from abstract to concrete models incorporating knowledge and runtime states. We demonstrate modelling three trust concepts and verifying associated trust properties in MAS. The formal, correctness-by-construction approach allows to deduce guarantees about trustworthy autonomy in human-AI partnerships. Overall, our contribution facilitates rigorous verification of trust in multiagent systems.

## 1 Introduction

Trust is a crucial contextual concept in multiagent systems (MAS), representing the cognitive state of a trustor towards a trustee [3]. While there are various accounts of trust, in this work, we focus on trust with respect to accomplishing tasks. While trust modelling in MAS has historically relied on reasoning about past behaviours[9], recent work emphasises integrating current context rather than fully depending on history. This involves verifying what agents can actually deliver based on their present capabilities, beyond reputations. We argue for complementing offline safety assurances with online trust verification for autonomous systems. Consider an autonomous delivery vehicle (ADV) tasked with transporting goods. Offline verification during design suffices for basic safety and whether the ADV is reliable in general (regardless of their current state and how they can perform in the context). However, assessing trust online for a particular delivery also requires checking the ADV's abilities given its current battery, payload etc. against user requirements.

Trust modelling in MAS, and what we introduce as "actual trust", entails representing different aspects like agents' abilities, knowledge and commitments. To model such a multidimensional notion, refinement techniques like Event-B [1] allow correct-by-construction modelling [7, 5]. Our key contribution is a refinement-based approach that supports formally verifying various trust concepts. We demonstrate formally modelling three trust notions relating to agent abilities, knowledge and commitments. The automated consistency guarantees complement offline assurance for trustworthy autonomy and human-AI partnerships [10, 11]. This work is an initial step on modelling trust using Event-B's refinement strategy that practically enables step-wise verification of actual trust between agents in autonomous systems.

## 2 Actual Trust: Power, Knowledge, and Commitments

In modelling and reasoning about trust, it is key to distinguish what an agent may rely on due to past behaviour of another agent and their *typical* behaviour from what in a given situation agents are *actually*

able to deliver. While the former category of trusting has a retrospective view, and uses history to reason about trust [9], the latter form of trust is to reason about what the other agent can *actually* deliver and has basis in what is known in the theory of causality as *actual causality* [6]. In this work, we focus on the latter notion, refer to it as *actual trust* and understand it as a relational notion between two agents or agent groups  $i$  (as the trustor) and  $j$  (as the trustee) and say in a particular multiagent system  $M$ ,  $i$  trusts  $j$  with respect to task  $t$  only if  $i$  is able to verify that  $j$  is able and committed to deliver  $t$ . To model and verify our notion of actual trust, as knowledge of another agents' ability and commitment to ensure a particular task, it is important to highlight how it relates to its key components conceptually.

**Trusting for ability to materialise eventualities:** In contrast to purely history-oriented perspectives to trust, that look at the history and trust an agent to behave similar to its past behaviour, we deem that trusting needs to be fine-tuned based on the current state of the system and actions agents are able to execute and what agents intend to deliver. For instance, even if an ADV was successful in former deliveries, it may be suffering from a low battery now and unable to deliver tasks. So, one should fine tune trust in the agent's power to deliver based on the current situation.

**Trust as an epistemic state:** We understand actual trust as an epistemic notion meaning that it is essentially about knowledge of the trustor on how another agent relates to a particular event. Recalling the running example, the user needs to reason about abilities of an ADV, consider its publicly announced intentions, and verify if the ADV can be trusted for a particular delivery. This form of trust allows specification of trust in different contexts and for different types requirements and knowledge levels. For instance, a given ADV  $j$  may be seen as "trusted for delivering 5kg of groceries" but this trust may not extend when it comes to passenger pickup..

**Public commitments as a proxy to intentions:** When we are dealing with autonomous AI agents, we need to consider that being able to deliver a task fundamentally differs from delivering the task. Imagine that an ADV  $v$  with a full battery and ability to deliver some goods is located relatively close to an agent  $i$  with a delivery task  $t$ . In this case,  $i$  can't simply assume that  $v$  can be trusted for delivering  $t$  as it may be already committed to deliver tasks other than  $t$  or is in the middle of other plan executions. To handle this, we use notion of publicly-announced commitments as a proxy to model what agents intend to bring about.<sup>1</sup>

### 3 Refinement-Based Trust Formal Modelling and Verification

**Background knowledge:** Event-B [1] is a refinement-based formal method for system development. The mathematical language of Event-B is based on set theory and first order logic<sup>a</sup>. An Event-B model consists of two parts: *contexts* for static data and *machines* for dynamic behaviour. An Event-B model is constructed by making progressive refinements starting from an initial abstract model which may have more general behaviours and gradually introducing more detail that constrains the behaviour towards the desired system. Each refinement step is verified to be a valid refinement of the previous step.

<sup>a</sup>Please refer to Event-B Language user manual [https://wiki.event-b.org/index.php/Event-B\\_Language](https://wiki.event-b.org/index.php/Event-B_Language) for extra support to understand the presented model.

Benefiting from refinement technique in building Event-B formal model, instead of one single-layer complex design model of system, we propose to gradually introduce different concepts of actual trust

<sup>1</sup>Note that assuming full access to agents' intentions is against separation of concerns, privacy, and encapsulation as key design principles in safe and responsible AI and software development.

through refinement steps. Figure 1 presents our vision idea of applying refinement-based development to model actual trust in autonomous systems<sup>2</sup>. Left side illustrates the trust relationship between trustor and trustee, while right side presents the structure of our proposed Event-B formal model, including three levels of refinements: machines ( $M0$ ,  $M1$  and  $M2$ ) and associated contexts ( $cntx0$ ,  $cntx1$  and  $cntx2$ ).

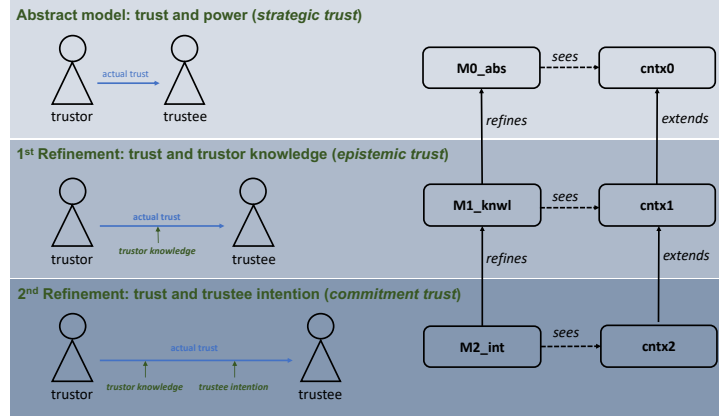


Figure 1: Trust Modelling and the Refinement Strategy

In line with trust key components and first-order building blocks presented in Section 2, starting from the top level, trust is first modelled as an abstract relationship between two agents, trustor and trustee,  $M0\_abs$ ; followed by first refinement level where *trustor knowledge* is introduced,  $M1\_knwl$ . Then *trustee intention* is introduced in a further refinement level,  $M2\_int$ .

### 3.1 Modelling trust in Event-B

Agents and tasks are defined as a set in the context  $cntx0$ , which is partitioned to two sub-sets: trustors and trustees:

**Background knowledge:** An Event-B context contains carrier sets  $s$ , constants  $c$ , and axioms  $A(c)$  that constrain the carrier sets and constants.

```
CONTEXT cntx0
SETS AGENTS TASKS CONSTANTS trustors trustees
AXIOMS axm1 : trustors ⊆ AGENTS // Definition: Subset ⊆
        axm2 : trustees ⊆ AGENTS
        axm3 : partition (AGENTS, trustors, trustees)
```

**Definition 1.** Abstract (strategic) trust: agent  $i$  weakly (abstraction) trusts  $j$  regarding task  $t$  if  $j$  has an action or a sequence of actions available to it to ensure  $e$ . We are operating in a cooperative setting, hence assuming that agents share information and have perfect knowledge of themselves as well as other agents' abilities. We define trust in terms of the ability to deliver  $t$ .

Trust is modelled as a three-dimension relation variable between a trustor, a set of trustees and a task, in the abstract machine. In  $M0\_abs$ , an invariant,  $inv1$ , is specifying the *agent\_task* variable as a function between a subset of trustees and a single task, indicating the task that can be delivered by a subset of

<sup>2</sup>Note that because of space limitation, the Event-B model of trust is not fully presented here. And for simplicity, in purpose of demonstrating the vision idea, we model trust in its simplest definition.

trustee agents. And *inv2* is specifying the *trustor\_trustee\_task* variable indicating the relation between a trustor and a pair of *agent\_task*.

**Background knowledge:** An Event-B machine contains variables  $v$ , invariant predicates  $I(v)$  that constrain the variables, and events. In Event-B, a machine corresponds to a transition system where *variables* represent the states and *events* specify the transitions.

```
@inv1: agent_task ∈ P(trustees) ↔ tasks // specifies which agents are able to deliver which task
// Definition: Powerset: P(S) = {s | s ⊆ S}
// Definition: A function (agent_task) is a relation with the restriction that each element of the
// domain (P(trustees)) is related to a unique element in the range (tasks); a many to one mapping
// Definition: Set membership ∈
@inv2: trustor_trustee_task ∈ trustors ↔ agent_task // specifies set of triples i ↦ (j ↦ t), when
// agent i can trust a set of agents j to deliver task t
```

The event *trust* is adding a new triple to the *trustor\_trustee\_task* variable, *act1*. While *grd1* – 3 is checking the type of the event parameters, *grd4* is indicating the above definition, ensuring *j* is able to deliver task *t*; guards *grd5* – 8 are described later.

**Background knowledge:** An event in a machine, comprises a guard denoting its enabling-condition and an action describing how the variables are modified when the event is executed. In general, an event *e* has the following form, where *t* are the event parameters,  $G(t, v)$  is the guard of the event, and  $v := E(t, v)$  is the action of the event: **e == any t where  $G(t, v)$  then  $v := E(t, v)$  end**

```
event trust any i j t
where @grd1 : i ∈ trustors
      @grd2 : j ∈ P(trustees)
      @grd3 : t ∈ tasks
      @grd4 : t ∈ agent_task[{j}] // j is able to deliver task t
      // Definition: relational image: r[S] = {y | ∃x. x ∈ S ∧ x ↦ y ∈ r} where S is a set
      @grd5 : i ∉ j // to preserve inv3
      // Definition: Set non-membership ∉
      @grd6 : j ≠ ∅ // abstract guard to preserve inv4
      @grd7 : j ⊆ knowledge[{i}] // refining guard to preserve inv4
      @grd8 : commitments[i ↦ (j ↦ t)] = {TRUE} // refining guard to preserve inv4
then @act1: trustor_trustee_task := trustor_trustee_task ∪ {i ↦ (j ↦ t)}
// Definition: Union ∪
```

**Running example:** For instance, for an agent *i* and an ADV *j* and task of “delivering 5kg of groceries”, *i* can trust *j* only if “delivering 5kg of groceries” is within the allocated tasks to *j*: *grd3*. Then, *act1* will add a new triple of  $(i, j, t)$  to the variable set *trustor\_trustee\_task*.

### 3.2 Modelling verifiable trust properties

To propose the idea of formal verification of properties of trust in autonomous systems, here we present two invariants, specifying two fundamental trust properties. *inv3* is specifying that an agent *i* would not trust itself to deliver a specific task *t*. And *inv4* is specifying avoiding trust deadlock, that for each trustor *i* and task *t*, there is always a non-empty subset of trustees *j* that can deliver *t*.

```

@inv3:  $\forall i, j. i \in \text{trustors} \wedge j \in \mathbb{P}(\text{trustees}) \wedge i \in \text{dom}(\text{trustor\_trustee\_task}) \Rightarrow i \notin j$ 
// Definition: Conjunction  $\wedge$ , Universal quantification  $\forall$ , Implication  $\Rightarrow$ 
// Definition: Domain:  $\text{dom}(r) \forall r. r \in S \leftrightarrow T \Rightarrow \text{dom}(r) = \{x. (\exists y. x \mapsto y \in r)\}$  where S and T are sets
@inv4:  $\forall i, t. i \in \text{trustors} \wedge t \in \text{tasks} \Rightarrow (\exists j. j \in \mathbb{P}(\text{trustees}) \wedge j \neq \emptyset)$ 
// Definition: Existential quantification  $\exists$ 

```

### 3.3 Verifying trust properties

**Background knowledge:** Event-B is supported by the Rodin tool set [2], an extensible open source toolkit which includes facilities for modelling, verifying the consistency of models using theorem proving and model checking techniques, and validating models with simulation-based approaches.

One of the generated proof obligations (PO) for an Event-B model, is "invariant preservation":

$e/v/INV$  (where  $e$  is the event name, and  $v$  is the invariant name)

$INV$  PO ensures that the property specified in the invariant  $INV$  is preserved by event  $e$ . To preserve the trust properties defined in  $inv3$  and  $inv4$ , the event  $trust$  is guarded by  $grd5$  and  $grd6$ , see above. Two POs  $trust/inv3/INV$  and  $trust/inv4/INV$  are generated and automatically discharged by Rodin tool.

### 3.4 Refining trust

Next, we introduce the refined notion of epistemic trust in which agents' knowledge is integrated.

**Definition 2.** Refined (epistemic) trust: for a stronger notion of trust we require a variable of *knowledge* specifying the knowledge relationship between two agents  $i, j$ , indicating whether  $i$  is fully aware of  $j$ 's abilities.

Refining model  $M1\_knl$  introduces the *knowledge* variable to model the knowledge of trustors about trustees:

```
@inv1:  $\text{knowledge} \in \text{trustors} \leftrightarrow \text{trustees}$ 
```

```
// Definition: A relation (knowledge) is a set of ordered pairs; a many to many mapping.
```

**Running example:** For instance, a given ADV  $j$  may be seen as "trusted for delivering 5kg of groceries" by an agent  $i$  who is fully aware of  $j$ 's abilities but not by agent  $v$  who is not aware of  $j$  and that  $j$  has the capacity to ensure  $t$ .

**Definition 3.** Refined (commitment) trust: for a stronger notion of trust we require a variable of *commitments* specifying a function that takes a trust triple  $(i, j, t)$  and determines whether agent  $j$  is committed to deliver task  $t$  for agent  $i$ . We refine the trust model, not only in terms of the ability, but also the commitment to deliver  $t$ .

And refining model  $M2\_int$  introduces the *commitment* variable to model the intention of trustees to deliver the associate task (for simplicity in this paper, we model commitment as a Boolean indicating whether an agent(s) as trustee intends to deliver the associated task or not):

```
@inv1:  $\text{commitments} \in \text{trustor\_trustee\_task} \rightarrow \text{BOOL}$ 
```

$inv4$  is refined to include the knowledge property in  $M1\_knl$  and commitment specification in  $M2\_int$ :

```

@inv4:  $\forall i, t. i \in \text{trustors} \wedge t \in \text{tasks} \Rightarrow$ 
   $(\exists j. j \in \mathbb{P}(\text{trustees}) \wedge j \neq \emptyset \wedge (j \mapsto t) \in \text{agent\_task} \wedge$ 
   $j \subseteq \text{knowledge}[\{i\}] \wedge \text{commitments}[i \mapsto (j \mapsto t)] = \{\text{TRUE}\} \wedge$ 
   $i \mapsto (j \mapsto t) \in \text{trustor\_trustee\_task})$ 

```



And refining event *trust* includes extra guards *grd7* and *grd8* to preserve *inv4*, see above. Not providing these guards results in failed generated *INV* POs.

**Running example:** For instance, for an agent *i* and an ADV *j* and task of “delivering 5kg of groceries”, *i* can trust *j* only if “delivering 5kg of groceries” is within the allocated tasks to *j*: *grd3* (verified in the abstract machine), and *i* is fully aware of *j*’s abilities: *grd7* (verified in the machine *M1\_knwl*) and *j* is committed to deliver 5kg of groceries to *i*: *grd8* (verified in the machine *M2\_int*).

**Model checking trust properties:** The presented Event-B model can be model checked by instantiating the context elements, for example for the ADV system. Also the scenario checker integrated in Rodin can demonstrate difference scenarios of the desire system. Due to the concise nature of this work and space limitation, we are unable to include the model checking experience.

## 4 Concluding Remarks and Future Directions

The step-wise refinement approach presented in this paper, demonstrates three notions of actual trust, and two verifiable trust properties. The model can simply refined to include more notions and properties. This paper elaborates on how the autonomous system research can benefit from refinement-based formal methods in terms of modelling trust. The abstraction technique aids the modelling and verification process in step-wise manner, where instead of a single complex model, the formal model is gradually built through refinement levels, hence easier to be understood and proved. Also the Event-B formal method provides the verification techniques (theorem proving and ProB model checking [8]) in each refinement level, to ensure the trustworthiness of autonomous systems.

**Contributions to Autonomous Systems (AS):** In AS, replacing human decision-making with machine decision-making results in challenges associated with stakeholders’ trust. Trustworthiness of an AS is key to its wide-spread adoption by society. To develop a trusted AS, it is important to understand how different stakeholders perceive an AS as trusted, and how the context of application affects their perceptions. The translation of trust issues into formalised solutions is challenging due to trust dynamics. In this work, we try to advance in this direction by utilising the ability of Event-B as a refinement-based formal method to manage the lack of information when modelling trust in multi-agent systems. High-level model aids to abstract away the uncertain/unknown trust specifications. We introduced the notion of actual trust versus statistical trust, toward trusting to an AS due to its safety checks, like, inherent uncertainties in the environment, diversity in the requirements and needs of different users and contexts of application. We formalised the notion of actual trust using Event-B formal modelling followed by verifying the safety properties of it. The actual trust notions is modelled and verified in three levels: strategic trust, epistemic trust and commitment trust.

**Future directions:** This formal modelling and verification approach for trust in autonomous systems can be extended in several directions. One avenue is via integrating Event-B models with Alternating-Time Temporal Logic [13] to allow more expressive temporal specifications and model checking of trust properties, e.g., in the context of connected mobility systems. Further research can also investigate gradation of trust (as a quantifiable notion) and formally relating trust and neighbouring notions in multiagent settings such as responsibility [12]. Quantifying trust based on strategy lengths and information-theoretic notions may also complement the approach pursued here. Overall, rigorous formal methods can provide significant assurances about trustworthy autonomy and human-AI partnerships, especially for safety-critical applications.

*Acknowledgements:* This work was supported by the UK Engineering and Physical Sciences Research Council (EPSRC) through a Turing AI Fellowship (EP/V022067/1) on Citizen-Centric AI Systems (<https://ccaais.ac.uk/>) and the UKRI Trustworthy Autonomous Systems Hub (EP/V00784X/1).

## References

- [1] J-R. Abrial (2010): *Modeling in Event-B: System and Software Engineering*. Cambridge University Press, doi:10.1017/S0956796812000081.
- [2] J-R Abrial, M. Butler, S. Hallerstede, T.S. Hoang, F. Mehta & L. Voisin (2010): *Rodin: An Open Toolset for Modelling and Reasoning in Event-B*. *Software Tools for Technology Transfer* 12(6), pp. 447–466, doi:10.1007/s10009-010-0145-y.
- [3] Cristiano Castelfranchi & Rino Falcone (2020): *Trust: Perspectives in cognitive science*. *The Routledge Handbook of Trust and Philosophy*, pp. 214–228, doi:10.4324/9781315542294-17.
- [4] Mehdi Dastani & Vahid Yazdanpanah (2023): *Responsibility of AI systems*. *Ai & Society* 38(2), pp. 843–852, doi:10.1007/s00146-022-01481-4.
- [5] Hang-Jiang Gao, Zheng Qin, Lei Lu, Li-Ping Shao & Xing-Chen Heng (2007): *Formal specification and proof of multi-agent applications using event b*. *Information Technology Journal* 6(7), pp. 1181–1189, doi:10.3923/itj.2007.1181.1189.
- [6] Joseph Y Halpern (2016): *Actual causality*. MIT Press, doi:10.7551/mitpress/10809.001.0001.
- [7] Arnaud Lanoix (2008): *Event-B Specification of a Situated Multi-Agent System: Study of a Platoon of Vehicles*. In: *Second IEEE/IFIP International Symposium on Theoretical Aspects of Software Engineering, TASE 2008, June 17-19, 2008, Nanjing, China*, IEEE Computer Society, pp. 297–304, doi:10.1109/TASE.2008.39.
- [8] Michael Leuschel & Michael Butler (2008): *ProB: An Automated Analysis Toolset for the B Method*. *Software Tools for Technology Transfer (STTT)* 10(2), pp. 185–203, doi:10.1007/s10009-007-0063-9.
- [9] Sarvapali D Ramchurn, Dong Huynh & Nicholas R Jennings (2004): *Trust in multi-agent systems*. *The knowledge engineering review* 19(1), pp. 1–25, doi:10.1017/S0269888904000116.
- [10] Sarvapali D Ramchurn, Sebastian Stein & Nicholas R Jennings (2021): *Trustworthy human-AI partnerships*. *Iscience* 24(8), p. 102891, doi:10.1016/j.isci.2021.102891.
- [11] Sebastian Stein & Vahid Yazdanpanah (2023): *Citizen-Centric Multiagent Systems*. In: *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pp. 1802–1807. Available at <https://dl.acm.org/doi/10.5555/3545946.3598843>.
- [12] Vahid Yazdanpanah & Mehdi Dastani (2016): *Quantified degrees of group responsibility*. In: *Coordination, Organizations, Institutions, and Norms in Agent Systems*, Springer, pp. 418–436, doi:10.1007/978-3-319-42691-4\_23.
- [13] Chenyang Zhu, Michael Butler, Corina Cirstea & Thai Son Hoang (2023): *A fairness-based refinement strategy to transform liveness properties in Event-B models*. *Science of Computer Programming* 225, p. 102907, doi:10.1016/j.scico.2022.102907.

# Comparing Differentiable Logics for Learning Systems: A Research Preview\*

Thomas Flinkow

Department of Computer Science  
Maynooth University  
Maynooth, Ireland

thomas.flinkow@mu.ie

Barak A. Pearlmutter

Department of Computer Science and Hamilton Institute  
Maynooth University  
Maynooth, Ireland

barak@pearlmutter.net

Rosemary Monahan

rosemary.monahan@mu.ie

Extensive research on formal verification of machine learning (ML) systems indicates that learning from data alone often fails to capture underlying background knowledge. A variety of verifiers have been developed to ensure that a machine-learned model satisfies correctness and safety properties, however, these verifiers typically assume a trained network with fixed weights. ML-enabled autonomous systems are required to not only detect incorrect predictions, but should also possess the ability to self-correct, continuously improving and adapting. A promising approach for creating ML models that inherently satisfy constraints is to encode background knowledge as logical constraints that guide the learning process via so-called differentiable logics. In this research preview, we compare and evaluate various logics from the literature in weakly-supervised contexts, presenting our findings and highlighting open problems for future work. Our experimental results are broadly consistent with results reported previously in literature; however, learning with differentiable logics introduces a new hyperparameter that is difficult to tune and has significant influence on the effectiveness of the logics.

## 1 Introduction

Advancements in machine learning (ML) in the past few years indicate great potential for applying ML to various domains. Autonomous systems are one such application domain, but using ML components in such a safety-critical domain presents unique new challenges for formal verification. These include (1) ML failing to learn background knowledge from data alone [16], (2) neural networks being susceptible to adversarial inputs, and (3) a lack of specifications, generally and especially when continuous learning is permitted [3]. Addressing these challenges is even more important and more difficult when the ML-enabled autonomous system is permitted to continue to learn after deployment, either to adapt to changing environments or to correct and improve itself when errors are detected [2].

A multitude of neural network verifiers (c.f. [10]) have been presented in the past few years; prominent examples include Reluplex [6], NNV [12], and MN-BaB [4]. These solvers use techniques such as satisfiability and reachability analysis, and can verify a variety of properties. However, these verifiers typically assume trained networks with fixed weights and do not target the learning process itself [7]. One step in the direction of correct-by-construction neural networks are so-called *differentiable logics*, which transform a logical constraint  $\phi$  into an additional logical loss term  $\mathcal{L}_L(\phi)$  to minimise when learning, where the logical loss of a constraint  $\phi$  is combined with standard cross-entropy loss as  $\mathcal{L} = \mathcal{L}_{CE} + \lambda \mathcal{L}_L(\phi)$ . In order to translate logical constraints into loss terms, a mapping must be defined that allows for real-valued truth values, and that is differentiable almost everywhere for use with standard gradient-based methods. In the following, we give a brief overview of two of these mappings (so-called *differentiable logics*) from popular literature, namely DL2 and fuzzy logics.

---

\*This publication has emanated from research conducted with the financial support of Science Foundation Ireland under Grant number 20/FFP-P/8853.

Table 1: The t-norms, t-conorms and implications used in our experiments.

Name	T-norm (Conjunction)	T-conorm (Disjunction)	Implication
Gödel	$T_G(x, y) = \min(x, y)$	$S_G(x, y) = \max(x, y)$	$I_G(x, y) = \begin{cases} 1, & \text{if } x < y \\ y \end{cases}$ $I_{KD}(x, y) = \max(\bar{x}, y)$
Łukasiewicz	$T_{LK}(x, y) = \max(0, x + y - 1)$	$S_{LK}(x, y) = \min(1, x + y)$	$I_{LK}(x, y) = \min(\bar{x} + y, 1)$
Yager	$T_{YG}(x, y) = \max(0, \overline{\bar{x} \Delta_p \bar{y}})$	$S_{YG}(x, y) = \min(1, x \Delta_p y)$	$I_{YG}(x, y) = \begin{cases} 1, & \text{if } x = y = 0 \\ y^x \end{cases}$
Product	$T_P(x, y) = xy$	$S_{PS}(x, y) = x + y - xy$	$I_{GG}(x, y) = \begin{cases} 1, & \text{if } x < y \\ y/x \end{cases}$ $I_{RC}(x, y) = \bar{x} + xy$

where  $u \Delta_p v := \sqrt[p]{|u|^p + |v|^p}$  is the  $p$ -norm Pythagorean sum,  $p \geq 1$ , and  $\bar{u} := 1 - u$

## 1.1 DL2

DL2 (“Deep Learning with Differentiable Logics”) [5] is a system for querying and training neural networks with logic. It maps absolute truth to 0 and other degrees of truth to positive values up to  $\infty$  based on the following elementary translation rules:

$$\mathcal{L}_{DL2}(x \leq y) := \max(x - y, 0), \quad \mathcal{L}_{DL2}(x \wedge y) := \mathcal{L}_{DL2}(x) + \mathcal{L}_{DL2}(y), \quad \mathcal{L}_{DL2}(x \vee y) := \mathcal{L}_{DL2}(x) \cdot \mathcal{L}_{DL2}(y). \quad (1)$$

Additionally, there is  $\mathcal{L}_{DL2}(x \neq y) := \xi[x = y]$ , where  $\xi > 0$  denotes a constant that was found to not have significant influence [5], and  $[\cdot]$  being the indicator (in Knuth’s notation). From these, other rules such as  $\mathcal{L}_{DL2}(x < y) := \mathcal{L}_{DL2}(x \leq y \wedge x \neq y)$  can be derived. Negation is handled by pushing the negation inwards to the level of comparison, e.g.  $\mathcal{L}_{DL2}(\neg(x \leq y)) := \mathcal{L}_{DL2}(y < x)$ . DL2 does not have a separate translation for implication, instead translating implication as  $\mathcal{L}_{DL2}(x \rightarrow y) := \mathcal{L}_{DL2}(\neg x \vee y)$ .

## 1.2 Fuzzy Logics

Whereas DL2 was designed specifically for deep learning contexts, fuzzy logics are logical systems that have been studied extensively and happen to be suitable for use as differentiable logics due to their many-valued nature, with operators that are often differentiable almost everywhere. Fuzzy logics express degrees of truth in the unit interval  $[0, 1]$ , with absolute truth mapped to 1. We use  $\mathcal{L}_L(\phi) := 1.0 - \mathcal{L}_{FL}(\phi)$  for the fuzzy logic loss in our implementation to address the inverse notion of truth.

Fuzzy logics are based on functions  $T : [0, 1]^2 \rightarrow [0, 1]$  that are commutative, associative, monotonic, and satisfy  $T(1, y) = y$ . These are called triangular norms (abbreviated as *t-norms*) and generalise conjunction. A t-conorm (also called *s-norm*) generalises disjunction and can be obtained from a t-norm using  $S(x, y) = 1 - T(1 - x, 1 - y)$ . From a t-conorm  $S$  and fuzzy negation  $N$ , one obtains a so-called  $(S, N)$ -implication (which generalises material implication) as  $I(x, y) := S(N(x), y)$ . Examples of  $(S, N)$ -implications are the Kleene-Dienes implication  $I_{KD}$  and Reichenbach implication  $I_{RC}$ , both with the standard negation  $N(x) = 1 - x$ . Other implications generalise the intuitionistic implication and are called *R-implications*, because they use the t-norm residuum  $R(x, y) = \sup\{t \in [0, 1] \mid T(x, t) \leq y\}$ . Example *R-implications* are the Gödel implication  $I_G$  and Goguen implication  $I_{GG}$ . The Łukasiewicz implication  $I_{LK}$  is both an  $(S, N)$ -implication and an *R-implication*. Other implications are neither—the Yager implication

$I_{YG}$ , for example, is an  $f$ -generated implication that is obtained using  $f(x) = -\ln x$  in  $I(x, y) := f^{-1}(xf(y))$  (with the understanding that  $0 \cdot \infty = 0$ ).

Additionally, [14] propose *sigmoidal* implications in order to prove the derivatives of the original implication, while preserving its characteristics. In Eq. (2),  $\sigma(x) := 1/(1 + \exp(-x))$  denotes the standard sigmoidal function and  $s$  is a parameter controlling the steepness. We use the sigmoidal implication in our experiments with  $I_{RC}$  and  $s = 9$ , as suggested by [14].

$$(I(x, y))_s := \frac{(1 + \exp(s/2))\sigma(sI(x, y) - s/2) - 1}{\exp(s/2) - 1} \quad (2)$$

Lastly, given a fuzzy implication  $I$  and bijection  $\phi : [0, 1]^2 \rightarrow [0, 1]$ , [1] show that the function  $(I(x, y))_\phi := \phi^{-1}I(\phi(x), \phi(y))$  is also a fuzzy implication. We use this in our experiments with the Reichenbach implication  $I_{RC}$  and  $\phi(x) = x^2$ .

Table 1 lists the definitions of the mentioned t-norms, t-conorms and implications, and Fig. 3 (Appendix A) displays plots of the implications.

**Mapping atomic terms** DL2 is designed for atomic terms that are inequalities or comparisons. As seen in Eq. (1), it provides the loss translation  $\mathcal{L}_{DL2}(x \leq y) := \max(x - y, 0)$  for comparison.

Fuzzy logics typically do not define fuzzy comparison operators. However, [11] introduce a mapping  $\mathcal{L}_{FL}(x \leq y) := 1 - \max\left(\frac{x-y}{x+y}, 0\right)$  for fuzzy logics. Fuzzy logics requires the truth values of the atomic terms  $x, y$  to be mapped into  $[0, 1]$  by some oracle. Because the atomic terms in our constraints are comparisons, we change this mapping from  $\mathcal{L}_{FL}(x \leq y) : [0, 1]^2 \rightarrow [0, 1]$  to  $\mathcal{L}_{FL}(x \leq y) : \mathbb{R}^2 \rightarrow [0, 1]$ , allowing us to forgo the need for an external oracle. The mapping is shown in Eq. (3) below, where we use  $\varepsilon = 0.05$ .

$$\mathcal{L}_{FL}(x \leq y) := 1 - \frac{\max(x - y, 0)}{|x| + |y| + \varepsilon} \quad (3)$$

Note that the fuzzy logic mapping  $\mathcal{L}_{FL}(x \leq y)$  has a property we intuitively might wish to hold: for example, we might want  $21 \leq 20$  to be as much of a violation as  $21000 \leq 20000$ . This cannot be achieved in DL2, where the violation depends only on the absolute difference.

## 2 Comparing Differentiable Logics: Experimental Setup

Our comparison experiment<sup>1</sup> is implemented in PyTorch and based on the original experiment in [5]. We train on the Fashion-MNIST, CIFAR-10, and GTSRB data sets with various constraints. In order to create meaningful scenarios where learning with logical constraints surpasses the baseline (learning from data alone), we train with a fraction of the data sets, namely 10 % for Fashion-MNIST, 50 % for CIFAR-10, and 90 % for GTSRB (as it consists of more classes, with more imperfect data). Additionally, we introduce 10 % label noise (training with incorrect labels) for all data sets, and apply various image manipulation techniques, such as random cropping, flipping, and colour changes. A batch size of 256 was used for all datasets.

The goal of our experiment is to compare various differentiable logics, including DL2 and popular fuzzy logics, and investigate which logic performs most favourable, focusing specifically on implication and conjunction, as these have noticeable consequences for the learning process: As pointed out by [14],

<sup>1</sup> Available on <https://github.com/tflinkow/dl-comparison>.

background knowledge and constraints are most often of the form “if  $A$ , then  $B$ ”. Choosing a suitable implication that performs well is thus an important task to guarantee best learning.

In [15], the authors introduce the *shadow-lifting* property for a conjunction, which requires the truth value of a conjunction to increase when the truth value of a conjunct increases. This property seems highly desirable for learning, as it allows for gradual improvement. For example, the formula  $0.1 \wedge 1.0$  should be more true than  $0.1 \wedge 0.2$ , but the Gödel t-norm  $T_G(x, y) = \min(x, y)$  would yield the same truth value in both cases. DL2 uses addition for conjunction, trivially satisfying shadow-lifting. The only t-norm to satisfy the shadow-lifting property is the product t-norm  $T_P(x, y) = xy$ . However, as noted by [14], its derivative will be low if  $x$  and  $y$  are both low, making it hard for the learning process to make progress.

## 2.1 Constraints

**Universal quantification** In [5], the authors categorise constraints into two distinct schemes: training set constraints, which relate sampled inputs  $\mathbf{x}$  and  $\mathbf{x}'$  from the training set, and global constraints, which concern inputs in the  $\varepsilon$ -ball around a particular input. They use a PGD-based approach for universally quantified constraints and are thus limited to robustness properties, as noted by [11]. In [14], the authors relax infinite quantifiers by assuming minibatches to be subsets of an independent distribution and using finite conjunction for universal quantifiers for the minibatch, thus losing soundness. [11] provide a semantics for quantifiers, independent of the concrete differentiable logic and going beyond robustness via expectation minimisation of a probability distribution.

We do not consider global constraints in our experiment and only utilise finite universal quantification via repeated application of conjunction. This is permitted due to t-norms being associative and commutative by definition, and DL2’s use of addition for conjunction.

**Investigating implication** Limited to only training set constraints, the original DL2 experiment has shown some constraints to already be satisfied in the baseline experiments, where learning with logics would only provide minor improvements. Their robustness constraint, for example, was already 94.5% satisfied on Fashion-MNIST for the baseline, compared to 98.36% with DL2.<sup>2</sup>

They also use a class-similarity constraint for the CIFAR-10 network to encode domain knowledge such as “a car is more similar to a truck than to a dog”. In their fully-supervised experiment, even the baseline experiment satisfied this constraint quite well, and DL2 was able to only improve constraint accuracy from 93.67% to 99.68%. As all fuzzy implications  $I(x, y)$  by definition behave the same for  $x = 0$  or  $x = 1$ , the original constraint<sup>3</sup> is not suitable to compare different implications. Our modified constraint is shown in Eq. (4)<sup>4</sup> and replaces the binary decision by a soft one, checking whether the network output for label  $l$  is greater than or equal to  $1/\#\text{classes}$ .

$$\text{CSim}(\mathcal{N}, \mathbf{x}, \text{Labels}) := \bigwedge_{(l_1, l_2, l_3) \in \text{Labels}} ((\mathcal{N}(\mathbf{x})_{l_1} \geq 1/10) \longrightarrow (\mathcal{N}(\mathbf{x})_{l_2} \geq \mathcal{N}(\mathbf{x})_{l_3})). \quad (4)$$

As explained above, we translate the conjunction using repeated application of the product t-norm for all fuzzy logics in order to only investigate the different mappings for implication.

<sup>2</sup> The only training set constraint where DL2 could significantly improve constraint accuracy in the original experiments (from 5.62% to 99.78% in Fashion-MNIST) is the Lipschitz constraint  $\text{Lipschitz}(\mathcal{N}, \mathbf{x}, \mathbf{x}', L) := \|\mathcal{N}(\mathbf{x}) - \mathcal{N}(\mathbf{x}')\|_2 \leq L\|\mathbf{x} - \mathbf{x}'\|_2$ . We do not include this constraint in our experiments as it does not use conjunction nor implication.

<sup>3</sup>

$$\text{CSim}(\mathcal{N}, \mathbf{x}, \text{Labels}) := \bigwedge_{(l_1, l_2, l_3) \in \text{Labels}} ((\text{argmax} \mathcal{N}(\mathbf{x}) = l_1) \longrightarrow (\mathcal{N}(\mathbf{x})_{l_2} \geq \mathcal{N}(\mathbf{x})_{l_3})).$$

<sup>4</sup> The definition of Labels is shown in Eq. (6) for Fashion-MNIST, and in Eq. (7) for CIFAR-10, both in Appendix A.

**Investigating conjunction** For investigating the shadow-lifting effect, we utilise the German traffic sign recognition benchmark (GTSRB) dataset and use a property that forces the network to make confident, strong decisions by requiring all elements of a group of classes to be either very likely or very unlikely. Groups consist of classes of a similar type (e.g. speed limit signs)<sup>5</sup>.

$$\text{Group}(\mathcal{N}, \mathbf{x}, \varepsilon, \text{Groups}) := \bigwedge_{\{g_i\} \in \text{Groups}} \left( \sum_i \mathcal{N}(\mathbf{x})_{g_i} \leq \varepsilon \vee \sum_i \mathcal{N}(\mathbf{x})_{g_i} \geq 1 - \varepsilon \right). \quad (5)$$

We use the probabilistic sum t-conorm  $S_{\text{PS}}$  for fuzzy disjunctions in order to only focus on the conjunction.

### 3 Results

Table 2a shows the results obtained from running the class-similarity constraint experiment on the Fashion-MNIST and CIFAR-10 networks, and Table 2b shows the results obtained from running the group constraint on GTSRB. For each of these, the displayed prediction and constraint accuracy are obtained by taking the largest of their products from the last 10 epochs. Additionally, Figs. 1 and 2 (Appendix A) show how prediction and constraint accuracy change over time.

What immediately stands out is that when training with any logic, constraint accuracy is significantly improved, while prediction accuracy is always slightly reduced, compared to the baseline experiment. This could be because our constraints fail to capture useful background knowledge that would help with predictions, or, as [5] note, we might observe a similar phenomenon as reported in [13], where adversarial robustness conflicts with standard generalisation.

Our observed results are broadly consistent with trends previously reported in literature. The Gödel and Goguen implications perform badly as we expect many wrong inferences (due to the non-existing derivative of  $I_G$  with respect to  $x$ , and due to the Goguen implication’s singularity as  $x \rightarrow 0$ ), which manifest in the table as reduced prediction and constraint accuracy.

Comparing DL2 and fuzzy logics, our results indicate that for implication, fuzzy logics are the better choice — for CIFAR-10, even  $I_{\text{KD}}$  performs better than DL2, although both are rewriting  $x \rightarrow y$  to  $\neg x \vee y$ . This difference could be due to a multitude of reasons; the mappings  $\mathcal{L}(x \leq y)$  and  $\mathcal{L}(x \wedge y)$  are very different for DL2 and fuzzy logics, as is their range. Albeit the most likely reason is a sub-optimal choice of hyperparameter  $\lambda$ , explained in more detail in the next paragraph. For the group constraint, DL2 performs slightly better than any of the fuzzy logics, although the differences between the logics are overall not as noticeable compared to the class-similarity constraint.

**Hyperparameter  $\lambda$**  Learning with logical loss introduces the hyperparameter  $\lambda$ , which is the logical weight for the total loss calculation. Finding a suitable logical weight  $\lambda$  is crucial for achieving good results, as choosing a sub-optimal value can potentially result in operators that are supposed to perform badly (such as the Gödel implication) performing even better than logics that are supposed to perform best. Our strategy to approximate good values of  $\lambda$  was to run the same experiment (data set, constraint, logic) for the same number of epochs for each logical weight value  $\lambda \in \{0, 0.2, 0.4, 0.6, 0.8, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$ . We then selected the value that yielded the highest combined prediction and constraint accuracy. Tables 2a and 2b show the value of  $\lambda$  we chose for each run, and Figs. 4 to 6 (Appendix A) show the prediction and constraint accuracies for each value of  $\lambda$ .

This approach is very expensive and not feasible for real-world application. Unfortunately, extrapolating from running experiments at a smaller number of epochs does not necessarily transfer over to running

<sup>5</sup> The definition of the set of sets Groups is shown in Eq. (8) (Appendix A).

Table 2: Results. P/C is prediction / constraint accuracy in %.

	(a) Class-Similarity constraint.						(b) Group constraint.			
	Fashion-MNIST			CIFAR-10			GTSRB			
	P	C	$\lambda$	P	C	$\lambda$	P	C	$\lambda$	
Baseline	77.55	84.31	–	79.06	48.65	–	Baseline	89.93	49.97	–
DL2	77.88	89.15	0.6	78.55	52.21	0.4	DL2	<b>88.16</b>	<b>77.25</b>	7.0
$I_G$	63.46	91.30	3.0	77.65	81.56	1.2	$T_G$	88.38	74.20	5.0
$I_{KD}$	75.39	80.59	0.8	78.82	72.94	0.6	$T_{LK}$	85.26	78.43	5.0
$I_{LK}$	64.64	97.28	4.0	76.06	87.75	6.0	$T_{RC}$	86.52	77.56	5.0
$I_{GG}$	67.25	95.54	3.0	74.83	88.67	10.0	$T_{YG}$	87.47	76.47	5.0
$I_{RC}$	76.79	92.56	0.8	79.14	80.51	0.8				
$(I_{RC})_{s=9}$	77.06	93.63	0.8	78.30	80.87	0.8				
$(I_{RC})_{\phi=x^2}$	<b>76.88</b>	<b>95.94</b>	1.0	<b>78.31</b>	<b>90.74</b>	1.6				
$I_{YG}$	74.14	80.15	1.0	77.81	73.19	0.8				

the experiment at the desired number of epochs. Further, as can be seen in Figs. 4 to 6 (Appendix A), the resulting graphs are non-monotonic, making it difficult to predict the prediction and constraint accuracy one would obtain with other values of  $\lambda$ .

## 4 Future Work

Our experiments have shown that learning with differentiable logics can generally improve how much a ML model satisfies a constraint. Imposing logical constraints on the training process in this manner could be a step in the direction of verified ML, allowing the use of continuous-learning in self-improving ML-enabled autonomous systems. It has to be noted that in contrast to formal verifiers, training with logical loss does not formally guarantee properties to hold in all possible cases.

We highlight a few more areas for future work in the following.

**Reusing logical constraints during inference** Because of the differentiable logics acting as a regulariser during training, any logical constraints imposed on the learning process are unavailable during inference. The trained model can therefore not make use of the logical constraints to check its predictions, for example to attach confidence scores to its predictions.

**Probabilistic logics** Despite expressing satisfaction of formulas on  $[0, 1]$ , fuzzy logics are inherently not probabilistic, having been designed instead for reasoning in the presence of vagueness. We point to DeepProbLog [9] as one example for a probabilistic logic for use with deep learning. In the context of neural networks, which often output probabilities, it could be more natural to reason about probabilities instead of vagueness, especially for constraints that include probabilities [3].

**Properties** A common problem with verifying ML is the lack of specifications, as noted by [8, 3]. Most properties in the literature are limited to robustness against slight perturbations, although differentiable logics can not only relate network inputs and outputs, but could also refer to the inner workings of the neural network, such as weights and activation states. A related area is to investigate whether learning with logical loss can be used to show that desired properties continue to hold when retraining the network.



## References

- [1] Michał Baczyński & Balasubramaniam Jayaram (2008): *Fuzzy Implications*. *Studies in Fuzziness and Soft Computing* v. 231, Springer Verlag, Berlin, doi:10.1007/978-3-540-69082-5.
- [2] Chih-Hong Cheng & Rongjie Yan (2021): *Continuous Safety Verification of Neural Networks*. In: *2021 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 1478–1483, doi:10.23919/DATE51398.2021.9473994.
- [3] Marie Farrell, Anastasia Mavridou & Johann Schumann (2023): *Exploring Requirements for Software that Learns: A Research Preview*. In Alessio Ferrari & Birgit Penzenstadler, editors: *Requirements Engineering: Foundation for Software Quality - 29th International Working Conference, REFSQ 2023, Barcelona, Spain, April 17-20, 2023, Proceedings, Lecture Notes in Computer Science* 13975, Springer, pp. 179–188, doi:10.1007/978-3-031-29786-1\_12.
- [4] Claudio Ferrari, Mark Niklas Muller, Nikola Jovanovic & Martin Vechev (2022): *Complete Verification via Multi-Neuron Relaxation Guided Branch-and-Bound*, doi:10.48550/arXiv.2205.00263.
- [5] Marc Fischer, Mislav Balunovic, Dana Drachler-Cohen, Timon Gehr, Ce Zhang & Martin Vechev (2019): *DL2: Training and Querying Neural Networks with Logic*. In: *Proceedings of the 36th International Conference on Machine Learning*, PMLR, pp. 1931–1941.
- [6] Guy Katz, Clark Barrett, David L. Dill, Kyle Julian & Mykel J. Kochenderfer (2017): *Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks*. In Rupak Majumdar & Viktor Kunčák, editors: *Computer Aided Verification*, Lecture Notes in Computer Science, Springer International Publishing, Cham, pp. 97–117, doi:10.1007/978-3-319-63387-9\_5.
- [7] M.Z. Kwiatkowska (2019): *Safety Verification for Deep Neural Networks with Provable Guarantees*. In: *Leibniz International Proceedings in Informatics, LIPIcs*, 140, doi:10.4230/lipics.concur.2019.1.
- [8] Martin Leucker (2020): *Formal Verification of Neural Networks?* In Gustavo Carvalho & Volker Stolz, editors: *Formal Methods: Foundations and Applications*, Lecture Notes in Computer Science, Springer International Publishing, Cham, pp. 3–7, doi:10.1007/978-3-030-63882-5\_1.
- [9] Robin Manhaeve, Sebastijan Dumancic, Angelika Kimmig, Thomas Demeester & Luc De Raedt (2018): *DeepProbLog: Neural Probabilistic Logic Programming*. In: *Advances in Neural Information Processing Systems*, 31, Curran Associates, Inc.
- [10] Mark Niklas Müller, Christopher Brix, Stanley Bak, Changliu Liu & Taylor T. Johnson (2022): *The Third International Verification of Neural Networks Competition (VNN-COMP 2022): Summary and Results*, doi:10.48550/arXiv.2212.10376.
- [11] Natalia Ślusarz, Ekaterina Komendantskaya, Matthew Daggitt, Robert Stewart & Kathrin Stark (2023): *Logic of Differentiable Logics: Towards a Uniform Semantics of DL*. In: *EPiC Series in Computing*, 94, EasyChair, pp. 473–493, doi:10.29007/c1nt.
- [12] Hoang-Dung Tran, Xiaodong Yang, Diego Manzananas Lopez, Patrick Musau, Luan Viet Nguyen, Weiming Xiang, Stanley Bak & Taylor T. Johnson (2020): *NNV: The Neural Network Verification Tool for Deep Neural Networks and Learning-Enabled Cyber-Physical Systems*. In Shuvendu K. Lahiri & Chao Wang, editors: *Computer Aided Verification*, Lecture Notes in Computer Science, Springer International Publishing, Cham, pp. 3–17, doi:10.1007/978-3-030-53288-8\_1.
- [13] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner & Aleksander Madry (2018): *Robustness May Be at Odds with Accuracy*. In: *International Conference on Learning Representations*.
- [14] Emile van Krieken, Erman Acar & Frank van Harmelen (2022): *Analyzing Differentiable Fuzzy Logic Operators*. *Artificial Intelligence* 302, p. 103602, doi:10.1016/j.artint.2021.103602. arXiv:2002.06100.
- [15] Peter Varnai & Dimos V. Dimarogonas (2020): *On Robustness Metrics for Learning STL Tasks*. In: *2020 American Control Conference (ACC)*, pp. 5394–5399, doi:10.23919/ACC45564.2020.9147692.

- [16] Shiqi Wang, Kexin Pei, Justin Whitehouse, Junfeng Yang & Suman Jana (2018): *Efficient Formal Safety Analysis of Neural Networks*. In: *Advances in Neural Information Processing Systems*, 31, Curran Associates, Inc.

## A Appendix

### A.1 Constraints

The full set of labels used in the class-similarity constraint for Fashion-MNIST is shown in Eq. (6), the set of labels for CIFAR-10 is shown in Eq. (7). Note that we have a conjunct for each class so as to rule out cases where the implication would vacuously be true, which does not necessarily reflect real world constraints.

$$\text{Labels}_{\text{Fashion-MNIST}} := \left\{ \begin{array}{l} (\text{T-shirt/top, Shirt, Ankle boot}), \\ (\text{Trouser, Dress, Bag}), \\ (\text{Pullover, Shirt, Sandal}), \\ (\text{Dress, Coat, Bag}), \\ (\text{Coat, Pullover, Shirt}), \\ (\text{Sandal, Sneaker, Dress}), \\ (\text{Shirt, Pullover, Sneaker}), \\ (\text{Sneaker, Sandal, Trouser}), \\ (\text{Bag, Sandal, Dress}), \\ (\text{Ankle boot, Sneaker, T-shirt/top}) \end{array} \right\} \quad (6)$$

$$\text{Labels}_{\text{CIFAR-10}} := \left\{ \begin{array}{l} (\text{airplane, ship, dog}), \\ (\text{automobile, truck, cat}), \\ (\text{bird, airplane, dog}), \\ (\text{cat, dog, frog}), \\ (\text{deer, horse, truck}), \\ (\text{dog, cat, bird}), \\ (\text{frog, ship, truck}), \\ (\text{horse, deer, airplane}), \\ (\text{ship, airplane, deer}), \\ (\text{truck, automobile, airplane}) \end{array} \right\} \quad (7)$$

The definitions of the groups used in the group constraint Eq. (5) for GTSRB are given in Eq. (8).

$$\begin{aligned} \text{Group}_{\text{Speed Limits}} &:= \left\{ \begin{array}{l} \text{limit 20km/h,} \\ \text{limit 30km/h,} \\ \text{limit 50km/h,} \\ \text{limit 60km/h,} \\ \text{limit 70km/h,} \\ \text{limit 80km/h,} \\ \text{end of limit 80km/h,} \\ \text{limit 100km/h,} \\ \text{limit 120km/h} \end{array} \right\}, & \text{Group}_{\text{Mandatory Actions}} &:= \left\{ \begin{array}{l} \text{turn right ahead,} \\ \text{turn left ahead,} \\ \text{ahead only,} \\ \text{go straight or right,} \\ \text{go straight or left,} \\ \text{keep right,} \\ \text{keep left,} \\ \text{roundabout} \end{array} \right\} \\ \text{Group}_{\text{Prohibitions}} &:= \left\{ \begin{array}{l} \text{no passing,} \\ \text{no passing for trucks,} \\ \text{no way,} \\ \text{no way one-way,} \\ \text{end of no passing,} \\ \text{end of no passing for trucks} \end{array} \right\}, & \text{Group}_{\text{Warnings}} &:= \left\{ \begin{array}{l} \text{caution general,} \\ \text{caution curve left,} \\ \text{caution curve right,} \\ \text{caution curvy,} \\ \text{caution bumps,} \\ \text{caution slippery,} \\ \text{caution narrow road,} \\ \text{road work,} \\ \text{pedestrians,} \\ \text{children crossing,} \\ \text{wild animals crossing} \end{array} \right\} \\ \text{Groups}_{\text{GTSRB}} &:= \{ \text{Group}_{\text{Speed Limits}}, \text{Group}_{\text{Prohibitions}}, \text{Group}_{\text{Mandatory Actions}}, \text{Group}_{\text{Warnings}} \} \quad (8) \end{aligned}$$

### A.2 Plots of Prediction and Constraint Accuracy Over Time

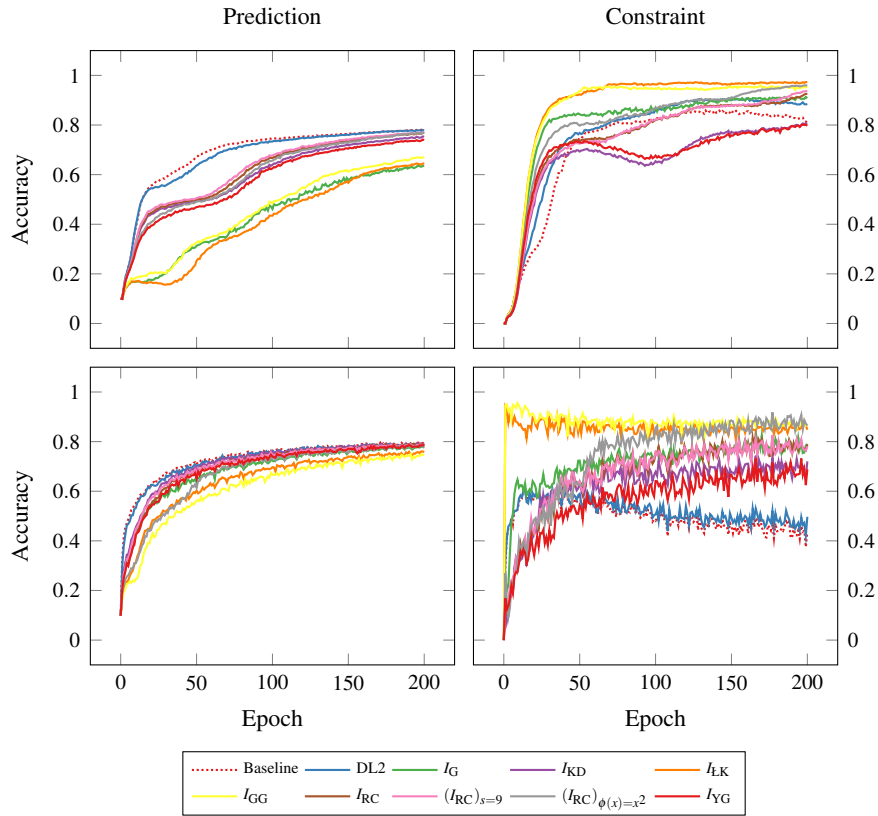


Figure 1: The figure shows how prediction accuracy (left column) and constraint accuracy (right column) change over time when training with the class-similarity constraint for 200 epochs with different logics on Fashion-MNIST (top row) and CIFAR-10 (bottom row).

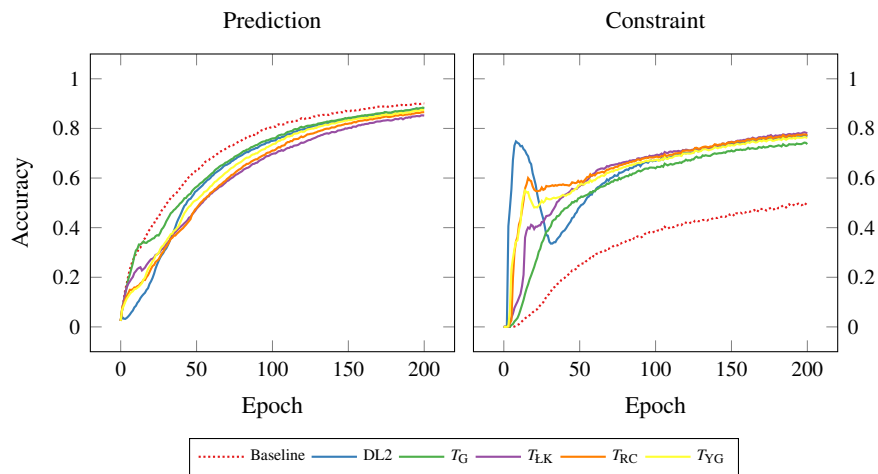


Figure 2: The figure shows how prediction accuracy (left column) and constraint accuracy (right column) change over time when training with the group constraint for 200 epochs with different logics on GTSRB.

### A.3 Runtime Overhead

Table 3: Average epoch train times in seconds.

(a) Class-Similarity constraint.			(b) Group constraint.	
	Fashion-MNIST	CIFAR-10	GTSRB	
Baseline	0.6 s	3.4 s	Baseline	1.1 s
DL2	0.6 s	3.4 s	DL2	1.1 s
$I_G$	0.6 s	3.8 s	$T_G$	1.1 s
$I_{KD}$	0.7 s	4.0 s	$T_{LK}$	1.1 s
$I_{LK}$	0.7 s	4.0 s	$T_{RC}$	1.1 s
$I_{GG}$	0.7 s	4.0 s	$T_{YG}$	1.2 s
$I_{RC}$	0.7 s	4.0 s		
$(I_{RC})_{s=9}$	0.7 s	4.1 s		
$(I_{RC})_{\phi=x^2}$	0.7 s	4.1 s		
$I_{YG}$	0.7 s	4.1 s		

The experiments were conducted on a system with a 4.5 GHz AMD Ryzen 9 77950X and a GeForce RTX 4090. The average epoch training times are shown in Table 3. In contrast to the original DL2 experiments [5], we did not observe significant overhead in our experiment. We note that this is due to our reuse of the model output from the cross-entropy loss calculation in the forward pass for our logical loss calculation. If our constraints required a separate forward pass, we would most likely see the same overhead as observed in the DL2 experiments.

### A.4 Plots of the Fuzzy Logic Implications

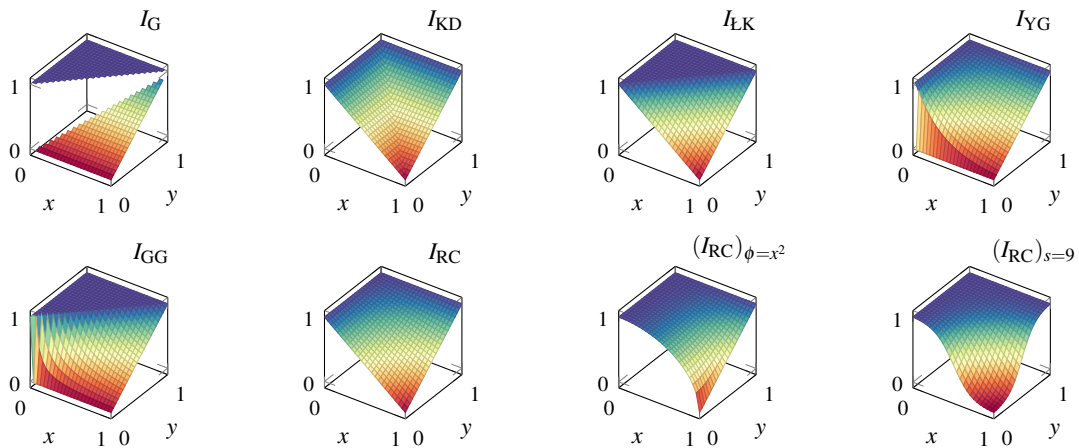


Figure 3: The fuzzy logic implications  $I(x,y)$  used in our experiments to map the logical statement  $x \rightarrow y$  into real-valued loss. Formal definitions are collected in Table 1.

### A.5 Plots of Prediction and Constraint Accuracy for Different Values of $\lambda$

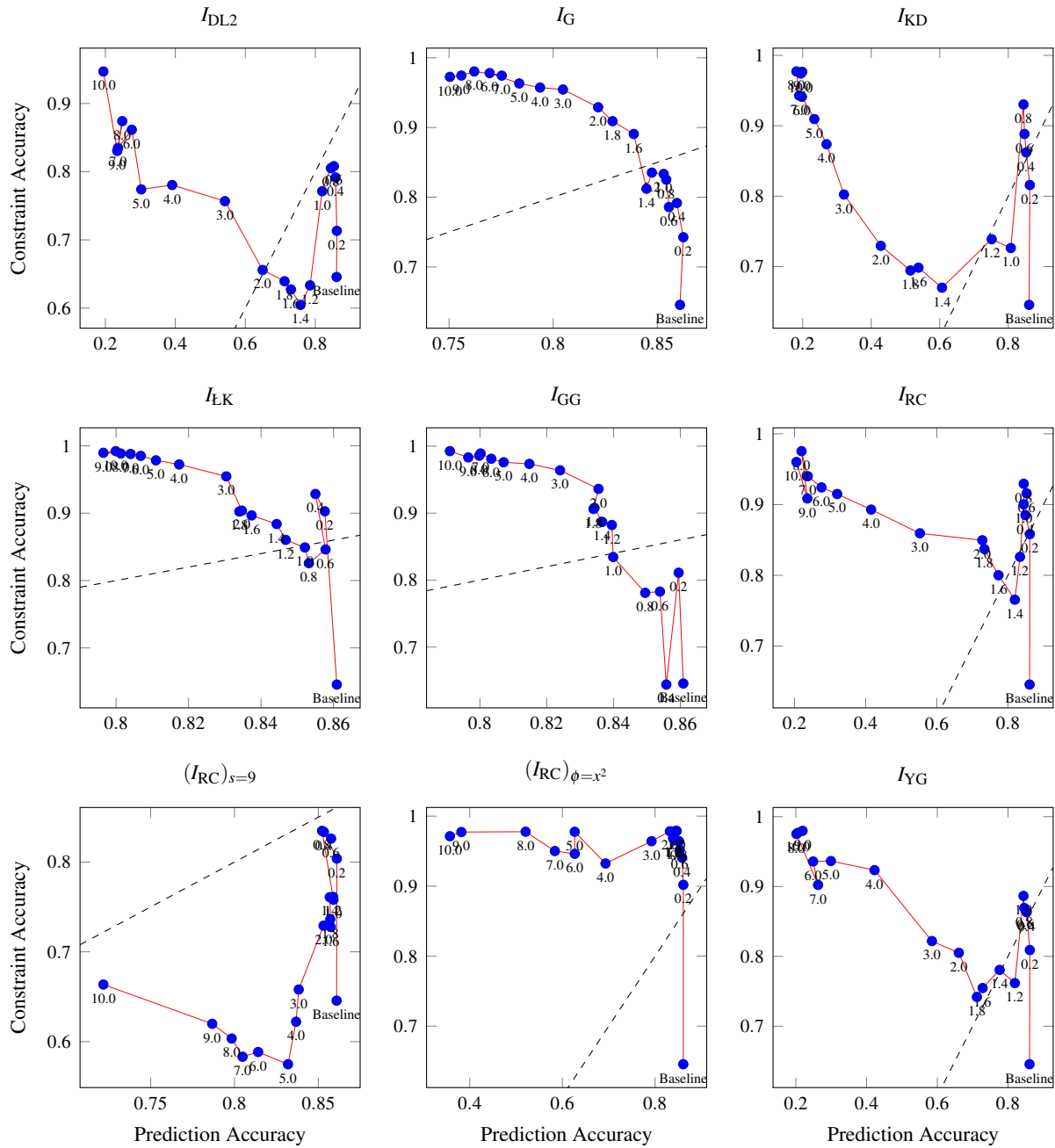


Figure 4: The figure displays prediction and constraint accuracy obtained when training with varying values of  $\lambda$  with the class-similarity constraint on Fashion-MNIST for 200 epochs.

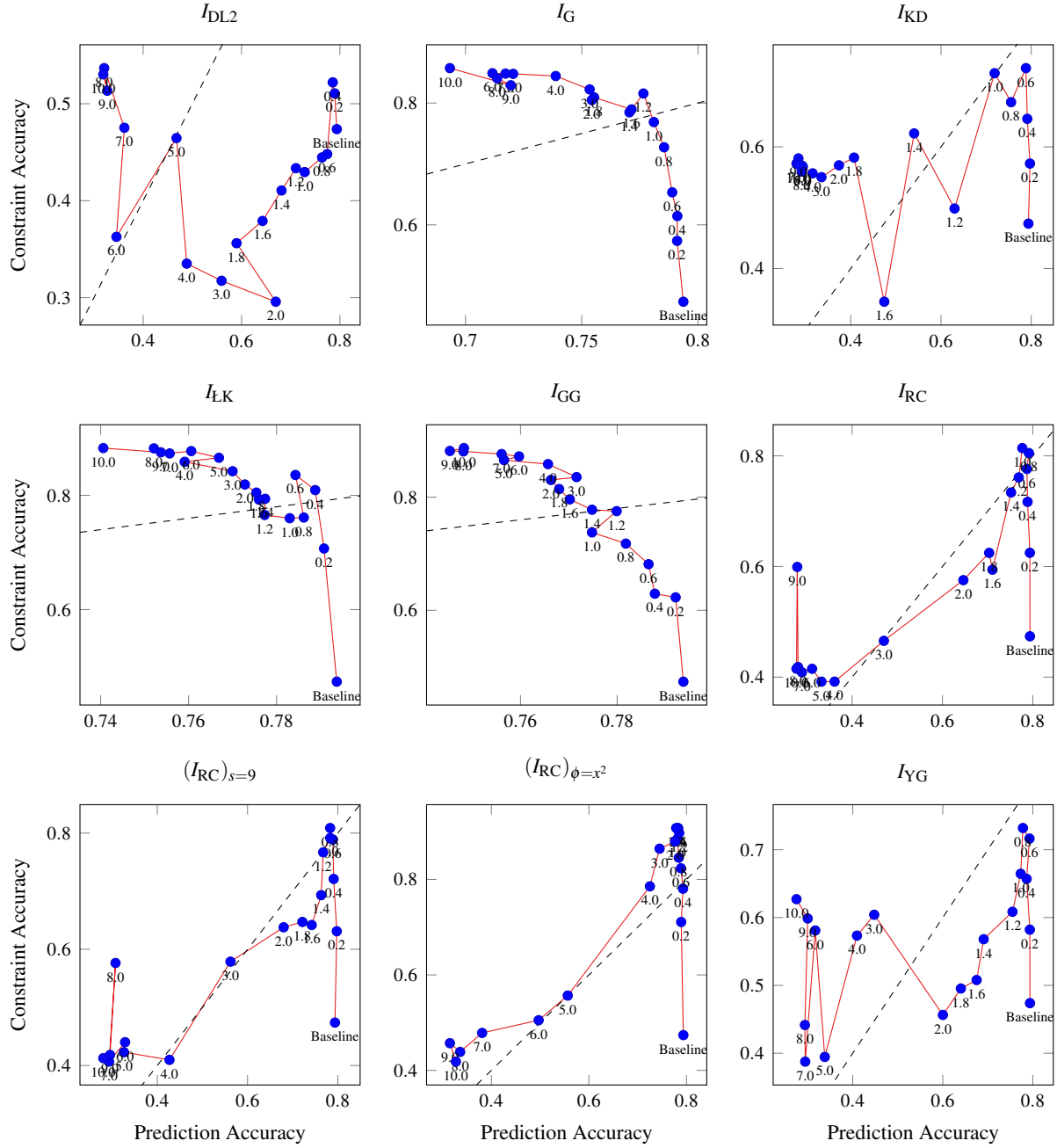


Figure 5: The figure displays prediction and constraint accuracy obtained when training with varying values of  $\lambda$  with the class-similarity constraint on CIFAR-10 for 200 epochs.

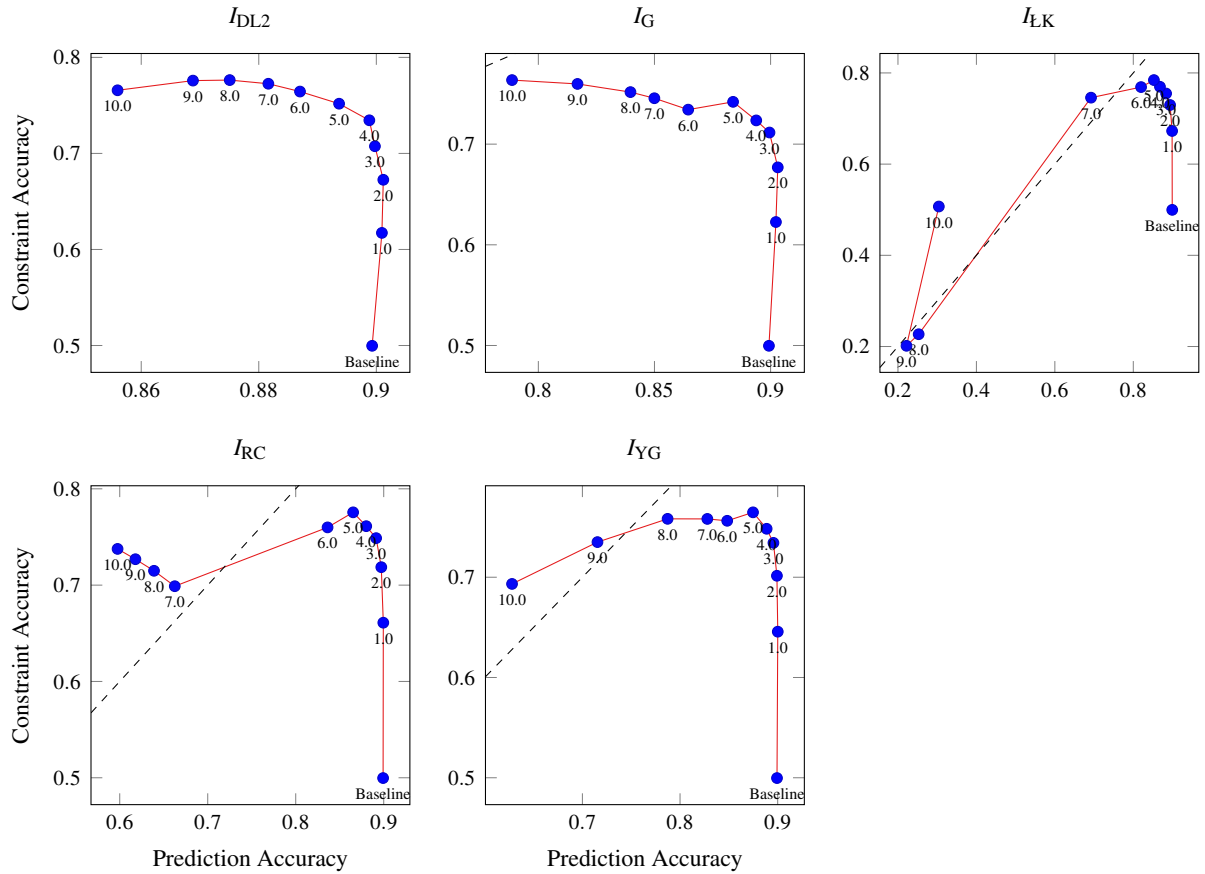


Figure 6: The figure displays prediction and constraint accuracy obtained when training with varying values of  $\lambda$  with the group constraint on GTSRB for 200 epochs.

# Extending Neural Network Verification to a Larger Family of Piece-wise Linear Activation Functions

László Antal

RWTH Aachen University  
Aachen, Germany

antal@cs.rwth-aachen.de

Hana Masara

RWTH Aachen University  
Aachen, Germany

hana.masara@rwth-aachen.de

Erika Ábrahám

RWTH Aachen University  
Aachen, Germany

abraham@cs.rwth-aachen.de

In this paper, we extend an available neural network verification technique to support a wider class of *piece-wise linear* activation functions. Furthermore, we extend the algorithms, which provide in their original form exact respectively over-approximative results for bounded input sets represented as star sets, to allow also *unbounded* input sets. We implemented our algorithms and demonstrated their effectiveness in some case studies.

## 1 Introduction

In the area of artificial intelligence, *feed-forward neural networks (FNNs)* [32] enjoy increasing popularity. FNNs can be trained to learn a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  from a set of input-output samples, and predict outputs also for previously unseen inputs. This way, FNNs can tackle problems that would otherwise require very complex solutions [45].

Nowadays, a wide range of applications use FNNs, such as autonomous vehicles [31], speech- and object-recognition systems [22, 13] or robot vision [33], just to mention a few. While FNNs are impressively effective, their reliability in safety-critical situations is still questionable [9, 16, 24]. Hence, verification methods play an important role in providing guarantees about their behavior. In this work, we focus on the *reachability problem* for FNNs, which is the problem of determining which output values (*reachable set*) an FNN computes for inputs from a given set.

**Related work.** The application of formal methods [61, 7, 20, 60] to verify the safety of neural networks began with [41]. Since then, the verification of neural networks has gained significant attention from the formal methods research community [55, 64, 12, 6, 59, 34, 14, 56, 4, 19, 28, 50].

Some of the available approaches encode the verification problems as logical formulae and use SMT-solvers for their solution [27, 63, 28, 12, 59]. Another common technique is reachable set calculation [23, 44, 36, 55, 64, 56] using an abstract representation like star sets [57] or symbolic intervals [29].

This paper builds on previous work [3, 54, 57], which solves the reachability problem using *star sets* to represent subsets of  $\mathbb{R}^k$  for any  $k \in \mathbb{N}$  with  $k > 0$ , like sets of input and output values. The authors present two methods, one with exact computations and one which over-approximates the reachable set.

**Contributions.** Our contributions in this paper are the following:

1. We extend the set of activation functions supported by [57, 54] to cover the piece-wise linear functions *leaky ReLU*, *hard tanh*, *hard sigmoid* and the *unit step*; while some of these functions have already been included in the respective algorithms, no complete formalizations were available, which we provide in this paper. Furthermore, we support more general, *parameterized* versions



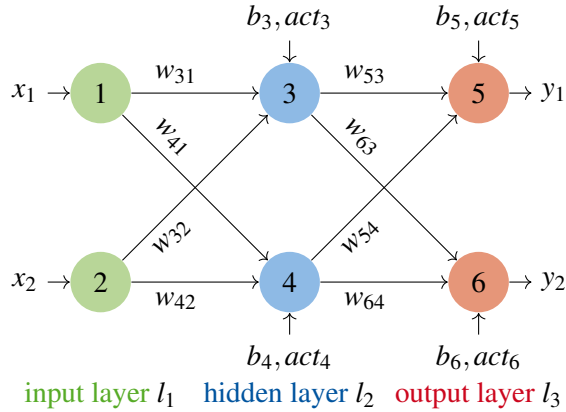


Figure 1: Illustration of a feed-forward fully-connected neural network, consisting of one input layer (green), one hidden layer (blue), and one output layer (red).

of the aforementioned activation functions. For each of the above, we present the reachability analysis algorithm using both the exact and the over-approximative methods.

2. While previous work was restricted to bounded input sets, we provide extensions to allow also *unbounded* input sets.
3. Using the open-source library HyPro<sup>1</sup> [47] for the star-set representation, we developed a C++ *implementation* of both the exact and the over-approximative analysis methods, covering all the above activation functions. This includes also an extension of HyPro with an NNET parser to input FNN models in NNET file format.
4. We propose some novel benchmarks (thermostat and sonar classifier) with the aim of supporting the formal methods community. Using our implementation, we provide *experimental evaluation* on the two proposed benchmarks and two other existing benchmarks discussing the results.

**Outline.** The rest of this paper is structured as follows. We present in Section 2 the fundamentals of this work, including feedforward neural networks (FNN), star sets, and reachability analysis of FNN with the rectified linear unit (ReLU) activation function. Then, in Section 3, we propose an exact and over-approximate analysis method for several other activation functions, considering both bounded and unbounded input sets. Afterwards, in Section 4, we present and evaluate experimental results on four different benchmarks. Finally, in Section 5 we conclude the paper and discuss future work.

## 2 Preliminaries

We use  $\mathbb{N}$  to denote the set of all natural numbers including 0 and  $\mathbb{R}$  for the reals, and consider elements from  $\mathbb{R}^n$  (for any  $n \in \mathbb{N}$ ) to be column vectors.

### 2.1 Feedforward Neural Networks

A *feedforward neural network (FNN)* [30, 51] is a directed weighted graph annotated with some data. It has a finite set of nodes called *neurons*, which are grouped into  $k \in \mathbb{N}_{\geq 2}$  disjoint non-empty ordered sets  $l_1, \dots, l_k$  called *layers*. We call  $l_1$  the *input layer*,  $l_k$  the *output layer*, while the others are *hidden layers*. Let  $\langle i \rangle$  denote the size  $|l_i|$  of layer  $i = 1, \dots, k$ . There is a directed edge from each neuron  $n$  in each

<sup>1</sup>Implementation available online at <https://github.com/hypro/hypro>. For reproducing the experimental results, please check the Case Studies/Neural Network Verification subsection of HyPro’s GitHub page: see the README.md file.

non-output layer  $l_{i-1}$  to each neuron  $n'$  in the next layer  $l_i$ , weighted by  $w_{n',n} \in \mathbb{R}$ ; let  $\mathbf{W}_i \in \mathbb{R}^{(i) \times (i-1)}$  be the matrix whose entry in row  $r$  and column  $c$  is the weight of the edge from the  $c$ th neuron in layer  $i-1$  to the  $r$ th neuron in layer  $i$ . In addition, each neuron  $n$  in each non-input layer is annotated with a *bias*  $b_n \in \mathbb{R}$  and an *activation function*  $act_n : \mathbb{R} \rightarrow \mathbb{R}$ ; for layer  $i$  with neurons  $l_i = \{n_1, \dots, n_{(i)}\}$ , let  $\mathbf{b}_i = (b_{n_1}, \dots, b_{n_{(i)}})^T$  and  $\mathbf{act}_i : \mathbb{R}^{(i)} \rightarrow \mathbb{R}^{(i)}$  with  $\mathbf{act}_i(\mathbf{y}) = (act_{n_1}(y_1), \dots, act_{n_{(i)}}(y_{(i)}))^T$  for any input  $\mathbf{y} = (y_1, \dots, y_{(i)})^T \in \mathbb{R}^{(i)}$ . A frequently used activation function is the *Rectified Linear Unit (ReLU)*, defined as  $ReLU(x) = \max(0, x)$  for  $x \in \mathbb{R}$ . An example FNN is shown in Figure 1.

For an *input*  $\mathbf{x}_1 = (x_1, \dots, x_{(1)}) \in \mathbb{R}^{(1)}$ , the *state*  $\mathbf{x}_i$  of each non-input layer  $l_i$  is defined recursively as

$$\mathbf{x}_i = \mathbf{act}_i(\mathbf{W}_i \mathbf{x}_{i-1} + \mathbf{b}_i). \quad (1)$$

Thus an FNN can be seen as a function  $f : \mathbb{R}^{(1)} \rightarrow \mathbb{R}^{(k)}$ , assigning to each input the output layer's state, which we call the *output*. For a given FNN and a set  $\mathcal{R}_1$  of possible inputs, the *FNN reachability problem* is the problem to compute all possible states for each of the layers  $1 < i \leq k$  [54]:

$$\mathcal{R}_i = \{ \mathbf{act}_i(\mathbf{W}_i \mathbf{x}_{i-1} + \mathbf{b}_i) \mid \mathbf{x}_{i-1} \in \mathcal{R}_{i-1} \}. \quad (2)$$

Solving the FNN reachability problem allows us to check properties of interest, e.g. safety properties (whether the output set is disjoint from a set of unsafe outputs) or stability (whether the distance between possible outputs is below a threshold for a given input set).

In this work, as input set we consider convex polyhedra  $\mathcal{R}_1 = \{ \mathbf{x} \in \mathbb{R}^{(1)} \mid \mathbf{A}\mathbf{x} \leq \mathbf{c} \}$  for some  $m \in \mathbb{N}_{\geq 1}$ ,  $\mathbf{A} \in \mathbb{R}^{m \times (1)}$  and column vector  $\mathbf{c} \in \mathbb{R}^m$ .

## 2.2 Stars

To compute  $\mathcal{R}_i$  via Equation 2, the two main operations that need to be applied on state sets are the activation function  $\mathbf{act}_i$  and *affine transformations* using the weights  $\mathbf{W}_i$  and biases  $\mathbf{b}_i$  of the layer  $i$ . For implementing these calculations efficiently, different state set representations have been proposed [48]. Under these, star sets (or short stars) turned out to be exceptionally good candidates, for their efficient handling of affine transformations and half-space intersections (see Propositions 2.2, 2.3 and 2.4).

For any  $n, m \in \mathbb{N}$ , an  $(n, m)$ -dimensional *star* is a tuple  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  of (i) a *center*  $\mathbf{c} \in \mathbb{R}^n$ , (ii) a *generator matrix*  $\mathbf{V} \in \mathbb{R}^{n \times m}$  whose columns  $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(m)} \in \mathbb{R}^n$  are called the *basis vectors* or *generators* and (iii) a *predicate*  $P \subseteq \mathbb{R}^m$ . The star  $\theta$  represents the set  $[\theta] = \{ \mathbf{c} + \sum_{j=1}^m (\alpha_j \mathbf{v}^{(j)}) \mid (\alpha_1, \dots, \alpha_m)^T \in P \}$ .

As in [54], we restrict  $P$  to be a convex polyhedron  $P = \{ \alpha \in \mathbb{R}^m \mid \mathbf{C}\alpha \leq \mathbf{d} \}$  for some  $p \in \mathbb{N}$ ,  $\mathbf{C} \in \mathbb{R}^{p \times m}$  and  $\mathbf{d} \in \mathbb{R}^p$ . The following star properties, whose proofs are included in Appendix A.1, will be used to solve the FNN reachability problem.

**Proposition 2.1** (Convex polyhedra). For any  $m, p \in \mathbb{N}$ ,  $\mathbf{C} \in \mathbb{R}^{p \times m}$  and  $\mathbf{d} \in \mathbb{R}^p$ , the convex polyhedron  $\mathcal{P} = \{ \mathbf{x} \in \mathbb{R}^m \mid \mathbf{C}\mathbf{x} \leq \mathbf{d} \}$  can be represented by a star.

**Proposition 2.2** (Affine transformation). Assume an  $(n, m)$ -dimensional star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  and let  $\mathbf{W} \in \mathbb{R}^{k \times n}$  and  $\mathbf{b} \in \mathbb{R}^k$ . Then the affine transformation  $\{ \mathbf{W}\mathbf{x} + \mathbf{b} \mid \mathbf{x} \in [\theta] \}$  of  $[\theta]$  is represented by  $\bar{\theta} = \langle \bar{\mathbf{c}}, \bar{\mathbf{V}}, P \rangle$  with  $\bar{\mathbf{c}} = \mathbf{W}\mathbf{c} + \mathbf{b}$  and  $\bar{\mathbf{V}} \in \mathbb{R}^{k \times m}$  with columns  $\mathbf{W}\mathbf{v}^{(1)}, \dots, \mathbf{W}\mathbf{v}^{(m)}$ .

**Proposition 2.3** (Intersection with halfspace). Assume an  $(n, m)$ -dimensional star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  and a half-space  $\mathcal{H} = \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{h}^T \mathbf{x} \leq g \}$  with some  $\mathbf{h} \in \mathbb{R}^n$  and  $g \in \mathbb{R}$ . Then the intersection  $[\theta] \cap \mathcal{H}$  is represented by the star  $\bar{\theta} = \langle \mathbf{c}, \mathbf{V}, P \cap P' \rangle$  with  $P' = \{ \alpha \in \mathbb{R}^m \mid (\mathbf{h}^T \mathbf{V})\alpha \leq g - \mathbf{h}^T \mathbf{c} \}$ .

**Proposition 2.4** (Emptiness check). A star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  is empty if and only if  $P$  is empty.

**Proposition 2.5** (Bounding box). Assume an  $(n, m)$ -dimensional star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  with  $\mathbf{c} = (\mathbf{c}_1, \dots, \mathbf{c}_n)^T$ , and let  $\mathbf{V}_{(i)}$  be the  $i^{\text{th}}$  row of  $\mathbf{V}$ . Let furthermore  $B = \{ (x_1, \dots, x_n)^T \in \mathbb{R}^n \mid \bigwedge_{i=1}^n lb_i \leq x_i \leq ub_i \}$  with  $lb_i = \mathbf{c}_i + \min_{\alpha \in P} \mathbf{V}_{(i)}\alpha$  and  $ub_i = \mathbf{c}_i + \max_{\alpha \in P} \mathbf{V}_{(i)}\alpha$  for  $i = 1, \dots, n$ . Then  $[\theta] \subseteq B$ .

### 2.3 Reachability Analysis for FNNs with ReLUs

Next, we present two algorithms proposed in [54] to solve the reachability problem for FNNs with the ReLU activation function for bounded polyhedral input sets. The first algorithm is exact and thus complete, whereas the second algorithm over-approximates reachability. We note that alongside ReLU, [57] includes some other activation functions but no complete formalizations were available. In Section 3, we will extend these algorithms to support further and more general piece-wise linear activation functions and unbounded input sets.

#### Exact Analysis

The exact algorithm first constructs a star from the input set which is required to be a polyhedron (see Proposition 2.1). Then, correspondingly to Equation 2 it propagates the star through the network, layer-by-layer, until we get the output set  $\mathcal{R}_k$ . This propagation involves two main operations.

(1) For each non-input layer  $i$  and each star representing possible states of the previous layer, to compute the reachable states of layer  $i$ , we first apply an affine transformation on the star, using the weight matrix  $\mathbf{W}_i$  and the bias vector  $\mathbf{b}_i$ . Thus, from a star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  we obtain a new star  $\theta' = \langle \mathbf{c}', \mathbf{V}', P \rangle$  with  $\mathbf{c}' = \mathbf{W}_i \mathbf{c} + \mathbf{b}_i$  and  $\mathbf{V}' = \mathbf{W}_i \mathbf{V}$  (see Proposition 2.2). Note that during the affine transformation the predicate does not change.

(2) Then the non-linear activation function is applied on the intermediate star  $\theta'$  dimension-wise to represent  $\mathcal{R}_i = \mathbf{act}_{n_{(i)}}(\dots \mathbf{act}_{n_1}([\theta']) \dots)$ , where,  $n_1, \dots, n_{(i)}$  are the neurons in layer  $i$ . Since we consider the ReLU activation function, the  $\mathbf{act}_{n_j}(\cdot)$  operation at neuron  $n_j$  is defined as  $\text{ReLU}(x_j) = \max(0, x_j)$ ; instead of  $\mathbf{act}_{n_j}(\cdot)$  we also write  $\mathbf{act}_j^{\text{R}}(\cdot)$  to denote that the ReLU function is applied in dimension  $j$  (i.e. at the  $j$ th neuron of a layer). To compute  $\mathbf{act}_j^{\text{R}}(\theta)$  for a star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$ , the star  $\theta$  is decomposed into two stars  $\theta_1 = \langle \mathbf{c}, \mathbf{V}, P_1 \rangle$  and  $\theta_2 = \langle \mathbf{c}, \mathbf{V}, P_2 \rangle$  such that  $[\theta_1] = [\theta] \cap \{(x_1, \dots, x_n) \in \mathbb{R}^n \mid x_j < 0\}$  and  $[\theta_2] = [\theta] \cap \{(x_1, \dots, x_n) \in \mathbb{R}^n \mid x_j \geq 0\}$  (see Proposition 2.3). On the negative branch, i.e., when  $x_j < 0$ , the ReLU function sets the corresponding values to zero. Thus all the resulting elements of the star  $\theta_1$  should have the value zero in dimension  $j$ . It affects the star as a projection to 0 in dimension  $j$ . We can obtain this result by applying the mapping matrix  $\mathbf{M} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{j-1}, \mathbf{0}, \mathbf{e}_{j+1}, \dots, \mathbf{e}_n]$  on  $\theta_1$ , where  $\mathbf{e}_i \in \mathbb{R}^n$  is the  $i$ th  $n$ -dimensional unit vector (with 1 at position  $i$  and 0s otherwise). On the positive branch  $x_j \geq 0$ , the ReLU function does not change the set elements of  $\theta_2$ . Thus, the application of ReLU results in the union of two stars  $\mathbf{act}_j^{\text{R}}(\theta) = \langle \mathbf{M}\mathbf{c}, \mathbf{M}\mathbf{V}, P_1 \rangle \cup \langle \mathbf{c}, \mathbf{V}, P_2 \rangle$ . Note that if the values in  $[\theta]$  in the given dimension  $j$  are purely positive or purely negative, then the result of  $\mathbf{act}_j^{\text{R}}(\theta)$  is just a single star.

#### Over-approximate Analysis

While the exact algorithm is complete, it suffers from scalability issues since the number of stars grows during the analysis *exponentially* with the number of neurons. To tackle this problem, one solution is to side-step to over-approximative computations, which makes the analysis more *scalable*, however, it sacrifices the *completeness* of the method.

The over-approximate method from [54] also builds on Equation 2, but the application of the activation functions is different: the original  $\mathbf{act}_j^{\text{R}}(\cdot)$  operation is replaced by an over-approximating  $\overline{\mathbf{act}}_j^{\text{R}}(\cdot)$  which produces only a single star as output as follows. A new variable  $\alpha_{m+1}$  and three more constraints are added to the predicate  $P$  of the star, with the purpose of capturing the over-approximation of the ReLU function at neuron  $n_j$  (see Figure 2).

The three new constraints are:  $\alpha_{m+1} \geq 0$ ,  $\alpha_{m+1} \geq x_j$ , and  $\alpha_{m+1} \leq \frac{u(x_j-l)}{u-l}$ , where  $l$  and  $u$  are the lower and upper bounds, respectively, for variable  $x_j$  in  $[\theta]$  (see Proposition 2.5). Finally, since we want the variable  $\alpha_{m+1}$  to hold the over-approximation of  $x_j$ , after introducing the new variable and constraints to the predicate, we need to update the center  $\mathbf{c}$  and basis  $\mathbf{V}$  of the star  $\theta$  correspondingly. First, the old values of  $x_j$  are projected out using the mapping matrix  $\mathbf{M} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{j-1}, \mathbf{0}, \mathbf{e}_{j+1}, \dots, \mathbf{e}_n]$ . Then, a new generator vector  $\mathbf{e}_j$  is added to the basis, to link  $x_j$  to  $\alpha_{m+1}$ .

Formally, for an  $(n, m)$ -dimensional star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  we define  $\overline{\text{act}}_j^R(\theta) = \langle \bar{\mathbf{c}}, \bar{\mathbf{V}}, \bar{P} \rangle$ , where  $\bar{\mathbf{c}} = \mathbf{M}\mathbf{c}$ ,  $\bar{\mathbf{V}} = [\mathbf{M}\mathbf{v}^{(1)}, \mathbf{M}\mathbf{v}^{(2)}, \dots, \mathbf{M}\mathbf{v}^{(m)}, \mathbf{e}_j]$  and  $\bar{P} = \{(\alpha_1, \dots, \alpha_{m+1}) \in \mathbb{R}^{m+1} \mid (\alpha_1, \dots, \alpha_m) \in P \wedge \alpha_{m+1} \geq 0 \wedge \alpha_{m+1} \geq x_j \wedge \alpha_{m+1} \leq \frac{u(x_j-l)}{u-l}\}$ .

In case  $l \geq 0$  or  $u \leq 0$ , the introduction of a new variable is not necessary and we can proceed in a similar way as in the exact case, i.e., for positive domain we keep the set as it is, for negative domain we project out the variable  $x_j$ . Note that this over-approximation method is the least conservative that we can achieve using convex, linear constraints.

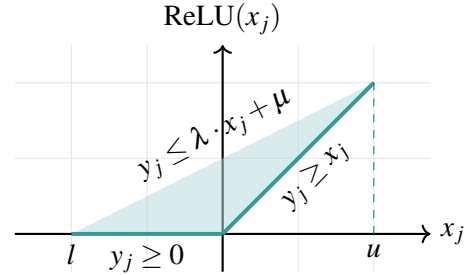


Figure 2: Relaxation of the ReLU function with  $\lambda = \frac{u}{u-l}$  and  $\mu = -\frac{lu}{u-l}$  [49]. Dark lines represent the exact set, the light area shows the approximate set.

### 3 FNN Reachability Analysis for Piece-wise Linear Activation Functions

Neural networks offer flexibility in choosing different activation functions. In this work, we present the extension of the reachability analysis algorithm to implement the *leaky rectified linear unit (leaky ReLU)*, *hard hyperbolic tangent (HardTanh)*, *hard sigmoid (HardSigmoid)*, and *unit step* activation functions. Below we define each of these functions and their application to a given star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$ .

#### 3.1 Unbounded Input Sets

During the analysis of an FNN, it may happen that one or more variables  $x_j$  of a star  $\theta$  become unbounded. That is, it has no lower bound (i.e.,  $l = -\infty$ ) or it has no upper bound (i.e.,  $u = \infty$ ). In the following, we present how to handle unbounded input sets as well.

Essentially, the exact reachability analysis of any piece-wise linear activation function presented in this paper does not change in case of unbounded input sets. The same steps are applied as per the exact analysis of bounded sets, i.e., (1) splitting the input set into multiple subsets based on the cases of the activation function, and (2) applying the corresponding transformations for each subset.

Conversely, in case of unbounded input, the over-approximate analysis does work differently, since the convex relaxations presented for bounded input need to be changed. In the rest of this paper, for each activation function, we show how the convex relaxations can be adjusted to achieve the *tightest* possible relaxation in case of unbounded inputs. Note that we distinguish for each function three cases of unboundedness of a variable  $x_j$ , either it has no lower bound ( $l = -\infty$  and  $u \in \mathbb{R}$ ), it has no upper bound ( $l \in \mathbb{R}$  and  $u = \infty$ ), or it has neither of the bounds ( $l = -\infty$  and  $u = \infty$ ).

Our implementation currently does not support unbounded input sets, so the presented methods for unbounded inputs are only theoretical results. Furthermore, the evaluated benchmarks also do not utilize unbounded sets.

### 3.2 Leaky ReLU Layer

Due to the dead neuron problem [10, 42] caused by the ReLU function, its alternative, the leaky ReLU function proposed by Mass et al.[37], is used in many applications.

**Definition 3.1** (Leaky ReLU [65]). The *leaky ReLU* activation function with scaling parameter  $\gamma \in (0, 1) \subset \mathbb{R}$  is defined for each  $x \in \mathbb{R}$  as

$$\text{LeakyReLU}(x) = \max(\gamma \cdot x, x) = \begin{cases} x & \text{if } x > 0 \\ \gamma \cdot x & \text{otherwise.} \end{cases} \quad (3)$$

#### Exact Analysis

The application of the leaky ReLU activation function is similar to the previously presented algorithm for the ReLU activation function, but they handle the negative inputs differently: While the ReLU function completely projects the input to zero, the leaky ReLU just scales the input down by  $\gamma \in (0, 1)$ . Thus, the application  $\text{act}_j^L(\theta)$  of leaky ReLU on a star  $\theta$  can be computed as follows. First we split the star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  into two subsets  $\theta_1 = \langle \mathbf{c}, \mathbf{V}, P_1 \rangle$  and  $\theta_2 = \langle \mathbf{c}, \mathbf{V}, P_2 \rangle$  with negative resp. non-negative  $x_j$ -values. Then we apply the corresponding transformations for both subsets. As previously, in the case of the positive subset  $\theta_2$ , no transformation is needed, since the leaky ReLU acts as an identity function for positive inputs. However, in case of the negative subset  $\theta_1$ , we apply the scaling matrix  $\mathbf{M} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \gamma \mathbf{e}_j, \dots, \mathbf{e}_{n-1}, \mathbf{e}_n]$ . Thus, the final result of the  $\text{act}_j^L(\cdot)$  operation at neuron  $n_j$  is the union of two stars:  $\text{act}_j^L(\theta) = \langle \mathbf{M}\mathbf{c}, \mathbf{M}\mathbf{V}, P_1 \rangle \cup \langle \mathbf{c}, \mathbf{V}, P_2 \rangle$ . The same observations apply here, that if the domain of a variable  $x_j$  is only negative (i.e.,  $u \leq 0$ ) or only positive (i.e.,  $l \geq 0$ ), the final result of the  $\text{act}_j^L(\cdot)$  operation is a single star: either  $\theta_1 = \langle \mathbf{M}\mathbf{c}, \mathbf{M}\mathbf{V}, P_1 \rangle$  or  $\theta_2 = \langle \mathbf{c}, \mathbf{V}, P_2 \rangle$ .

#### Over-approximate Analysis

The over-approximate analysis of the leaky ReLU is also similar to the one for ReLU. For bounded inputs, correspondingly to the Planet relaxation [49], we also try to find an enclosing triangle, which is the tightest convex, linear relaxation that we can achieve for leaky ReLUs (see Figure 3). The three constraints on the freshly introduced variable  $\alpha_{m+1}$  are the following: (1)  $\alpha_{m+1} \geq \gamma \cdot x_j$ , (2)  $\alpha_{m+1} \geq x_j$ , and (3)  $\alpha_{m+1} \leq \frac{u-\gamma l}{u-l}x_j + \frac{u \cdot l \cdot (\gamma-1)}{u-l}$ . At this point, the result of the  $\text{act}_j^L(\cdot)$  operation is a single star set with one more variable and three more constraints than the original input star. It is important to note: if the domain of variable  $x_j$  is fully positive (i.e.,  $l \geq 0$ ) or fully negative (i.e.,  $u \leq 0$ ), then the resulting star is the same as described for the exact approach. On the other hand, when there is an unbounded input set  $\theta$ , three cases are distinguished: (1)  $x_j \in (-\infty, u)$ , (2)  $x_j \in (l, \infty)$ , and (3)  $x_j \in (-\infty, \infty)$ . The analysis for unbounded input is similar to the bounded case but the introduced constraints change, as visualized in Figure 4. Note that these are the tightest linear, convex relaxations that can be achieved.

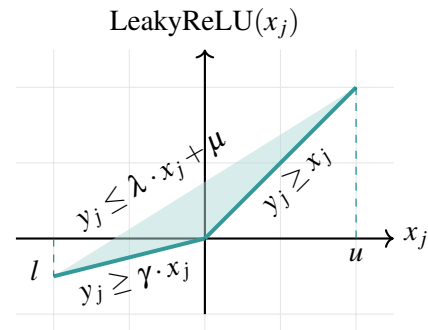


Figure 3: Relaxation for the leaky ReLU function. The dark line shows the exact set and the light area the approximate set. In the figure,  $\lambda = \frac{u-\gamma l}{u-l}$  and  $\mu = \frac{u \cdot l \cdot (\gamma-1)}{u-l}$ .

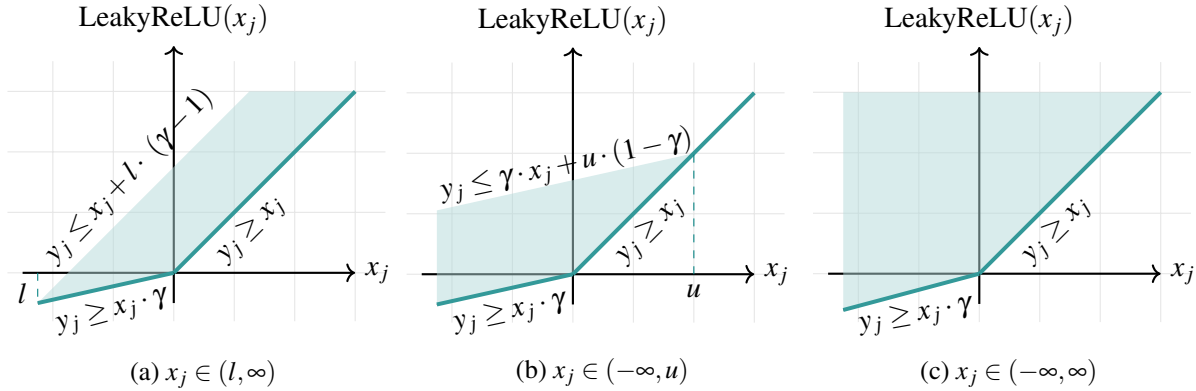


Figure 4: Convex relaxations of the leaky ReLU function with three cases of an unbounded input set.

### 3.3 Hard Tanh Layer

The hard hyperbolic tangent function, commonly known as the hard tanh function, is a linearized variant of the hyperbolic tangent activation function. In our work, we have generalized this function by introducing the parameters  $V_{\min}$  and  $V_{\max}$ , which replace the original values of  $-1$  and  $1$ , respectively [8]. This modification allows us to flexibly adapt the function according to our specific needs and requirements.

**Definition 3.2** (Hard Hyperbolic Tangent). The *hard hyperbolic tangent* (*HardTanh*) activation function with parameters  $V_{\min} \in \mathbb{R}$  and  $V_{\max} \in \mathbb{R}_{\geq V_{\min}}$  is defined for each  $x \in \mathbb{R}$  by

$$\text{HardTanh}(x) = \begin{cases} V_{\min} & \text{if } x < V_{\min} \\ x & \text{if } V_{\min} \leq x \leq V_{\max} \\ V_{\max} & \text{if } x > V_{\max} \end{cases} \quad (4)$$

#### Exact Analysis

For the analysis of FNNs with the hard tanh activation function at neuron  $n_j$ , which we denote as  $\text{act}_j^H(\cdot)$ , we split the result of the affine transformation  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  into three subsets:  $\theta_1 = \langle \mathbf{c}, \mathbf{V}, P_1 \rangle$  is the intersection of  $\theta$  with the hyperplanes  $V_{\min} \leq x_j \leq V_{\max}$ ,  $\theta_2 = \langle \mathbf{c}, \mathbf{V}, P_2 \rangle$  with  $x_j < V_{\min}$  and  $\theta_3 = \langle \mathbf{c}, \mathbf{V}, P_3 \rangle$  with  $x_j > V_{\max}$  (see Proposition 2.3).

According to Equation 4,  $\text{act}_j^H(\cdot)$  leaves the elements of  $\theta_1$  unchanged since  $x_j$  is in the range between  $V_{\min}$  and  $V_{\max}$ . For  $\theta_2$ , all of its elements get the value  $V_{\min}$  in dimension  $j$  since  $x_j < V_{\min}$ , hence, we project the star onto  $V_{\min}$  in the dimension  $j$ . To achieve this result, we apply the mapping matrix  $\mathbf{M} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{j-1}, \mathbf{0}, \mathbf{e}_{j+1}, \dots, \mathbf{e}_n]$ . Additionally, we set the  $j^{\text{th}}$  dimension of the center to  $V_{\min}$  by adding the shifting vector  $\mathbf{s}_{\min} = [0, \dots, V_{\min}, \dots, 0]^T$  to the center. For  $\theta_3$ , we do the same by mapping the set with the mapping matrix, but instead, we set the center to  $V_{\max}$  by adding the shifting vector  $\mathbf{s}_{\max} = [0, \dots, V_{\max}, \dots, 0]^T$  to it. Thus, we project the star onto  $V_{\max}$  in the dimension  $j$ . Accordingly, the  $\text{act}_j^H(\theta)$  operation at neuron  $j$  results in the union of three star sets:  $\text{act}_j^H(\cdot) = \langle \mathbf{c}, \mathbf{V}, P_1 \rangle \cup \langle \mathbf{M}\mathbf{c} + \mathbf{s}_{\min}, \mathbf{M}\mathbf{V}, P_2 \rangle \cup \langle \mathbf{M}\mathbf{c} + \mathbf{s}_{\max}, \mathbf{M}\mathbf{V}, P_3 \rangle$ .

Note that some of the intersections of the input star  $\theta$  with the halfspaces  $V_{\min} \leq x_j \leq V_{\max}$ ,  $x_j < V_{\min}$ , and  $x_j > V_{\max}$  may be empty (see Proposition 2.4). In that case, we can spare the computation for the empty subsets, and continue the reachability analysis only with the non-empty resulting stars.

### Over-approximate Analysis

In the over-approximate analysis, the  $\overline{\text{act}}_j^H(\theta)$  operation should yield a single star set. Thus we aim to find an enclosing triangle or trapezoid, which is the tightest convex, linear relaxation that we can achieve for hard tanh. For bounded inputs, we make a case distinction. If the lower bound (in the bounding box

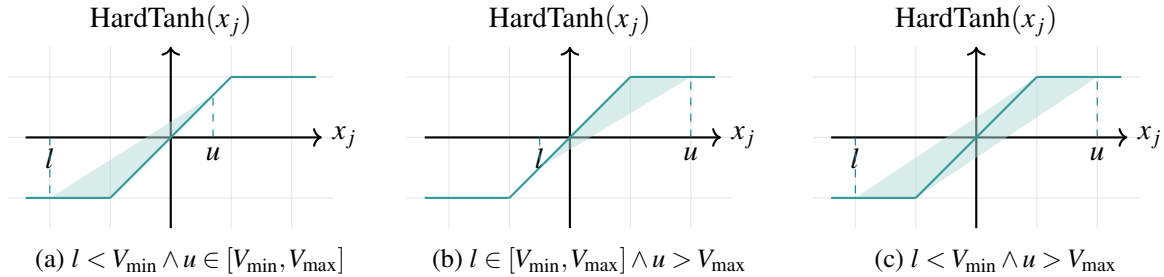


Figure 5: Relaxation for the hard tanh function. The dark line shows the exact set (non-convex) and the light area the approximate set (convex and linear).

of  $\theta$  in dimension  $j$ , see Proposition 2.5) is less than  $V_{\min}$ , and the upper bound is between  $V_{\min}$  and  $V_{\max}$ , the three constraints on the newly introduced variable  $\alpha_{m+1}$  are the following: (1)  $\alpha_{m+1} \geq x_j$ , (2)  $\alpha_{m+1} \geq V_{\min}$  and (3)  $\alpha_{m+1} \leq \frac{u_j - V_{\min}}{u_j - l_j} \cdot x_j - \frac{u_j \cdot (l_j - V_{\min})}{u_j - l_j}$ . For the opposite case, we introduce the new variable  $\alpha_{m+1}$  and three constraints: (1)  $\alpha_{m+1} \leq V_{\max}$ , (2)  $\alpha_{m+1} \leq x_j$  and (3)  $\alpha_{m+1} \geq -\frac{l_j - V_{\max}}{u_j - l_j} \cdot x_j - \frac{l_j \cdot (V_{\max} - u_j)}{u_j - l_j}$ . When the star is over  $V_{\min}$  and  $V_{\max}$  (i.e.,  $l < V_{\min} \wedge u > V_{\max}$ ), we introduce the new variable  $\alpha_{m+1}$  and four additional constraints: (1)  $\alpha_{m+1} \geq V_{\min}$ , (2)  $\alpha_{m+1} \leq V_{\max}$ , (3)  $\alpha_{m+1} \leq \frac{V_{\max} - V_{\min}}{V_{\max} - l_j} \cdot x_j - \frac{V_{\max} \cdot (l_j - V_{\min})}{V_{\max} - l_j}$  and (4)  $\alpha_{m+1} \geq \frac{V_{\min} - V_{\max}}{V_{\min} - u_j} \cdot x_j - \frac{V_{\min} \cdot (V_{\max} - u_j)}{V_{\min} - u_j}$ .

It is important to highlight that when the domain of variable  $x_j$  is between  $V_{\min}$  and  $V_{\max}$ , less than  $V_{\min}$  (i.e.,  $u < V_{\min}$ ) or greater than  $V_{\max}$  (i.e.,  $l > V_{\max}$ ), the result is again a single star and is computed the same way as described in the exact approach.

Furthermore, when dealing with an unbounded input set  $\theta$  we distinguish three cases, as mentioned earlier. These cases are as follows: (1)  $x_j \in (-\infty, u)$ , (2)  $x_j \in (l, \infty)$ , and (3)  $x_j \in (-\infty, \infty)$ . The cases (1) and (2) are again divided into two sub-cases, hence we obtain five different cases, each one presented in Table 1, coupled with the corresponding constraints and illustrations.

### 3.4 Hard Sigmoid Layer

The hard sigmoid activation function is a linearized variant of the sigmoid function. Since the hard sigmoid function has different variants in use [52, 2, 40], we generalize it by adding parameters.

**Definition 3.3** (Hard Sigmoid Function). The *hard sigmoid* (*HardSigmoid*) function with parameters  $V_{\min} \in \mathbb{R}$  and  $V_{\max} \in \mathbb{R}_{\geq V_{\min}}$  is defined for each  $x \in \mathbb{R}$  by

$$\text{HardSigmoid}(x) = \begin{cases} 0 & \text{if } x \leq V_{\min} \\ \frac{1}{V_{\max} - V_{\min}} \cdot x + \frac{V_{\min}}{V_{\min} - V_{\max}} & \text{if } V_{\min} < x < V_{\max} \\ 1 & \text{if } x \geq V_{\max} \end{cases} \quad (5)$$

Domain of $x_j$	Introduced constraints	Graphical illustration
$l = -\infty \wedge u \in [V_{\min}, V_{\max}]$	$\alpha_{m+1} \geq V_{\min}$ $\alpha_{m+1} \geq x_j$ $\alpha_{m+1} \leq u_j$	
$l = -\infty \wedge u > V_{\max}$	$\alpha_{m+1} \geq V_{\min}$ $\alpha_{m+1} \leq V_{\max}$ $\alpha_{m+1} \geq \frac{V_{\min} - V_{\max}}{V_{\min} - u_j} \cdot x_j - \frac{V_{\min} \cdot (V_{\max} - u_j)}{V_{\min} - u_j}$	
$l \in [V_{\min}, V_{\max}] \wedge u = \infty$	$\alpha_{m+1} \leq V_{\max}$ $\alpha_{m+1} \leq x_j$ $\alpha_{m+1} \geq l_j$	
$l < V_{\min} \wedge u = \infty$	$\alpha_{m+1} \leq V_{\max}$ $\alpha_{m+1} \geq V_{\min}$ $\alpha_{m+1} \leq \frac{V_{\max} - V_{\min}}{V_{\max} - l_j} \cdot x_j - \frac{V_{\max} \cdot (l_j - V_{\min})}{V_{\max} - l_j}$	
$l = -\infty \wedge u = \infty$	$\alpha_{m+1} \geq V_{\min}$ $\alpha_{m+1} \leq V_{\max}$	

Table 1: Approximation rules for the hard tanh function, when the input is unbounded. We distinguish five cases in total, for each we show the case itself, the introduced constraints and a graphical illustration.

## Exact Analysis

The analysis of the hard sigmoid works similarly to the one of the hard tanh function. The difference is that instead of the star remaining the same in the range between  $V_{\min}$  and  $V_{\max}$ , we scale the star according to Equation 5. To compute  $\mathbf{act}_j^S(\theta)$ , the star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  is partitioned into three subsets  $\theta_1$ ,  $\theta_2$  and  $\theta_3$ , covering the partitions with  $V_{\min} < x_j < V_{\max}$ ,  $x_j \leq V_{\min}$  respectively  $x_j \geq V_{\max}$ . We scale  $\theta_1$  by applying the scaling matrix  $\mathbf{M}_{sc} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \frac{1}{V_{\max} - V_{\min}} \mathbf{e}_j, \dots, \mathbf{e}_{n-1}, \mathbf{e}_n]$  and shift the center with the translation vector  $\mathbf{s}_{sc} = [0, \dots, \frac{V_{\min}}{V_{\min} - V_{\max}}, \dots, 0]^T$ . Furthermore, the elements of  $\theta_2$  are set to zero in dimension  $j$  by applying the mapping matrix  $\mathbf{M} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{j-1}, \mathbf{0}, \mathbf{e}_{j+1}, \dots, \mathbf{e}_n]$ . Finally, the elements of  $\theta_3$  are set to one by using the same projection matrix  $\mathbf{M}$ , plus setting the center to one by the shifting vector  $\mathbf{s}_{one} = [0, \dots, 1, \dots, 0]^T$ . Consequently, the result is the union of three stars:  $\mathbf{act}_j^S(\theta) = \langle \mathbf{M}_{sc} \mathbf{c} + \mathbf{s}_{sc}, \mathbf{M}_{sc} \mathbf{V}, P_1 \rangle \cup \langle \mathbf{M} \mathbf{c}, \mathbf{M} \mathbf{V}, P_2 \rangle \cup \langle \mathbf{M} \mathbf{c} + \mathbf{s}_{one}, \mathbf{M} \mathbf{V}, P_3 \rangle$ .

Again, when intersecting the star  $\theta$  with  $V_{\min} < x_j < V_{\max}$ ,  $x_j \leq V_{\min}$  respectively  $x_j \geq V_{\max}$ , certain resulting subsets may become empty (see 2.4) and thus their further processing can be omitted.



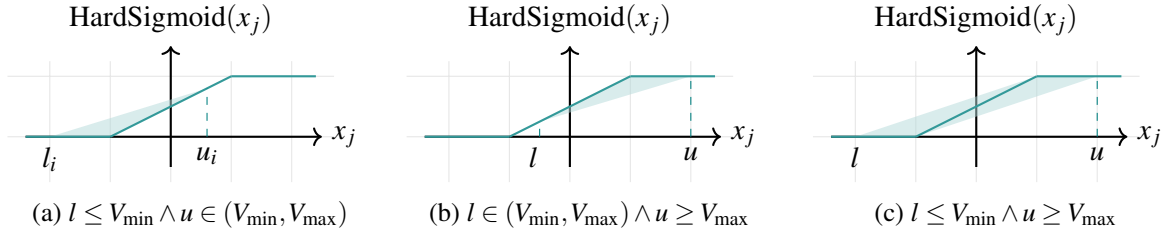


Figure 6: Relaxation for the hard tanh function. The dark line shows the exact set (non-convex) and the light area the approximate set (convex and linear).

### Over-approximate Analysis

Using the over-approximate analysis of hard sigmoid, we consider cases where a convex triangle or trapezoid is applicable based on the input. The  $\overline{\text{act}}_j^S(\theta)$  operation introduces a new variable  $\alpha_{m+1}$  regardless of which case occurs.

If the lower bound is less than  $V_{\min}$  and the upper bound is between  $V_{\min}$  and  $V_{\max}$ , then three new constraints are introduced: (1)  $\alpha_{m+1} \geq 0$ , (2)  $\alpha_{m+1} \geq \frac{1}{V_{\max}-V_{\min}} \cdot x_j + \frac{V_{\min}}{V_{\max}-V_{\min}}$ , and (3)  $\alpha_{m+1} \leq \frac{u \cdot (x_j - l)}{u - l}$ . In the dual scenario when the lower bound is between  $V_{\min}$  and  $V_{\max}$  while the upper bound exceeds  $V_{\max}$ , we encounter the constraints: (1)  $\alpha_{m+1} \leq 1$ , (2)  $\alpha_{m+1} \leq \frac{1}{V_{\max}-V_{\min}} \cdot x_j - \frac{V_{\min}}{V_{\max}-V_{\min}}$ , and (3)  $\alpha_{m+1} \geq \frac{l-1}{l-u} \cdot x_j + \frac{l \cdot (1-u)}{l-u}$ . Lastly, when in dimension  $j$  the star is between  $V_{\min}$  and  $V_{\max}$ , then we introduce four constraints: (1)  $\alpha_{m+1} \leq 1$ , (2)  $\alpha_{m+1} \geq 0$ , (3)  $\alpha_{m+1} \leq \frac{1}{V_{\max}-l} \cdot x_j - \frac{l}{V_{\max}-l}$ , and (4)  $\alpha_{m+1} \geq \frac{1}{V_{\min}-u} \cdot x_j - \frac{V_{\min}}{u-V_{\min}}$ . It is important to highlight that when the domain of variable  $x_j$  is between  $V_{\min}$  and  $V_{\max}$ , less than  $V_{\min}$  (i.e.,  $u \leq V_{\min}$ ) or greater than  $V_{\max}$  (i.e.,  $l \geq V_{\max}$ ), the resulting stars remain the same as described in the exact approach.

Furthermore, when dealing with an unbounded input set  $\theta$  we distinguish three cases, as mentioned earlier. These cases are as follows: (1)  $x_j \in (-\infty, u)$ , (2)  $x_j \in (l, \infty)$ , and (3)  $x_j \in (-\infty, \infty)$ . The cases (1) and (2) are again divided into two sub-cases, hence we obtain five different cases, each one presented in Table 2, coupled with the corresponding constraints and illustrations.

### 3.5 Unit Step Function Layer

The unit step activation function (also called the heaviside function) is widely used in neural networks. In this work, we generalize the unit step function, by introducing three parameters with commonly used values  $val = 0$ ,  $R_{\min} = 0$ , and  $R_{\max} = 1$ .

**Definition 3.4** (Unit Step [15]). The *unit step* function with *separator*  $val \in \mathbb{R}$ , *lower limit*  $R_{\min} \in \mathbb{R}$  and *upper limit*  $R_{\max} \in \mathbb{R}_{\geq R_{\min}}$  is defined for each  $x \in \mathbb{R}$  by

$$\text{UnitStep}(x) = \begin{cases} R_{\min} & \text{if } x < val \\ R_{\max} & \text{if } x \geq val \end{cases} \quad (6)$$

### Exact Analysis

The result  $\text{act}_j^U(\theta)$  of applying unit step on a star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  in dimension  $j$  is obtained as follows. First,  $\theta$  is decomposed into two parts  $\theta_1$  and  $\theta_2$  that result from the intersection of  $\theta$  with  $x_j < val$  resp.  $x_j \geq val$ . Then, the values in the  $j$ th dimension are set to  $R_{\min}$  and  $R_{\max}$ , respectively in the stars  $\theta_1$  and  $\theta_2$ .

Domain of $x_j$	Introduced constraints	Graphical illustration
$l = -\infty \wedge u \in [V_{\min}, V_{\max}]$	$\alpha_{m+1} \geq V_{\min}$ $\alpha_{m+1} \geq x_j$ $\alpha_{m+1} \leq u_j$	
$l = -\infty \wedge u > V_{\max}$	$\alpha_{m+1} \geq V_{\min}$ $\alpha_{m+1} \leq V_{\max}$ $\alpha_{m+1} \geq \frac{V_{\min} - V_{\max}}{V_{\min} - u_j} \cdot x_j - \frac{V_{\min} \cdot (V_{\max} - u_j)}{V_{\min} - u_j}$	
$l \in [V_{\min}, V_{\max}] \wedge u = \infty$	$\alpha_{m+1} \leq V_{\max}$ $\alpha_{m+1} \leq x_j$ $\alpha_{m+1} \geq l_j$	
$l < V_{\min} \wedge u = \infty$	$\alpha_{m+1} \leq V_{\max}$ $\alpha_{m+1} \geq V_{\min}$ $\alpha_{m+1} \leq \frac{V_{\max} - V_{\min}}{V_{\max} - l_j} \cdot x_j - \frac{V_{\max} \cdot (l_j - V_{\min})}{V_{\max} - l_j}$	
$l = -\infty \wedge u = \infty$	$\alpha_{m+1} \geq V_{\min}$ $\alpha_{m+1} \leq V_{\max}$	

Table 2: Approximation rules for hard sigmoid, when the input is unbounded. In total, we distinguish five cases, for each we show the case itself, the introduced constraints and a graphical illustration.

We achieve this by using the projection matrix  $\mathbf{M} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{j-1}, \mathbf{0}, \mathbf{e}_{j+1}, \dots, \mathbf{e}_n]$  and translation vectors  $\mathbf{s}_{\min} = [0, \dots, R_{\min}, \dots, 0]^\top$  and  $\mathbf{s}_{\max} = [0, \dots, R_{\min}, \dots, 0]^\top$ . The resulting stars are  $\langle \mathbf{M}\mathbf{c} + \mathbf{s}_{\min}, \mathbf{M}\mathbf{V}, P_1 \rangle$  and  $\langle \mathbf{M}\mathbf{c} + \mathbf{s}_{\max}, \mathbf{M}\mathbf{V}, P_2 \rangle$ . Note that if the domain  $(l, u)$  of  $x_j$  does not contain the value  $val$ , then the case splitting is not necessary and only one of the stars is the final result, correspondingly to the non-empty intersection with one of the halfspaces.

### Over-approximate Analysis

The over-approximate computation of the unit step function uses a linear, convex trapezoid as shown in Figure 7, which is again the tightest over-approximation that we can achieve. The  $\overline{\text{act}}_j^U(\theta)$  operation also introduces a new variable  $\alpha_{m+1}$  and, in this case, four new constraints, which define the trapezoid. The four constraints are as follows: (1)  $\alpha_{m+1} \geq R_{\min}$ , (2)  $\alpha_{m+1} \leq R_{\max}$ , (3)  $\alpha_{m+1} \leq \frac{R_{\max} - R_{\min}}{val - l} \cdot x_j + \frac{val \cdot R_{\min} - l \cdot R_{\max}}{val - l}$ , and (4)  $\alpha_{m+1} \geq \frac{R_{\max} - R_{\min}}{u - val} \cdot x_j + \frac{u \cdot R_{\min} - val \cdot R_{\max}}{u - val}$ . As previously, the result of  $\overline{\text{act}}_j^U(\theta)$  is a single star which over-approximates the exact resulting star(s). In case the domain of  $\theta$  in dimension  $j$  lies completely in either  $(-\infty, val]$  or  $[val, \infty)$ , then the resulting star is either  $\theta_1 = \langle \mathbf{s}_{\min} + \mathbf{M}\mathbf{c}, \mathbf{M}\mathbf{V}, P \rangle$

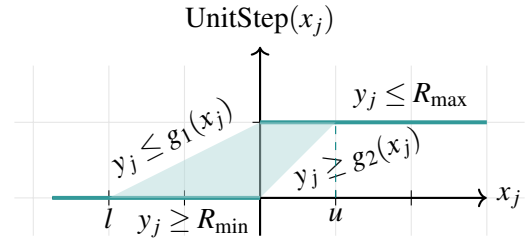


Figure 7: Relaxation for the unit step function. The dark line shows the exact set and the light area the approximate set. The constraints  $y_j \leq g_1(x_j)$  and  $y_j \geq g_2(x_j)$  correspond to relaxations (3) and (4).

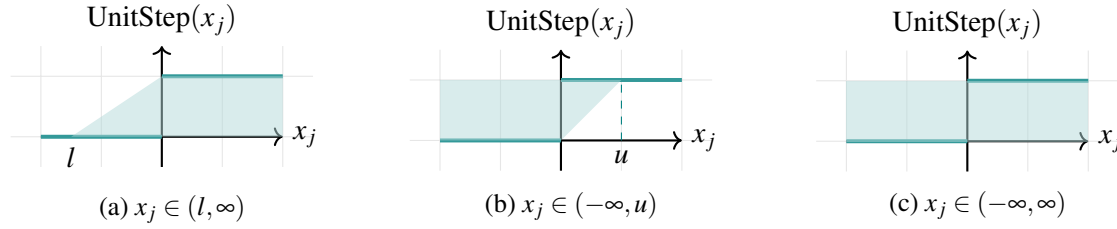


Figure 8: Convex relaxations of the unit stepfunction with three cases of an unbounded input set.

or  $\theta_2 = \langle \mathbf{s}_{\max} + \mathbf{M}\mathbf{c}, \mathbf{M}\mathbf{V}, P \rangle$ , respectively.

Finally, in case there is an unbounded input star  $\theta$ , again three cases are distinguished, as in case of LeakyReLU. The three cases are as follows: (1)  $x_j \in (-\infty, u)$ , (2)  $x_j \in (l, \infty)$ , and (3)  $x_j \in (-\infty, \infty)$ . The analysis for unbounded input is the same; the only aspect that changes is the introduced constraints. See the corresponding constraints for each case visualized in Figure 8.

## 4 Experimental Evaluation

We implemented our proposed algorithms using the open-source C++ tool HyPro [46] and evaluated them on four different benchmark families. The ACAS Xu and drone hovering benchmarks contain only ReLU activations while the thermostat and sonar classifier benchmarks use the unit step and hard sigmoid activation functions besides ReLU. Both the exact and the over-approximation approaches are evaluated. The evaluations were performed on RWTH Aachen University's HPC Cluster [58] using Rocky Linux 8 as the operating system. Each execution ran on an individual node equipped with 16GB RAM and two Intel Xeon Platinum 8160 "SkyLake" processors with a total of 16 cores. A 48-hour timeout was set for each experiment.

### 4.1 ACAS Xu

The Airborne Collision Avoidance System Xu (ACAS Xu) is a mid-air collision avoidance system focusing on unmanned aircrafts. The ACAS Xu networks (ACAS Xu DNNs) provide advisories for horizontal maneuvers to avoid collisions while minimizing unnecessary alerts. The ACAS Xu benchmark consists of a set of 45 feedforward neural networks, each with seven fully connected layers, comprising a combined count of 300 neurons. Each network possesses five inputs (see Figure 9) and five outputs. For further information about the ACAS Xu benchmark see [26, 27].

For our evaluation, we first compute the reachable set of the networks. Afterward, we check whether the reachable set is fully included in the safe zone. If yes then the FNN is safe, otherwise we can conclude unsafety only for the exact analysis. We check the *safety verification time* (VT) in seconds, using the ten safety properties  $\phi_1, \phi_2, \dots, \phi_{10}$  from [27].

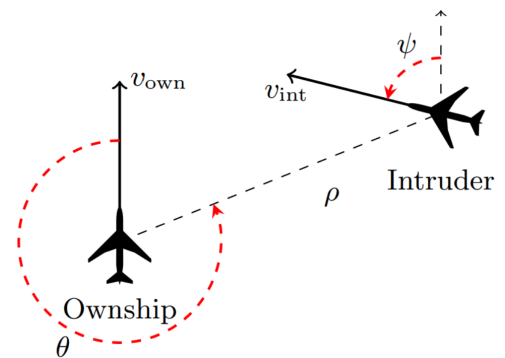


Figure 9: Vertical view of the inputs of ACAS Xu networks. [27]

According to the condensed results, which are shown in Table 3, we can conclude that the star set approach is able to correctly verify the safety properties. We marked with bold-face numbers, where the given property could be verified on all the relevant networks. In case of exact analysis of  $\phi_2$ , the verification results were correct, but in case of 3 networks, timeout occurred. The over-approximate analysis could verify correctly only a subset of the networks. We refer to the Appendix A.2 of this paper for the detailed results, where we show the reachability result and safety verification times of each property and network combinations. Regarding the running time of the reachability analysis: a meaningful comparison could have been made with the implementation provided in [57]; however, it is in Matlab and currently, we do not own a Matlab license.

prop.	Exact	Overapprox.
	AVG VT(s)	AVG VT(s)
$\phi_1$	<b>35244.9</b>	2293.4
$\phi_2$	44715.5	2316.2
$\phi_3$	<b>279.4</b>	12.4
$\phi_4$	<b>98.0</b>	11.4

Table 3: Average Verification results for properties  $\phi_1, \phi_2, \phi_3, \phi_4$  in seconds.

## 4.2 Drone Hovering

$AC_{x,y}$	Exact			Overapprox.		
	RT(s)	RES	CT(s)	RT(s)	RES	CT(s)
$AC1_1$	61.4	True	4.9	0.2	False	0.0
$AC1_2$	0.5	True	0.0	0.1	False	0.0
$AC2_1$	462.4	True	17.7	0.5	False	0.0
$AC2_2$	0.1	True	0.0	0.1	False	0.0
$AC3_1$	-	-	-	2.9	False	0.4
$AC3_2$	5.1	True	0.1	0.2	False	0.0
$AC4_1$	-	-	-	8.7	False	1.5
$AC4_2$	103.0	True	5.5	0.7	False	0.0
$AC5_1$	304.8	True	26.1	0.4	False	0.1
$AC5_2$	0.1	False	0.0	0.1	False	0.0
$AC6_1$	2631.7	True	84.4	0.7	False	0.1
$AC6_2$	0.1	False	0.0	0.1	False	0.0
$AC7_1$	-	-	-	2.5	False	0.1
$AC7_2$	4.1	True	0.1	0.3	False	0.0
$AC8_1$	-	-	-	65.9	False	19.8
$AC8_2$	0.8	False	0.0	0.6	False	0.0

Table 4: Evaluation results of the drones benchmark. The network is identified as  $AC_x$ , the lower-right index  $y$  shows the tested property.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result. True indicates that the given neural network was verified to be safe, with respect to the property. Conversely, False means that the network could not be verified as safe (either because of over-approximation error or due to the network being inherently unsafe). Cells with (-) indicate cases where timeout occurs.

Autonomous drone control revolves around launching a drone into the air and enabling it to hover at a desired altitude [18, 17]. This benchmark consists of eight neural networks. The first four consist of two, and the other four networks consist of three hidden layers, each followed by a ReLU activation function. For further info about the benchmark we refer to [38]. We compute the reachability set of the networks as well as the safety verification using our algorithm and measure the reachable set computation time and safety checking time in seconds. The networks are verified both with the exact and the over-approximation method. For each neural network we test two properties. The presented results in Table 4 show, as we would expect, that the over-approximative method is much faster compared to the exact algorithm. However, the exact method verifies almost every property while the over-approximate approach fails in all cases (though some were inherently unsafe). This confirms that the over-approximate

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
Set 1	4359	False	783	True	263	True	102	True
Set 2	206243	False	1284	True	245	True	100	True
Set 3	33945	True	3768	True	401	True	308	True
Set 4	7974	True	359	True	103	True	102	True

Table 5: Local adversarial robustness tests of the exact approach. *RT* is the reachable set computation time in milliseconds. *RES* is the safety verification result. True indicates that the neural network correctly classifies the input set, while False means that the network was unable to correctly classify the input set.

analysis is more scalable and has a smaller computational cost; however, it sacrifices the completeness of the method.

### 4.3 Thermostat Controller

This benchmark mentioned in the Master’s thesis [25] maintains the room temperature  $x$  between  $17^\circ\text{C}$  and  $23^\circ\text{C}$  using a thermostat. It achieves this by activating (mode on) and deactivating (mode off) the heater based on the sensed temperature. The neural network representing the thermostat’s controller is a feedforward neural network with four layers. The input consists of two neurons that express the temperature  $x \in \mathbb{R}$  and the current mode (on or off) as  $m \in \{0, 1\}$ . Furthermore, two hidden layers follow, each with ten neurons. Lastly, using the unit step activation function, the output layer predicts whether the heater should turn on or off, producing the control input  $Kh = 15$  or  $Kh = 0$ , respectively. We compute the reachable sets to verify the safety of the described NN using our reachability method.

We tested one safety property of the thermostat controller, the input temperature being between  $22^\circ$  and  $23^\circ$ , and the thermostat being turned on, i.e.,  $m = 1$ , the expected control output should be the turn off signal. However, the reachability analysis shows two resulting star sets representing the value 15, meaning that the neural network violates its safety specification. Therefore, we take those star sets and construct the complete counter input set to falsify the neural network, i.e., prove that it is unsafe. The construction of the complete counter input set works as explained in Theorem 2 of [54].

### 4.4 Sonar Binary Classifier

In this section, we evaluate the robustness of a neural network used for binary classification of a sonar dataset. This dataset describes sonar chirp returns bouncing off from different objects [5]. It contains 60 input variables representing the returned beams’ strength at different angles. The verified neural network should be capable of robust binary classification, distinguishing between rocks and metal cylinders. The neural network consists of one hidden layer with 60 neurons, followed by a ReLU activation and an output layer with a single neuron, followed by the composition of a hard sigmoid and a unit step activation function. The property we want to verify is the local robustness of the neural network. A neural network is  $\delta$ -locally-robust at input  $x$ , if for every  $x'$  such that  $\|x - x'\|_\infty \leq \delta$ , the network assigns the same output label to  $x$  and  $x'$ . Our focus lies in determining the robustness threshold that our verification method can provide for the network (i.e., finding the largest  $\delta$  for which the robustness property still holds).

We examine this problem on four inputs of the dataset and four  $\delta$  values. The first two inputs should output 1, which means a rock, and the next two 0, which means a metal cylinder. The True

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
Set 1	234	Inconclusive	205	True	163	True	103	True
Set 2	396	Inconclusive	279	True	157	True	103	True
Set 3	407	Inconclusive	367	True	177	True	174	True
Set 4	339	True	167	True	104	True	101	True

Table 6: Local adversarial robustness tests of the over-approximate approach. *RT* is the reachable set computation time in milliseconds. *RES* is the safety verification result. True indicates that the neural network correctly classifies the input set, while False means that the network was unable to correctly classify the input set. Additionally, Inconclusive is assigned when the reachability analysis algorithm cannot provide a conclusive answer due to the over-approximation errors.

results indicate correct classifications within the robustness threshold ( $\forall x'$  being correctly classified), False denotes incorrect predictions. Moreover, in the case of over-approximate analysis, Inconclusive means that the verification result is ambiguous due to over-approximation error. A comparison between the exact and over-approximative algorithms reveals that the exact algorithm proves network robustness in more cases. Furthermore, different input sets (meaning a single input and its  $\delta$  neighborhood) exhibit varying local robustness. For example, in Table 5, for Set 2, the optimal  $\delta$  value is between 0.01 and 0.001. Tables 5 and 6 are condensed versions of our experiments, to see the complete results, please check the Appendix A.3 of this paper.

## 5 Conclusion

In this paper, we proposed algorithms for star-based reachability analysis of various activation functions used in feed-forward neural networks. To ensure generality, we implemented the activation functions with flexibility for adaptation to specific use cases. We implemented an NNET parser in Hypro to simplify the incorporation of additional benchmarks. The presented evaluation results offer valuable insights into network behavior and safety.

As future work, we plan to integrate further layer types. Consequently, we are planning to integrate a more widely-used standard such as ONNX, for storing and parsing neural network inputs. Moreover, comprehensive experiments and evaluations will offer deeper insights into the performance, accuracy, and limitations of this analysis method when applied to neural networks with other activation functions and layer types, hence, exploring its effectiveness on a more realistic and diverse scale of benchmarks.

We also plan to integrate backpropagation methods using star-sets. Backpropagation is a widely used algorithm for training artificial neural networks, offering numerous advantages in efficient training, scalability, flexibility, and generalization capabilities [62]. Investigating the compatibility and benefits of incorporating backpropagation with star sets can significantly contribute to the advancement of safe and reliable neural networks.

Finally, we are planning to adapt abstraction refinement techniques (such as CEGAR), to reduce the over-approximation error during the reachable set analysis.

**Acknowledgements.** We are grateful to Dario Guidotti, Stefano Demarchi, and Armando Tacchella for generously sharing with us their drone hovering benchmark. This project has received funding from the European Union’s Horizon 2020 programme under the Skłodowska-Curie grant agreement No. 956200.

## References

- [1] Oludare Isaac Abiodun, Aman Jantan, Abiodun Esther Omolara, Kemi Victoria Dada, Abubakar Malah Umar, Okafor Uchenwa Linus, Humaira Arshad, Abdullahi Aminu Kazaure, Usman Gana & Muhammad Ubale Kiru (2019): *Comprehensive review of artificial neural network applications to pattern recognition*. *IEEE Access* 7, pp. 158820–158846, doi:10.1109/ACCESS.2019.2945545.
- [2] Sushma Priya Anthadupula & Manasi Gyanchandani (2021): *A Review and Performance Analysis of Non-Linear Activation Functions in Deep Neural Networks*. *Int. Res. J. Mod. Eng. Technol. Sci*, doi:10.1109/iscid.2009.214.
- [3] Stanley Bak & Parasara Sridhar Duggirala (2017): *Simulation-Equivalent Reachability of Large Linear Systems with Inputs*. In Rupak Majumdar & Viktor Kunčák, editors: *Computer Aided Verification*, pp. 401–420, doi:10.1007/978-3-319-63387-9\_20.
- [4] Akhilan Boopathy, Tsui-Wei Weng, Pin-Yu Chen, Sijia Liu & Luca Daniel (2019): *CNN-Cert: An Efficient Framework for Certifying Robustness of Convolutional Neural Networks*. *Proceedings of the AAAI Conference on Artificial Intelligence* 33(01), pp. 3240–3247, doi:10.1609/aaai.v33i01.33013240. Available at <https://ojs.aaai.org/index.php/AAAI/article/view/4193>.
- [5] Jason Brownlee (2022): *Binary Classification Tutorial with the Keras Deep Learning Library*. <https://machinelearningmastery.com/binary-classification-tutorial-with-the-keras-deep-learning-library/>. [Accessed : June 1, 2023].
- [6] Chih-Hong Cheng, Georg Nührenberg & Harald Ruess (2017): *Maximum resilience of artificial neural networks*. In: *Automated Technology for Verification and Analysis: 15th International Symposium, ATVA 2017, Pune, India, October 3–6, 2017, Proceedings 15*, Springer, pp. 251–268, doi:10.1007/978-3-319-68167-2\_18.
- [7] Edmund M Clarke & Jeannette M Wing (1996): *Formal methods: State of the art and future directions*. *ACM Computing Surveys (CSUR)* 28(4), pp. 626–643, doi:10.1145/242223.242257.
- [8] Ronan Collobert (2004): *Large scale machine learning*. Technical Report, Université de Paris VI.
- [9] Ekin Cubuk, Barret Zoph, Samuel Schoenholz & Quoc Le (2017): *Intriguing Properties of Adversarial Examples*.
- [10] Leonid Datta (2020): *A Survey on Activation Functions and their relation with Xavier and He Normal Initialization*.
- [11] Educative (2023): *What is the vanishing gradient problem?* <https://www.educative.io/answers/what-is-the-vanishing-gradient-problem>. [Accessed: May 12, 2023].
- [12] Ruediger Ehlers (2017): *Formal verification of piece-wise linear feed-forward neural networks*. In: *Automated Technology for Verification and Analysis: 15th International Symposium, ATVA 2017, Pune, India, October 3–6, 2017, Proceedings 15*, Springer, pp. 269–286, doi:10.1007/978-3-319-68167-2\_19.
- [13] Dumitru Erhan, Christian Szegedy, Alexander Toshev & Dragomir Anguelov (2014): *Scalable Object Detection Using Deep Neural Networks*. In: *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2155–2162, doi:10.1109/CVPR.2014.276.
- [14] Aymeric Fromherz, Klas Leino, Matt Fredrikson, Bryan Parno & Corina Pasareanu (2021): *Fast Geometric Projections for Local Robustness Certification*. In: *International Conference on Learning Representations*. Available at <https://openreview.net/forum?id=zWyluxjDdZJ>.
- [15] Osvaldo Gervasi, Beniamino Murgante, Antonio Laganà, David Taniar, Youngsong Mun & Marina L. Gavrilova, editors (2008): *Computational Science and Its Applications - ICCSA 2008, International Conference, Perugia, Italy, June 30 - July 3, 2008, Proceedings, Part I*. *Lecture Notes in Computer Science* 5072, Springer, doi:10.1007/978-3-540-69839-5.
- [16] Ian J. Goodfellow, Jonathon Shlens & Christian Szegedy (2014): *Explaining and Harnessing Adversarial Examples*. *CoRR* abs/1412.6572. Available at <https://api.semanticscholar.org/CorpusID:6706414>.

- [17] Dario Guidotti (2022): *Verification of Neural Networks for Safety and Security-critical Domains*. ISSN 1613-0073 CEUR Workshop Proceedings. Available at [https://ceur-ws.org/Vol-3345/paper10\\_RiCeRCa3.pdf](https://ceur-ws.org/Vol-3345/paper10_RiCeRCa3.pdf).
- [18] Dario Guidotti, Stefano Demarchi, Luca Pulina & Armando Tacchella (2022): *Evaluating Reachability Algorithms for Neural Networks on NeVer2*.
- [19] Patrick Henriksen & Alessio Lomuscio (2021): *DEEPSPLIT: An Efficient Splitting Method for Neural Network Verification via Indirect Effect Analysis*. In: *IJCAI*, pp. 2549–2555, doi:10.24963/ijcai.2021/351.
- [20] Michael Hinchey, Jonathan Bowen & Christopher Rouff (2006): *Introduction to Formal Methods*. Springer, doi:10.1007/1-84628-271-3\_2.
- [21] Geoffrey Hinton, Li Deng, Dong Yu, George E. Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N. Sainath & Brian Kingsbury (2012): *Deep Neural Networks for Acoustic Modeling in Speech Recognition: The Shared Views of Four Research Groups*. *IEEE Signal Processing Magazine* 29(6), pp. 82–97, doi:10.1109/MSP.2012.2205597.
- [22] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath et al. (2012): *Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups*. *IEEE Signal processing magazine* 29(6), pp. 82–97, doi:10.1109/MSP.2012.2205597.
- [23] Chao Huang, Jiameng Fan, Wenchao Li, Xin Chen & Qi Zhu (2019): *Reachnn: Reachability analysis of neural-network controlled systems*. *ACM Transactions on Embedded Computing Systems (TECS)* 18(5s), pp. 1–22, doi:10.1145/3358228.
- [24] Xiaowei Huang, Marta Kwiatkowska, Sen Wang & Min Wu (2017): *Safety verification of deep neural networks*. In: *International conference on computer aided verification*, Springer, pp. 3–29, doi:10.1007/978-3-319-63387-9\_1.
- [25] Ruoran Gabriela Jiang (2023): *Verifying ai-controlled hybrid systems*. Master’s thesis, RWTH Aachen University, Aachen, Germany. Available at [https://ths.rwth-aachen.de/wp-content/uploads/sites/4/master\\_thesis\\_jiang.pdf](https://ths.rwth-aachen.de/wp-content/uploads/sites/4/master_thesis_jiang.pdf).
- [26] Kyle D. Julian, Mykel J. Kochenderfer & Michael P. Owen (2019): *Deep Neural Network Compression for Aircraft Collision Avoidance Systems*. *Journal of Guidance, Control, and Dynamics* 42(3), pp. 598–608, doi:10.2514/1.g003724. Available at <https://doi.org/10.2514%2F1.g003724>.
- [27] Guy Katz, Clark Barrett, David L. Dill, Kyle Julian & Mykel J. Kochenderfer (2017): *Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks*. In Rupak Majumdar & Viktor Kunčák, editors: *Computer Aided Verification*, Springer International Publishing, Cham, pp. 97–117, doi:10.1007/978-3-319-63387-9\_5.
- [28] Guy Katz, Derek A Huang, Duligur Ibeling, Kyle Julian, Christopher Lazarus, Rachel Lim, Parth Shah, Shantanu Thakoor, Haoze Wu, Aleksandar Zeljić et al. (2019): *The marabou framework for verification and analysis of deep neural networks*. In: *International Conference on Computer Aided Verification*, Springer, pp. 443–452, doi:10.1007/978-3-030-25540-4\_26.
- [29] Philipp Kern, Marko Kleine Büning & Carsten Sinz (2022): *Optimized Symbolic Interval Propagation for Neural Network Verification*. In: *1st Workshop on Formal Verification of Machine Learning (WFVML 2022) colocated with ICML 2022: International Conference on Machine Learning*.
- [30] Kumar, Niranjana (2019): *Deep Learning: Feedforward Neural Networks Explained*. <https://medium.com/hackernoon/deep-learning-feedforward-neural-networks-explained-c34ae3f084f1>. [Accessed: May 03, 2023].
- [31] Sampo Kuutti, Richard Bowden, Yaochu Jin, Phil Barber & Saber Fallah (2021): *A Survey of Deep Learning Applications to Autonomous Vehicle Control*. *IEEE Transactions on Intelligent Transportation Systems* 22(2), pp. 712–733, doi:10.1109/TITS.2019.2962338.
- [32] Yann LeCun, Yoshua Bengio & Geoffrey Hinton (2015): *Deep learning*. *nature* 521(7553), pp. 436–444, doi:10.1038/nature14539.



- [33] Andy Lee (2015): *Comparing deep neural networks and traditional vision algorithms in mobile robotics*. Swarthmore University. Available at <https://api.semanticscholar.org/CorpusID:10011895>.
- [34] Changliu Liu, Tomer Arnon, Christopher Lazarus, Christopher Strong, Clark Barrett & Mykel J. Kochenderfer (2021): *Algorithms for Verifying Deep Neural Networks*. *Foundations and Trends in Optimization* 4(3-4), pp. 244–404, doi:10.1561/24000000035. Available at <http://theory.stanford.edu/~barrett/pubs/LAL+21.pdf>.
- [35] Weibo Liu, Zidong Wang, Xiaohui Liu, Nianyin Zeng, Yurong Liu & Fuad E. Alsaadi (2017): *A survey of deep neural network architectures and their applications*. *Neurocomputing* 234, pp. 11–26, doi:10.1016/j.neucom.2016.12.038. Available at <https://www.sciencedirect.com/science/article/pii/S0925231216315533>.
- [36] Alessio Lomuscio & Lalit Maganti (2017): *An approach to reachability analysis for feed-forward ReLU neural networks*.
- [37] Andrew L. Maas (2013): *Rectifier Nonlinearities Improve Neural Network Acoustic Models*. In: *Proceedings of the International Conference on Machine Learning*, pp. 1–6. Available at <https://www.semanticscholar.org/paper/Rectifier-Nonlinearities-Improve-Neural-Network-Maas/367f2c63a6f6a10b3b64b8729d601e69337ee3cc>.
- [38] Hana Masara (2023): *Star Set-based Reachability Analysis of Neural Networks with Differing Layers and Activation Functions*. Bachelor’s thesis, RWTH Aachen University, Aachen, Germany. Available at <https://ths.rwth-aachen.de/wp-content/uploads/sites/4/Thesis-Hana-Masara.pdf>.
- [39] ONNX. <https://onnx.ai/>. [Accessed: May 30, 2023].
- [40] (2021): *PyTorch: torch.nn.Hardsigmoid*. <https://pytorch.org/docs/stable/generated/torch.nn.Hardsigmoid.html>. Accessed: May 16, 2023.
- [41] Luca Pulina & Armando Tacchella (2010): *An Abstraction-Refinement Approach to Verification of Artificial Neural Networks*. In Tayssir Touili, Byron Cook & Paul Jackson, editors: *Computer Aided Verification*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 243–257, doi:10.1007/978-3-642-14295-6\_24.
- [42] Luthfi Ramadhan (2021): *Neural Network: The Dead Neuron*. <https://towardsdatascience.com/neural-network-the-dead-neuron-eaa92e575748>. [Accessed: May 14, 2023].
- [43] Waseem Rawat & Zenghui Wang (2017): *Deep convolutional neural networks for image classification: A comprehensive review*. *Neural computation* 29(9), pp. 2352–2449, doi:10.1162/neco\_a\_00990.
- [44] Wenjie Ruan, Xiaowei Huang & Marta Kwiatkowska (2018): *Reachability Analysis of Deep Neural Networks with Provable Guarantees*. In: *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, International Joint Conferences on Artificial Intelligence Organization, pp. 2651–2659, doi:10.24963/ijcai.2018/368. Available at <https://doi.org/10.24963/ijcai.2018/368>.
- [45] Wojciech Samek, Grégoire Montavon, Sebastian Lapuschkin, Christopher J Anders & Klaus-Robert Müller (2021): *Explaining deep neural networks and beyond: A review of methods and applications*. *Proceedings of the IEEE* 109(3), pp. 247–278, doi:10.1109/JPROC.2021.3060483.
- [46] Stefan Schupp, Erika Ábrahám, Ibtissem Makhoul & Stefan Kowalewski (2017): *HyPro: A C++ Library of State Set Representations for Hybrid Systems Reachability Analysis*. In Clark Barrett, Misty Davies & Temesghen Kahsay, editors: *NASA Formal Methods*, pp. 288–294, doi:10.1007/978-3-319-57288-8\_20.
- [47] Stefan Schupp, Erika Ábrahám, Ibtissem Ben Makhoul & Stefan Kowalewski (2017): *HyPro: A C++ library of state set representations for hybrid systems reachability analysis*. In: *NASA Formal Methods Symposium*, Springer, pp. 288–294, doi:10.1007/978-3-319-57288-8\_20.
- [48] Stefan Schupp, Goran Frehse & Erika Ábrahám (2019): *State set representations and their usage in the reachability analysis of hybrid systems*. Ph.D. thesis, RWTH Aachen University, doi:10.18154/RWTH-2019-08875. Available at <https://publications.rwth-aachen.de/record/767529/files/767529.pdf>.
- [49] Gagandeep Singh, Timon Gehr, Markus Püschel & Martin T. Vechev (2019): *An abstract domain for certifying neural networks*. *Proc. ACM Program. Lang.* 3(POPL), pp. 41:1–41:30, doi:10.1145/3290354. Available at <https://doi.org/10.1145/3290354>.

- [50] Bing Sun, Jun Sun, Ting Dai & Lijun Zhang (2021): *Probabilistic Verification of Neural Networks Against Group Fairness*. In: *Formal Methods: 24th International Symposium, FM 2021, Virtual Event, November 20-26, 2021, Proceedings*, Springer-Verlag, Berlin, Heidelberg, pp. 83–102, doi:10.1007/978-3-030-90870-6\_5. Available at [https://doi.org/10.1007/978-3-030-90870-6\\_5](https://doi.org/10.1007/978-3-030-90870-6_5).
- [51] Daniel Svozil, Vladimir Kvasnicka & Jiri Pospichal (1997): *Introduction to multi-layer feed-forward neural networks*. *Chemometrics and Intelligent Laboratory Systems* 39(1), pp. 43–62, doi:10.1016/S0169-7439(97)00061-0.
- [52] (2023): *TensorFlow Documentation*. [https://www.tensorflow.org/api\\_docs/python/tf/keras/activations/hard\\_sigmoid](https://www.tensorflow.org/api_docs/python/tf/keras/activations/hard_sigmoid). Accessed: May 16, 2023.
- [53] Dung Tran (2020): *Verification of Learning-enabled Cyber-Physical Systems*. Ph.D. thesis, Vanderbilt University Graduate School. Available at <http://hdl.handle.net/1803/15957>.
- [54] Hoang-Dung Tran, Diago Manzananas Lopez, Patrick Musau, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang & Taylor T. Johnson (2019): *Star-Based Reachability Analysis of Deep Neural Networks*. In Maurice H. ter Beek, Annabelle McIver & José N. Oliveira, editors: *Formal Methods – The Next 30 Years*, Springer International Publishing, Cham, pp. 670–686, doi:10.1007/978-3-030-30942-8\_39.
- [55] Hoang-Dung Tran, Patrick Musau, Diego Manzananas Lopez, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang & Taylor T Johnson (2019): *Parallelizable reachability analysis algorithms for feed-forward neural networks*. In: *2019 IEEE/ACM 7th International Conference on Formal Methods in Software Engineering (FormaliSE)*, IEEE, pp. 51–60, doi:10.1109/FormaliSE.2019.00012.
- [56] Hoang-Dung Tran, Neelanjana Pal, Diego Manzananas Lopez, Patrick Musau, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang, Stanley Bak & Taylor T Johnson (2021): *Verification of piecewise deep neural networks: a star set approach with zonotope pre-filter*. *Formal Aspects of Computing* 33, pp. 519–545, doi:10.1007/s00165-021-00553-4.
- [57] Hoang-Dung Tran, Neelanjana Pal, Diego Manzananas Lopez, Patrick Musau, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang, Stanley Bak & Taylor T. Johnson (2021): *Verification of piecewise deep neural networks: a star set approach with zonotope pre-filter*. *Formal Aspects of Computing* 33(4), pp. 519–545, doi:10.1007/s00165-021-00553-4. Available at <https://doi.org/10.1007/s00165-021-00553-4>.
- [58] RWTH University (2023): *RWTH High Performance Computing (Linux)*. <https://help.itc.rwth-aachen.de/service/rhr4fjjuutt/>. [Accessed: July 24, 2023].
- [59] Shiqi Wang, Kexin Pei, Justin Whitehouse, Junfeng Yang & Suman Jana (2018): *Efficient Formal Safety Analysis of Neural Networks*. In: *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS’18*, Curran Associates Inc., Red Hook, NY, USA, pp. 6369–6379.
- [60] Jeannette M Wing (1990): *A specifier’s introduction to formal methods*. *Computer* 23(9), pp. 8–22, doi:10.1109/2.58215.
- [61] Jim Woodcock, Peter Gorm Larsen, Juan Bicarregui & John Fitzgerald (2009): *Formal methods: Practice and experience*. *ACM computing surveys (CSUR)* 41(4), pp. 1–36, doi:10.1145/1592434.1592436.
- [62] Logan G. Wright, Tatsuhiko Onodera, Martin M. Stein, Tianyu Wang, Darren T. Schachter, Zoey Hu & Peter L. McMahon (2022): *Deep physical neural networks trained with backpropagation*. *Nature* 601(7894), pp. 549–555, doi:10.1038/s41586-021-04223-6. Available at <https://doi.org/10.1038/s41586-021-04223-6>.
- [63] Haoze Wu, Aleksandar Zeljić, Guy Katz & Clark Barrett (2022): *Efficient Neural Network Analysis with Sum-of-Infeasibilities*. In Dana Fisman & Grigore Rosu, editors: *Tools and Algorithms for the Construction and Analysis of Systems*, Springer International Publishing, Cham, pp. 143–163, doi:10.1007/978-3-030-99524-9\_8.
- [64] Weiming Xiang, Hoang-Dung Tran & Taylor T. Johnson (2018): *Output Reachable Set Estimation and Verification for Multilayer Neural Networks*. *IEEE Transactions on Neural Networks and Learning Systems* 29(11), pp. 5777–5783, doi:10.1109/TNNLS.2018.2808470.

- [65] Jin Xu, Zishan Li, Bowen Du, Miaomiao Zhang & Jing Liu (2020): *Reluplex made more practical: Leaky ReLU*. In: *2020 IEEE Symposium on Computers and Communications (ISCC)*, pp. 1–7, doi:10.1109/ISCC50000.2020.9219587.

## A Supplementary Material

### A.1 Formal proofs

**Proposition A.1** (Convex polytopes as stars). For any  $m, p \in \mathbb{N}$ ,  $\mathbf{C} \in \mathbb{R}^{p \times m}$  and  $\mathbf{d} \in \mathbb{R}^p$ , the convex polyhedron  $\mathcal{P} = \{\mathbf{x} \in \mathbb{R}^m \mid \mathbf{C}\mathbf{x} \leq \mathbf{d}\}$  can be represented by a star.

*Proof.* It is straightforward to obtain an equivalent starset  $\theta$  of the polytope  $\mathcal{P}$ , using the null vector as center, i.e.,  $\mathbf{c} = [0, 0, \dots, 0]^\top$ , the set of  $n$  unit vectors  $\mathbf{e}_i$  for the basis, i.e.  $\mathbf{V} = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  (i.e., the generator matrix  $\mathbf{V} = \mathbb{I}_n$ ), and the predicate  $P$  in the form of  $\alpha \in P \equiv \mathbf{C}\alpha \leq \mathbf{d}$ .  $\square$

**Proposition A.2** (Affine transformation). Assume an  $(n, m)$ -dimensional star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  and let  $\mathbf{W} \in \mathbb{R}^{k \times n}$  and  $\mathbf{b} \in \mathbb{R}^k$ . Then the affine transformation  $\{\mathbf{W}\mathbf{x} + \mathbf{b} \mid \mathbf{x} \in [\theta]\}$  of  $[\theta]$  is represented by  $\bar{\theta} = \langle \bar{\mathbf{c}}, \bar{\mathbf{V}}, P \rangle$  with  $\bar{\mathbf{c}} = \mathbf{W}\mathbf{c} + \mathbf{b}$  and  $\bar{\mathbf{V}} \in \mathbb{R}^{k \times m}$  with columns  $\mathbf{W}\mathbf{v}^{(1)}, \dots, \mathbf{W}\mathbf{v}^{(m)}$ .

*Proof.* Using the definition of the resulting star set after applying the affine transformation, we have  $\bar{\theta} = \{\mathbf{y} \mid \mathbf{y} = \mathbf{W}(\mathbf{c} + \sum_{j=1}^m (\alpha_j \mathbf{v}^{(j)})) + \mathbf{b} \text{ such that } \alpha \in P\}$ . That means,  $\bar{\theta}$  is another star, having the center  $\bar{\mathbf{c}} = \mathbf{W}\mathbf{c} + \mathbf{b}$  and generator vectors  $\bar{\mathbf{V}} = \{\mathbf{W}\mathbf{v}^{(1)}, \mathbf{W}\mathbf{v}^{(2)}, \dots, \mathbf{W}\mathbf{v}^{(m)}\}$ . Note that the predicate does not change during the computation of the affine mapping of a star.  $\square$

**Proposition A.3** (Intersection with halfspace). Assume an  $(n, m)$ -dimensional star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  and a half-space  $\mathcal{H} = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{h}^T \mathbf{x} \leq g\}$  with some  $\mathbf{h} \in \mathbb{R}^n$  and  $g \in \mathbb{R}$ . Then the intersection  $[\theta] \cap \mathcal{H}$  is represented by the star  $\bar{\theta} = \langle \mathbf{c}, \mathbf{V}, P \cap P' \rangle$  with  $P' = \{\alpha \in \mathbb{R}^m \mid (\mathbf{h}^T \mathbf{V})\alpha \leq g - \mathbf{h}^T \mathbf{c}\}$ .

*Proof.* The resulting star is  $\bar{\theta} = \{\mathbf{x} \mid \mathbf{x} = \mathbf{c} + \sum_{j=1}^m (\alpha_j \mathbf{v}^{(j)}) \text{ s. t. } (\alpha_1, \dots, \alpha_m)^T \in P \wedge \mathbf{h}^T \mathbf{x} \leq g\}$ . Since  $\mathbf{x} = \mathbf{c} + \sum_{j=1}^m (\alpha_j \mathbf{v}^{(j)})$ , the new constraint can be written as  $\mathbf{h}^T (\mathbf{c} + \mathbf{V}\alpha) \leq g$ , where  $\alpha = [\alpha_1, \dots, \alpha_m]^\top$ . Consequently, the new predicate is  $P \cap P'$ ,  $P'(\alpha) = (\mathbf{h}^T \mathbf{V})\alpha \leq g - \mathbf{h}^T \mathbf{c}$ .  $\square$

**Proposition A.4** (Emptiness checking). A star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  is empty if and only if  $P$  is empty.

*Proof.* It is straightforward to see that only the predicate restricts the elements of a star. In other words, if the predicate does not allow any solution (i.e., it's *empty*), then the star set is empty as well.  $\square$

**Proposition A.5** (Bounding box). Assume an  $(n, m)$ -dimensional star  $\theta = \langle \mathbf{c}, \mathbf{V}, P \rangle$  with  $\mathbf{c} = (\mathbf{c}_1, \dots, \mathbf{c}_n)^T$ , and let  $\mathbf{V}_{(i)}$  be the  $i^{\text{th}}$  row of  $\mathbf{V}$ . Let furthermore  $B = \{(x_1, \dots, x_n)^T \in \mathbb{R}^n \mid \bigwedge_{i=1}^n lb_i \leq x_i \leq ub_i\}$  with  $lb_i = \mathbf{c}_i + \min_{\alpha \in P} \mathbf{V}_{(i)}\alpha$  and  $ub_i = \mathbf{c}_i + \max_{\alpha \in P} \mathbf{V}_{(i)}\alpha$  for  $i = 1, \dots, n$ . Then  $[\theta] \subseteq B$ .

*Proof.* According to the star set's definition  $\mathbf{x}_i = \mathbf{c}_i + \sum_{j=1}^m \alpha_j \mathbf{v}_i^{(j)}$ . That is, if we want to find the lower (or upper) bound of  $\mathbf{x}_i$ , we have to find the solution of  $\text{minimize}_{\mathbf{x} \in \theta} \mathbf{x}_i$  (or  $\text{maximize}_{\mathbf{x} \in \theta} \mathbf{x}_i$ , respectively). Using the definition of the star set, we get  $\mathbf{c}_i + \text{minimize}_{\alpha \in P} \mathbf{V}_{(i)}\alpha$  (or  $\mathbf{c}_i + \text{maximize}_{\alpha \in P} \mathbf{V}_{(i)}\alpha$ ).  $\square$



## A.2 ACAS Xu Detailed Results

$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{2,1}$	15837.45	False	325.53	193197	1338.12	False	43.21	1
$N_{2,2}$	36112.55	False	1174.41	472257	1005.57	False	68.81	1
$N_{2,3}$	18532.01	False	303.91	194275	962.63	False	74.02	1
$N_{2,4}$	8254.17	False	299.08	114155	1391.11	False	73.80	1
$N_{2,5}$	56241.64	False	2275.48	677510	1790.96	False	111.03	1
$N_{2,6}$	25991.47	False	650.43	309631	2044.11	False	45.39	1
$N_{2,7}$	65508.57	False	1544.92	679523	3366.64	False	161.01	1
$N_{2,8}$	53195.81	False	1260.82	585647	3677.42	False	135.44	1
$N_{2,9}$	-	-	-	-	1855.61	False	64.66	1
$N_{3,1}$	15559.97	False	727.35	252793	939.06	False	54.55	1
$N_{3,2}$	12911.84	False	295.19	181433	2362.05	False	72.16	1
$N_{3,3}$	24675.76	True	508.77	341669	900.05	False	33.53	1
$N_{3,4}$	7393.05	False	205.80	133782	1500.14	False	95.91	1
$N_{3,5}$	23852.32	False	795.67	365066	1723.58	False	67.01	1
$N_{3,6}$	70608.24	False	2823.21	1003886	2694.87	False	120.82	1
$N_{3,7}$	41256.50	False	1233.86	475299	4753.87	False	150.12	1
$N_{3,8}$	37253.31	False	754.01	472562	1304.76	False	38.45	1
$N_{3,9}$	48309.81	False	1152.44	379221	2498.83	False	54.83	1
$N_{4,1}$	66720.59	False	854.57	402853	973.59	False	74.23	1
$N_{4,2}$	85507.29	True	1130.23	484555	1139.11	False	74.50	1
$N_{4,3}$	10627.68	False	303.07	138170	1096.26	False	38.32	1
$N_{4,4}$	12923.14	False	357.63	143424	1257.89	False	77.31	1
$N_{4,5}$	75982.74	False	1223.96	457447	3312.96	False	93.14	1
$N_{4,6}$	-	-	-	-	1802.32	False	68.50	1
$N_{4,7}$	107331.18	False	2132.42	652417	781.43	False	32.85	1
$N_{4,8}$	67596.46	False	1893.68	515113	2661.60	False	186.54	1
$N_{4,9}$	-	-	-	-	13969.59	False	29.64	1
$N_{5,1}$	30332.56	False	407.09	201773	1291.09	False	60.81	1
$N_{5,2}$	43636.56	False	486.83	261011	1173.13	False	36.83	1
$N_{5,3}$	10740.98	False	191.73	125860	1072.01	False	56.17	1
$N_{5,4}$	6221.00	False	138.55	81364	1415.43	False	109.63	1
$N_{5,5}$	31217.25	False	386.70	213059	1275.20	False	61.52	1
$N_{5,6}$	97829.27	False	1149.06	622084	2196.38	False	94.40	1
$N_{5,7}$	57210.41	False	918.01	355919	2347.15	False	92.50	1
$N_{5,8}$	101123.47	False	1521.47	544813	3443.00	False	106.33	1
$N_{5,9}$	78557.70	False	1164.09	619284	3216.64	False	92.62	1

Table 7: Verification results for property  $P_2$  on 36 ACAS Xu networks. ( $RT$ ) is the reachable set computation time, and ( $CT$ ) is the safety checking time, both in seconds. ( $RES$ ) is the safety verification result. ( $OSS$ ) describes the number of the output star sets. The cells with (-) correspond to networks in which our algorithm was not able to compute the reachability set successfully in less than 48 hours.

$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{1,1}$	3013.64	True	45.71	39835	823.72	False	7.97	1
$N_{1,2}$	3575.72	True	53.74	45648	2214.74	False	15.48	1
$N_{1,3}$	11037.81	True	200.90	114287	3211.44	False	16.37	1
$N_{1,4}$	13111.44	True	267.78	154529	2915.33	False	21.64	1
$N_{1,5}$	9756.54	True	196.13	122297	1618.80	False	9.97	1
$N_{1,6}$	35718.94	True	823.34	376647	1385.90	False	13.06	1
$N_{1,7}$	4712.34	True	85.86	66416	1228.45	False	13.30	1
$N_{1,8}$	8279.50	True	174.76	110139	2226.28	False	38.84	1
$N_{1,9}$	9136.22	True	189.03	135645	2999.83	False	24.42	1
$N_{2,1}$	15355.00	True	325.53	193197	1538.98	False	11.57	1
$N_{2,2}$	34071.21	True	617.56	472257	1147.34	False	19.50	1
$N_{2,3}$	13319.28	True	253.42	194275	956.66	False	18.24	1
$N_{2,4}$	8124.85	True	167.98	114155	1141.23	False	21.41	1
$N_{2,5}$	53191.97	True	1103.15	677510	1489.67	False	22.42	1
$N_{2,6}$	23772.93	True	417.89	309631	1972.41	False	10.89	1
$N_{2,7}$	53504.01	True	1016.92	679523	3330.97	False	39.63	1
$N_{2,8}$	48084.68	True	777.85	585647	4132.09	False	38.59	1
$N_{2,9}$	86837.06	True	1395.51	910575	2225.92	False	19.55	1
$N_{3,1}$	14553.80	True	497.52	252793	1130.78	False	22.86	1
$N_{3,2}$	16570.78	True	359.74	181433	2678.57	False	19.91	1
$N_{3,3}$	28386.63	True	649.15	341669	994.11	False	8.74	1
$N_{3,4}$	8765.59	True	154.49	133782	1662.63	False	23.46	1
$N_{3,5}$	28583.49	True	641.96	365066	1884.56	False	15.86	1
$N_{3,6}$	65843.71	True	1595.55	1003886	2494.56	False	28.20	1
$N_{3,7}$	47664.54	True	1094.01	475299	4453.61	False	36.35	1
$N_{3,8}$	38414.57	True	652.14	472562	1501.02	False	11.65	1
$N_{3,9}$	33949.16	True	746.01	379221	3221.20	False	17.70	1
$N_{4,1}$	35166.18	True	603.51	402853	1007.71	False	19.63	1
$N_{4,2}$	41913.56	True	748.73	484555	1322.22	False	21.65	1
$N_{4,3}$	9966.29	True	164.86	138170	1256.42	False	10.84	1
$N_{4,4}$	12343.79	True	213.54	143424	1295.82	False	20.12	1
$N_{4,5}$	39853.04	True	861.19	457447	3795.18	False	28.86	1
$N_{4,6}$	134265.70	True	2703.07	1296311	1816.04	False	18.24	1
$N_{4,7}$	61325.43	True	942.40	652417	999.71	False	10.53	1
$N_{4,8}$	32988.07	True	775.16	515113	2781.59	False	49.10	1
$N_{4,9}$	87048.69	True	1456.81	984701	14379.08	False	7.58	1
$N_{5,1}$	13215.58	True	262.90	201773	1284.01	False	14.30	1
$N_{5,2}$	17025.91	True	347.06	261011	1025.24	False	7.94	1
$N_{5,3}$	10237.07	True	219.97	125860	1089.35	False	13.62	1
$N_{5,4}$	5989.44	True	106.26	81364	1228.81	False	22.77	1
$N_{5,5}$	14346.32	True	239.07	213059	1361.83	False	17.98	1
$N_{5,6}$	102872.82	True	1120.75	622084	2093.55	False	22.16	1
$N_{5,7}$	31414.81	True	595.66	355919	2211.67	False	21.59	1
$N_{5,8}$	95265.60	True	886.15	544813	3226.16	False	24.57	1
$N_{5,9}$	95694.38	True	1002.30	619284	3546.82	False	24.30	1

Table 8: Verification results for property  $P_1$  on 45 ACAS Xu networks.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result.  $OSS$  describes the number of the output star sets.

$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{1,1}$	1768.51	True	206.94	71930	33.88	False	1.45	1
$N_{1,2}$	1647.32	True	115.58	40273	34.65	False	4.28	1
$N_{1,3}$	442.65	True	29.13	12444	28.54	False	1.87	1
$N_{1,4}$	224.24	True	3.64	4346	12.95	True	0.09	1
$N_{1,5}$	246.37	True	4.09	4820	12.65	True	0.11	1
$N_{1,6}$	65.95	True	0.91	1281	3.79	True	0.03	1
$N_{2,1}$	485.40	True	17.82	16382	23.79	False	1.84	1
$N_{2,2}$	178.63	True	6.48	6924	10.60	False	1.04	1
$N_{2,3}$	316.67	True	10.55	10694	22.85	False	1.21	1
$N_{2,4}$	20.16	True	0.21	351	1.84	True	0.07	1
$N_{2,5}$	111.27	True	2.09	2466	8.17	True	0.11	1
$N_{2,6}$	13.62	True	0.12	255	3.79	True	0.03	1
$N_{2,7}$	60.04	True	0.91	1229	5.07	True	0.15	1
$N_{2,8}$	17.78	True	0.17	329	3.84	True	0.01	1
$N_{2,9}$	9.46	True	0.09	189	0.81	True	0.00	1
$N_{3,1}$	153.28	True	8.77	5999	5.19	True	0.73	1
$N_{3,2}$	2018.12	True	88.51	37541	27.15	False	1.97	1
$N_{3,3}$	390.10	True	12.46	7935	23.19	True	1.89	1
$N_{3,4}$	76.08	True	1.53	2100	29.34	False	2.27	1
$N_{3,5}$	44.62	True	1.19	1042	6.18	True	0.30	1
$N_{3,6}$	99.32	True	1.46	1868	15.38	False	1.73	1
$N_{3,7}$	4.20	True	0.04	107	1.39	True	0.01	1
$N_{3,8}$	33.03	True	0.55	669	5.84	True	0.28	1
$N_{3,9}$	42.36	True	0.82	1223	3.04	True	0.12	1
$N_{4,1}$	50.28	True	1.66	2298	4.80	False	0.90	1
$N_{4,2}$	627.65	True	19.89	18088	16.52	False	1.07	1
$N_{4,3}$	976.11	True	25.53	21237	19.26	False	1.15	1
$N_{4,4}$	34.30	True	0.39	560	2.86	True	0.06	1
$N_{4,5}$	9.23	True	0.23	361	2.41	True	0.03	1
$N_{4,6}$	107.72	True	1.86	2533	39.84	True	0.64	1
$N_{4,7}$	51.25	True	0.61	948	4.47	True	0.08	1
$N_{4,8}$	35.98	True	0.38	576	3.01	True	0.02	1
$N_{4,9}$	36.69	True	0.38	616	7.89	True	0.14	1
$N_{5,1}$	328.52	True	11.35	9556	12.10	False	0.66	1
$N_{5,2}$	61.58	True	2.10	2126	5.45	True	0.88	1
$N_{5,3}$	72.33	True	4.63	2906	10.64	True	3.13	1
$N_{5,4}$	33.32	True	0.86	765	4.98	True	0.09	1
$N_{5,5}$	44.61	True	1.02	1310	7.65	True	0.54	1
$N_{5,6}$	63.87	True	0.71	1166	10.97	True	0.56	1
$N_{5,7}$	4.93	True	0.04	88	0.81	True	0.01	1
$N_{5,8}$	134.31	True	2.11	2406	9.51	True	0.15	1
$N_{5,9}$	4.04	True	0.05	111	1.86	True	0.01	1

Table 9: Verification results for property  $P_3$  on 45 ACAS Xu networks.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result.  $OSS$  describes the number of the output star sets.



$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{1,1}$	424.83	True	29.77	19142	12.57	False	4.78	1
$N_{1,2}$	366.72	True	23.35	13143	18.92	False	4.37	1
$N_{1,3}$	277.95	True	13.40	9837	22.69	False	1.40	1
$N_{1,4}$	24.99	True	1.42	1184	5.78	False	0.67	1
$N_{1,5}$	208.17	True	6.32	6608	6.67	False	0.39	1
$N_{1,6}$	117.16	True	3.32	4443	9.87	True	0.92	1
$N_{2,1}$	123.79	True	5.38	5066	12.60	False	2.10	1
$N_{2,2}$	149.75	True	6.18	4500	14.89	False	2.42	1
$N_{2,3}$	27.12	True	0.99	1087	3.62	True	0.72	1
$N_{2,4}$	22.95	True	0.51	913	8.44	True	0.06	1
$N_{2,5}$	88.98	True	2.45	3419	11.61	True	0.45	1
$N_{2,6}$	46.37	True	0.97	1462	15.22	True	0.46	1
$N_{2,7}$	18.88	True	0.29	555	5.68	True	0.08	1
$N_{2,8}$	126.04	True	1.03	1805	51.33	False	1.45	1
$N_{2,9}$	8.05	True	0.06	157	1.85	True	0.01	1
$N_{3,1}$	160.56	True	5.17	4281	9.46	True	1.15	1
$N_{3,2}$	231.23	True	14.38	8708	4.15	True	1.19	1
$N_{3,3}$	25.16	True	1.25	1201	2.44	True	0.16	1
$N_{3,4}$	31.60	True	1.22	1214	4.10	True	0.27	1
$N_{3,5}$	122.59	True	5.01	3630	26.23	True	1.13	1
$N_{3,6}$	62.39	True	1.21	1495	14.23	True	0.85	1
$N_{3,7}$	55.24	True	0.48	862	5.02	True	0.12	1
$N_{3,8}$	20.56	True	0.56	542	8.46	False	0.35	1
$N_{3,9}$	148.09	True	1.57	2684	15.36	True	0.95	1
$N_{4,1}$	19.95	True	0.87	848	1.34	True	0.35	1
$N_{4,2}$	38.94	True	1.41	1348	9.42	True	2.55	1
$N_{4,3}$	78.86	True	3.72	3725	12.85	True	4.57	1
$N_{4,4}$	58.67	True	0.82	1253	19.74	False	1.14	1
$N_{4,5}$	45.51	True	0.71	1255	8.60	True	0.23	1
$N_{4,6}$	87.67	True	1.65	2366	11.26	True	0.56	1
$N_{4,7}$	6.78	True	0.10	216	3.65	True	0.08	1
$N_{4,8}$	79.35	True	1.01	1591	9.29	True	0.09	1
$N_{4,9}$	139.12	True	1.75	2566	9.37	True	0.20	1
$N_{5,1}$	166.59	True	9.26	6932	9.45	True	0.63	1
$N_{5,2}$	117.59	True	6.30	4361	4.22	True	0.56	1
$N_{5,3}$	42.58	True	2.18	1662	7.46	True	0.71	1
$N_{5,4}$	40.08	True	1.24	1088	5.68	True	0.16	1
$N_{5,5}$	52.80	True	0.89	1549	7.65	True	0.24	1
$N_{5,6}$	27.96	True	0.50	677	6.74	True	0.47	1
$N_{5,7}$	6.08	True	0.06	157	2.63	True	0.04	1
$N_{5,8}$	24.29	True	0.28	502	12.76	True	0.73	1
$N_{5,9}$	34.99	True	0.55	992	5.26	True	0.15	1

Table 10: Verification results for property  $P_4$  on 42 ACAS Xu networks.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result.  $OSS$  describes the number of the output star sets.

$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{1,1}$	2933.99	True	352.76	59734	461.70	False	8.63	1

Table 11: Verification results for property  $P_5$  on 1 ACAS Xu network.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result.  $OSS$  describes the number of the output star sets.

$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{1,1}$	44026.93	True	1159.88	187775	1083.38	False	17.91	1

Table 12: Verification results for property  $P_6$  on 1 ACAS Xu network.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result.  $OSS$  describes the number of the output star sets.

$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{11}$	-	-	-	-	1520.91	False	147.38	1

Table 13: Verification results for property  $P_7$  on 1 ACAS Xu network.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result.  $OSS$  describes the number of the output star sets. The cells with (-) correspond to networks in which our algorithm was not able to compute the reachability set successfully in less than 48 hours

$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{2,9}$	-	-	-	-	1560.52	False	46.37	1

Table 14: Verification results for property  $P_8$  on 1 ACAS Xu network.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result.  $OSS$  describes the number of the output star sets. The cells with (-) correspond to networks in which our algorithm was not able to compute the reachability set successfully in less than 48 hours.

$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{3,3}$	33727.84	False	458.09	338600	541.62	False	7.83	1

Table 15: Verification results for property  $P_9$  on 1 ACAS Xu network.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result.  $OSS$  describes the number of the output star sets.

$N_{x,y}$	Exact				Overapproximation			
	RT(s)	RES	CT(s)	OSS	RT(s)	RES	CT(s)	OSS
$N_{4,5}$	3281.44	True	181.59	41088	1087.41	False	17.68	1

Table 16: Verification results for property  $P_{10}$  on 1 ACAS Xu network.  $RT$  is the reachable set computation time, and  $CT$  is the safety checking time, both in seconds.  $RES$  is the safety verification result.  $OSS$  describes the number of the output star sets.

### A.3 Sonar Binary Classifier Detailed Results

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
1	4359	False	783	True	263	True	102	True
2	1848	False	100	True	101	True	101	True
3	206243	False	1284	True	245	True	100	True
4	49216	False	220	False	98	True	102	True
5	253978	False	1248	True	99	True	123	True
6	6452	False	622	True	276	True	101	True
7	25640	False	98	False	99	False	103	False
8	646937	False	18313	False	448	True	101	True
9	5149	False	100	False	101	False	103	False
10	6860	False	1757	False	240	True	107	True
11	10646	False	317	True	104	True	102	True
12	2542	False	281	True	100	True	100	True
13	125001	False	418	False	101	True	101	True
14	2022	False	122	False	99	False	101	False
15	14478	False	240	False	99	False	102	False
16	17343	False	263	True	99	True	102	True
17	21727	False	250	False	97	True	101	True
18	289475	False	1066	False	98	True	101	True
19	6654	True	102	True	99	True	100	True
20	-	-	46162	False	99	True	101	True
21	53343	False	99	True	100	True	101	True
22	26779	False	192	True	100	True	104	True

Continued on next page

Table 17 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
23	14711	False	444	True	98	True	101	True
24	11783	True	309	True	100	True	101	True
25	3145	True	255	True	100	True	102	True
26	3529	False	99	True	104	True	101	True
27	966607	False	16329	False	476	True	194	True
28	6241	False	119	False	100	True	107	True
29	180773	False	103	True	101	True	101	True
30	96317	False	459	True	251	True	102	True
31	6667	False	111	True	100	True	103	True
32	18427	False	253	True	100	True	102	True
33	132122	True	307	True	100	True	101	True
34	478050	False	6562	False	395	True	101	True
35	35769	False	223	True	101	True	101	True
36	34491	False	1266	False	282	True	101	True
37	47224	False	101	True	106	True	102	True
38	1127	False	102	True	101	True	101	True
39	10811	False	103	True	101	True	102	True
40	15289	True	140	True	101	True	102	True
41	3926	False	690	True	103	True	101	True
42	82744	True	444	True	100	True	102	True
43	15145	True	166	True	100	True	103	True
44	2509	True	100	True	105	True	100	True
45	30090	False	117	True	99	True	100	True
46	754722	False	789	False	100	True	104	True
47	105796	False	1501	False	101	True	101	True
48	5791	False	99	False	100	False	104	False
49	19318	False	103	False	100	True	100	True
50	24453	False	217	False	101	True	100	True
51	12330	False	436	True	99	True	101	True
52	337889	True	100	True	101	True	103	True
53	4480	False	112	True	116	True	101	True
54	8742	False	100	True	102	True	101	True
55	2912	False	257	False	100	True	100	True
56	9405	False	104	False	102	True	102	True
57	15708	False	118	False	100	True	101	True
58	8188	False	163	True	100	True	104	True
59	2663	False	102	True	103	True	101	True
60	14901	False	102	True	101	True	101	True
61	1608	False	221	True	100	True	101	True
62	2251	False	179	True	104	True	101	True
63	2915	False	103	True	101	True	102	True

Continued on next page

Table 17 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
64	177838	True	259	True	267	True	103	True
65	3152	True	114	True	103	True	103	True
66	11178	True	102	True	103	True	102	True
67	54060	True	585	True	161	True	115	True
68	6936	False	276	True	100	True	101	True
69	6557	False	389	True	103	True	101	True
70	16761	False	306	True	183	True	188	True
71	40704	False	275	True	100	True	99	True
72	7142	False	384	True	105	True	100	True
73	97419	False	1020	False	101	True	101	True
74	8259	False	1175	False	101	True	102	True
75	118488	True	100	True	102	True	101	True
76	38638	False	246	True	100	True	100	True
77	24556	False	101	True	100	True	101	True
78	88296	False	399	False	102	True	101	True
79	67012	False	244	True	101	True	116	True
80	35943	False	388	False	161	True	100	True
81	40865	False	103	True	103	True	102	True
82	-	-	401	True	100	True	101	True
83	1435255	False	943	True	100	True	101	True
84	879783	False	2954	False	104	True	101	True
85	-	-	466010	False	286	True	101	True
86	59286	False	210	True	100	True	100	True
87	181884	False	1866	True	100	True	101	True
88	30037	False	573	True	285	True	276	True
89	79237	False	449	True	103	True	101	True
90	15188	False	193	False	100	True	100	True
91	3272	True	131	True	100	True	101	True
92	13264	False	246	False	100	True	102	True
93	30169	False	283	False	103	False	101	False
94	3518	False	1950	False	255	True	251	True
95	2507	False	99	True	100	True	102	True
96	28120	True	174	True	111	True	102	True
97	167384	True	414	True	209	True	125	True
98	33945	True	3768	True	401	True	308	True
99	7974	True	359	True	103	True	102	True
100	12175	True	4981	True	233	True	102	True
101	1486	True	100	True	101	True	102	True
102	254	True	102	True	101	True	101	True
103	1281	True	208	True	115	True	102	True
104	4265	True	167	True	101	True	102	True

Continued on next page

Table 17 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
105	1894	True	99	True	107	True	101	True
106	18441	True	101	True	100	True	101	True
107	50497	False	341	True	104	True	101	True
108	4747	True	1125	True	101	True	101	True
109	74918	True	101	True	100	True	102	True
110	4394	True	763	True	268	True	101	True
111	11408	True	103	True	101	True	100	True
112	9260	True	163	True	99	True	101	True
113	2292	True	99	True	100	True	99	True
114	15426	True	544	True	101	True	101	True
115	13498	True	99	True	101	True	101	True
116	28109	True	100	True	103	True	101	True
117	9277	False	219	True	102	True	101	True
118	15982	True	148	True	100	True	101	True
119	1954	True	285	True	100	True	101	True
120	1786	True	283	True	102	True	102	True
121	3574	True	505	True	286	True	102	True
122	20352	True	100	True	100	True	101	True
123	4438	True	177	True	121	True	100	True
124	6523	True	101	True	119	True	101	True
125	4662	True	286	True	101	True	100	True
126	1024	True	170	True	171	True	100	True
127	61120	True	139	True	100	True	101	True
128	2719	True	278	True	100	True	100	True
129	11209	True	150	True	101	True	100	True
130	4386	True	100	True	107	True	100	True
131	2130	True	162	True	99	True	101	True
132	10842	True	179	True	100	True	100	True
133	8829	True	100	True	101	True	101	True
134	1257	True	340	True	101	True	101	True
135	13869	True	171	True	103	True	99	True
136	1025	True	131	True	100	True	101	True
137	1991	True	235	True	99	True	101	True
138	1061	True	99	True	102	True	104	True
139	5666	True	99	True	100	True	100	True
140	1339	True	490	True	99	True	103	True
141	1115	True	98	True	99	True	100	True
142	7571	True	159	True	100	True	101	True
143	463	True	99	True	102	True	104	True
144	1871	True	262	True	100	True	100	True
145	21578	True	341	True	100	True	101	True

Continued on next page

Table 17 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
146	21035	True	475	True	193	True	104	True
147	572	True	98	True	99	True	100	True
148	563	True	98	True	102	True	99	True
149	36945	True	99	True	130	True	101	True
150	6400	True	100	True	109	True	100	True
151	68362	True	102	True	101	True	100	True
152	5669	True	151	True	106	True	140	True
153	17600	True	156	True	100	True	101	True
154	11486	True	103	True	102	True	99	True
155	2891	True	237	True	100	True	101	True
156	68190	True	580	True	101	True	101	True
157	10525	True	99	True	100	True	102	True
158	1928	True	98	True	99	True	101	True
159	2401	True	272	True	102	True	100	True
160	6601	True	99	True	101	True	101	True
161	1920	True	223	True	100	True	100	True
162	1038	True	273	True	115	True	101	True
163	3557	True	99	True	101	True	101	True
164	25147	True	397	True	100	True	100	True
165	11193	True	98	True	100	True	105	True
166	1910	True	99	True	101	True	100	True
167	3131	True	320	True	116	True	114	True
168	1138	True	240	True	101	True	110	True
169	1829	True	164	True	105	True	100	True
170	1653	True	170	True	101	True	100	True
171	1018	True	169	True	163	True	107	True
172	736	True	99	True	106	True	100	True
173	4313	True	99	True	105	True	101	True
174	37584	True	99	True	102	True	102	True
175	538	True	192	True	197	True	101	True
176	10231	True	343	True	157	True	153	True
177	4989	True	111	True	100	True	106	True
178	2044	True	101	True	101	True	108	True
179	649	True	101	True	105	True	101	True
180	592	True	101	True	101	True	101	True
181	332	True	100	True	99	True	101	True
182	166	True	99	True	99	True	100	True
183	382	True	99	True	100	True	101	True
184	368	True	100	True	100	True	100	True
185	515	True	191	True	100	True	100	True
186	461	True	99	True	101	True	101	True

Continued on next page

Table 17 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
187	808	True	164	True	99	True	100	True
188	2518	True	98	True	109	True	102	True
189	608	True	100	True	100	True	101	True
190	1162	True	99	True	99	True	100	True
191	1297	True	100	True	100	True	101	True
192	2890	True	100	True	100	True	102	True
193	2765	True	99	True	100	True	102	True
194	2336	True	164	True	102	True	102	True
195	624	True	100	True	100	True	102	True
196	1340	True	102	True	99	True	100	True
197	5056	True	258	True	100	True	101	True
198	3146	True	261	True	101	True	102	True
199	1367	True	261	True	102	True	100	True
200	2288	True	235	True	102	True	101	True
201	2449	True	201	True	101	True	102	True
202	3083	True	99	True	100	True	101	True
203	2729	True	100	True	99	True	101	True
204	2680	True	288	True	285	True	101	True
205	864	True	100	True	99	True	100	True
206	1089	True	99	True	101	True	101	True
207	866	True	101	True	99	True	103	True
208	3258	True	168	True	99	True	101	True

Table 17: Local adversarial robustness tests of the exact approach. *RT* is the reachable set computation time in milliseconds. *RES* is the safety verification result. True indicates that the neural network correctly classifies the input set as expected, while False means that the neural network was unable to correctly classify the input set. Cells with (-) indicate cases where timeout occurs.

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
1	234	Inconclusive	205	True	163	True	103	True
2	278	True	103	True	102	True	103	True
3	396	Inconclusive	279	True	157	True	103	True
4	480	Inconclusive	128	True	101	True	104	True
5	391	Inconclusive	260	True	101	True	104	True
6	341	True	203	True	157	True	103	True
7	437	Inconclusive	104	False	101	False	104	False
8	403	Inconclusive	371	False	181	True	106	True
9	364	Inconclusive	102	False	101	False	103	False

Continued on next page



Table 18 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
10	351	Inconclusive	268	False	144	True	107	True
11	352	False	176	True	103	True	145	True
12	320	True	166	True	102	True	102	True
13	394	Inconclusive	137	False	102	True	104	True
14	289	Inconclusive	114	False	101	False	108	False
15	438	Inconclusive	214	False	101	False	102	False
16	362	False	134	True	101	True	103	True
17	321	Inconclusive	140	True	100	True	110	True
18	420	Inconclusive	193	False	100	True	104	True
19	299	True	107	True	100	True	102	True
20	454	Inconclusive	390	True	100	True	106	True
21	459	Inconclusive	105	True	101	True	103	True
22	498	True	143	True	101	True	103	True
23	369	Inconclusive	204	True	101	True	107	True
24	385	True	160	True	101	True	103	True
25	342	True	154	True	103	True	104	True
26	310	Inconclusive	101	True	103	True	105	True
27	554	Inconclusive	384	True	172	True	137	True
28	338	Inconclusive	112	True	102	True	103	True
29	564	Inconclusive	106	True	106	True	106	True
30	494	Inconclusive	167	True	158	True	104	True
31	296	Inconclusive	106	True	101	True	103	True
32	336	Inconclusive	164	True	106	True	107	True
33	512	True	156	True	102	True	104	True
34	474	Inconclusive	292	False	199	True	103	True
35	377	Inconclusive	140	True	102	True	104	True
36	436	Inconclusive	246	True	180	True	103	True
37	380	True	102	True	103	True	107	True
38	211	Inconclusive	104	True	103	True	106	True
39	436	True	101	True	111	True	104	True
40	345	True	105	True	103	True	104	True
41	321	Inconclusive	175	True	102	True	105	True
42	424	True	191	True	106	True	104	True
43	374	True	118	True	102	True	103	True
44	247	True	103	True	102	True	106	True
45	504	Inconclusive	108	True	108	True	113	True
46	574	Inconclusive	186	True	102	True	104	True
47	410	Inconclusive	253	False	102	True	106	True
48	379	Inconclusive	105	False	105	False	102	False
49	411	Inconclusive	105	False	103	True	104	True
50	396	Inconclusive	140	True	101	True	106	True

Continued on next page

Table 18 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
51	415	True	198	True	104	True	102	True
52	497	True	106	True	103	True	103	True
53	346	Inconclusive	111	True	106	True	111	True
54	339	Inconclusive	107	True	103	True	103	True
55	350	Inconclusive	145	False	102	True	104	True
56	321	Inconclusive	102	False	103	True	107	True
57	361	Inconclusive	105	True	107	True	103	True
58	323	True	120	True	102	True	103	True
59	329	False	102	True	103	True	106	True
60	443	True	102	True	104	True	103	True
61	298	False	143	True	102	True	103	True
62	431	False	128	True	104	True	107	True
63	442	False	103	True	110	True	104	True
64	434	True	153	True	156	True	104	True
65	337	True	107	True	104	True	105	True
66	402	True	103	True	103	True	119	True
67	441	True	207	True	131	True	104	True
68	425	False	157	True	109	True	102	True
69	378	False	181	True	102	True	126	True
70	426	Inconclusive	155	True	133	True	147	True
71	418	Inconclusive	164	True	105	True	103	True
72	416	Inconclusive	180	True	104	True	103	True
73	447	Inconclusive	203	False	104	True	105	True
74	298	Inconclusive	218	False	103	True	104	True
75	507	True	103	True	104	True	103	True
76	479	Inconclusive	145	True	105	True	103	True
77	502	Inconclusive	105	True	104	True	104	True
78	476	Inconclusive	176	True	105	True	102	True
79	505	Inconclusive	144	True	107	True	102	True
80	389	Inconclusive	186	True	128	True	103	True
81	495	Inconclusive	103	True	102	True	103	True
82	569	Inconclusive	176	True	102	True	104	True
83	541	Inconclusive	251	True	104	True	103	True
84	515	Inconclusive	337	True	104	True	104	True
85	578	Inconclusive	413	False	158	True	103	True
86	349	Inconclusive	130	True	103	True	103	True
87	414	False	251	True	104	True	103	True
88	367	Inconclusive	202	True	158	True	171	True
89	427	Inconclusive	161	True	102	True	102	True
90	326	Inconclusive	126	False	102	True	104	True
91	259	True	110	True	105	True	103	True

Continued on next page

Table 18 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
92	436	Inconclusive	150	True	103	True	104	True
93	410	Inconclusive	152	True	103	False	104	False
94	293	Inconclusive	208	True	147	True	146	True
95	286	Inconclusive	103	True	103	True	102	True
96	411	True	126	True	114	True	105	True
97	586	True	139	True	128	True	112	True
98	407	Inconclusive	367	True	177	True	174	True
99	339	True	167	True	104	True	101	True
100	405	Inconclusive	301	True	145	True	104	True
101	259	True	103	True	110	True	104	True
102	141	Inconclusive	106	True	107	True	103	True
103	359	True	142	True	101	True	103	True
104	317	True	127	True	104	True	105	True
105	313	True	101	True	108	True	103	True
106	329	Inconclusive	101	True	103	True	102	True
107	503	Inconclusive	151	True	101	True	104	True
108	342	Inconclusive	229	True	108	True	103	True
109	401	Inconclusive	102	True	102	True	103	True
110	383	Inconclusive	225	True	156	True	103	True
111	366	Inconclusive	102	True	141	True	102	True
112	299	True	119	True	129	True	102	True
113	305	Inconclusive	102	True	103	True	103	True
114	379	Inconclusive	190	True	104	True	105	True
115	378	Inconclusive	101	True	103	True	103	True
116	413	Inconclusive	102	True	103	True	104	True
117	377	Inconclusive	136	True	107	True	103	True
118	361	True	116	True	103	True	103	True
119	293	Inconclusive	159	True	103	True	105	True
120	260	True	149	True	109	True	103	True
121	320	Inconclusive	218	True	185	True	104	True
122	398	True	101	True	105	True	103	True
123	328	True	127	True	103	True	103	True
124	360	Inconclusive	102	True	102	True	103	True
125	328	True	165	True	103	True	103	True
126	251	Inconclusive	125	True	131	True	103	True
127	354	True	113	True	104	True	103	True
128	331	Inconclusive	150	True	108	True	104	True
129	264	True	120	True	103	True	102	True
130	314	True	102	True	102	True	102	True
131	233	True	126	True	138	True	103	True
132	285	True	123	True	101	True	103	True

Continued on next page

Table 18 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
133	359	True	103	True	106	True	104	True
134	185	True	157	True	105	True	103	True
135	403	True	126	True	105	True	102	True
136	260	True	112	True	103	True	102	True
137	232	True	144	True	103	True	105	True
138	226	True	102	True	102	True	103	True
139	380	True	101	True	104	True	102	True
140	256	True	232	True	102	True	103	True
141	293	True	101	True	107	True	105	True
142	352	True	118	True	101	True	102	True
143	197	True	101	True	104	True	104	True
144	296	True	158	True	104	True	102	True
145	359	True	149	True	103	True	102	True
146	374	True	182	True	145	True	103	True
147	195	True	100	True	103	True	106	True
148	192	True	101	True	104	True	101	True
149	484	True	102	True	134	True	102	True
150	376	True	102	True	109	True	105	True
151	564	Inconclusive	102	True	103	True	103	True
152	375	Inconclusive	116	True	103	True	138	True
153	407	Inconclusive	120	True	110	True	122	True
154	369	Inconclusive	101	True	104	True	102	True
155	303	Inconclusive	149	True	107	True	111	True
156	528	Inconclusive	196	True	103	True	104	True
157	355	True	100	True	103	True	102	True
158	315	True	102	True	102	True	104	True
159	302	True	141	True	103	True	104	True
160	304	Inconclusive	101	True	105	True	104	True
161	353	True	136	True	103	True	102	True
162	224	True	151	True	104	True	101	True
163	351	True	101	True	137	True	103	True
164	446	Inconclusive	155	True	102	True	103	True
165	401	True	100	True	105	True	103	True
166	276	Inconclusive	101	True	103	True	104	True
167	336	True	175	True	110	True	108	True
168	253	Inconclusive	159	True	104	True	102	True
169	281	Inconclusive	120	True	102	True	110	True
170	282	True	123	True	109	True	102	True
171	227	Inconclusive	125	True	122	True	102	True
172	224	True	101	True	101	True	106	True
173	323	True	101	True	116	True	103	True

Continued on next page

Table 18 – continued from previous page

	$\delta = 0.01$		$\delta = 0.001$		$\delta = 0.0001$		$\delta = 0.00001$	
	RT	RES	RT	RES	RT	RES	RT	RES
174	385	True	102	True	109	True	102	True
175	213	True	138	True	137	True	104	True
176	278	True	156	True	118	True	120	True
177	324	True	105	True	101	True	103	True
178	284	Inconclusive	101	True	104	True	103	True
179	193	Inconclusive	103	True	104	True	105	True
180	227	True	101	True	102	True	103	True
181	189	True	102	True	102	True	104	True
182	124	True	100	True	102	True	102	True
183	171	True	101	True	102	True	102	True
184	170	True	102	True	102	True	103	True
185	190	True	130	True	103	True	102	True
186	194	True	101	True	105	True	102	True
187	248	True	127	True	102	True	105	True
188	348	True	101	True	101	True	104	True
189	221	True	102	True	103	True	103	True
190	240	Inconclusive	102	True	102	True	103	True
191	246	True	102	True	102	True	103	True
192	293	Inconclusive	104	True	105	True	104	True
193	357	True	102	True	103	True	105	True
194	261	Inconclusive	121	True	103	True	103	True
195	208	Inconclusive	101	True	104	True	103	True
196	248	True	103	True	102	True	102	True
197	351	True	175	True	103	True	103	True
198	297	True	142	True	105	True	102	True
199	245	True	149	True	105	True	104	True
200	283	True	154	True	104	True	103	True
201	275	True	144	True	108	True	102	True
202	280	True	102	True	102	True	106	True
203	323	True	102	True	102	True	102	True
204	274	True	165	True	177	True	103	True
205	247	True	101	True	102	True	105	True
206	250	True	102	True	103	True	103	True
207	236	True	101	True	101	True	104	True
208	302	Inconclusive	129	True	101	True	102	True

Table 18: Local adversarial robustness tests of the over-approximate approach. *RT* is the reachable set computation time in milliseconds. *RES* is the safety verification result. True indicates that the neural network correctly classifies the input set as expected, while False means that the neural network was unable to correctly classify the input set. Additionally, Inconclusive is assigned when the reachability analysis algorithm cannot provide a conclusive answer due to the over-approximation errors.



# Certified Control for Train Sign Classification\*

Jan Roßbach<sup>id</sup>

Heinrich-Heine-Universität Düsseldorf  
Mathematisch-Naturwissenschaftliche Fakultät  
Institut für Informatik  
jan.rossbach@uni-duesseldorf.de

Michael Leuschel<sup>id</sup>

Heinrich-Heine-Universität Düsseldorf  
Mathematisch-Naturwissenschaftliche Fakultät  
Institut für Informatik  
leuschel@uni-duesseldorf.de

There is considerable industrial interest in integrating AI techniques into railway systems, notably for fully autonomous train systems. The KI-LOK research project is involved in developing new methods for certifying such AI-based systems. Here we explore the utility of a certified control architecture for a runtime monitor that prevents false positive detection of traffic signs in an AI-based perception system. The monitor uses classical computer vision algorithms to check if the signs – detected by an AI object detection model – fit predefined specifications. We provide such specifications for some critical signs and integrate a Python prototype of the monitor with a popular object detection model to measure relevant performance metrics on generated data. Our initial results are promising, achieving considerable precision gains with only minor recall reduction; however, further investigation into generalization possibilities will be necessary.

## 1 Introduction and Motivation

Artificial intelligence has been increasingly used in various sectors, including transportation [16]. One particular area where artificial intelligence (AI) has gained attention is the development of autonomous driving systems for railways [12]. The results already achieved in other transport sectors, mainly automotive, have encouraged the development of AI in the railway industry [19].

While this technology holds high economic interest, reliable certification methods are necessary to ensure safe and regulated access to these innovations [12]. Traditional verification approaches such as formal methods have faced difficulties in this area due to the opaque nature of AI, particularly in computer vision where class definitions for classification tasks based on raw pixel values have been considered challenging.

The KI-LOK<sup>1</sup> research project addresses these challenges by developing certification methodologies for autonomous AI-based railway systems. As part of this, a case study [5] on train movements during shunting movements is being analyzed. A formal B [1] model has been developed [5] to analyze the environment and ensure the safety of the deterministic steering system through model checking with the PROB [11] model checker. The safety of the system was found to be conditional on correct results from the AI-based perception system. In this work, we attempt to move towards verification of part of this perception system using a runtime monitor with a certified control [8] architecture. This architecture reduces the part of the system requiring formal verification compared to traditional monitor architectures putting a more formal analysis back into reach. In particular, we focus on a subset of the train sign classification component. It is responsible for detecting and classifying signs in the shunting yard to ensure safe train movements. False recognition of a 'track-free' (Sh1) signal has been determined to have

---

\*This research is part of the KI-LOK project funded by the “Bundesministerium für Wirtschaft und Energie”; grant # 19/21007E.

<sup>1</sup><https://ki-lok.itpower.de>

safety implications. We aim to significantly reduce or eliminate such false positives for some of the most critical classes by defining a sign-specific ontology and checking it at runtime. For this we introduce such a specification and show the potential performance gains by evaluating a prototype implementation in Python on a custom dataset.

## 2 Background and Related Work

The case study[5] being considered has been developed by Thales (now Ground Transportation Systems) and focuses on a train during shunting movements. The system includes an AI-based perception system and a deterministic steering system. The role of the perception system is to detect and classify obstacles (persons, animals, vehicles, ...) and railway infrastructure elements. The steering system then makes appropriate decisions about moving the locomotive based on that information.

There was a set of requirements provided with the case study, including the correct detection of several shunting train signs. In order to increase confidence in the perception system we aim to check the recognized signs with a runtime monitor. This will give strong confidence that detected signs are correct. In order to safeguard against unrecognized signs we will need to lean on other measures taken by the project, like a thorough environment ontology and systematic test case generation [4].

### 2.1 Certified Control

Certified Control[8] is an architectural framework for the real-time validation of autonomous systems. It distinguishes itself from conventional monitoring components by omitting its reliance on independent perception and instead counting on the controller to provide a *certificate* containing all essential information. This certificate serves as input for the runtime monitor, which assesses the accuracy of system behavior against specified criteria. By adopting this approach, the architecture establishes a trusted foundation that can potentially be subjected to a rigorous formal verification process.

The controller, which is not included in the *trusted base*, can utilize sophisticated algorithms such as neural networks without needing explicit formal verification. By separating the tasks of generating visual insights and ensuring safety, established verification methods can continue to be used with minimal adjustments. To accomplish this, a formal acceptance specification for the certificate is necessary to ensure compliance with safety requirements like *the detected lane lines are parallel* or *there are no objects on the track for 100m*. This reduces the amount of code needing verification and allows the AI components to go unverified.

While the effectiveness of this architecture in lane line detection for regular vehicles is promising [8], its applicability to other autonomous perception tasks such as sign classification and object detection remains uncertain. Therefore, we aim to investigate the applicability and effectiveness of such a certified control architecture in the context of the case studies train control perception system.

### 2.2 Related Work

Other attempts at verifying an autonomous train perception systems notably include [12]. The authors propose a multi-sensor pipeline relying on the statistical independence of the different perception mechanisms to control hazards and ensure suitable model performance. The goal is to show possible ways of certifying according to the ANSI/UL 4600 [6] standard, which provides a framework for integrating AI into fully autonomous systems. The standard gives practical guidelines and advice for a possible safety case, notably including the entire autonomy pipeline and AI algorithms. We also hope to provide methods



to aid with a verification according to this standard, while a full certification is currently out of reach. Other approaches to formal runtime monitor verification of AI systems have been done in the field of reinforcement learning using safety shields [10]. But these approaches focus on training agents to choose optimal policies depending on given environmental factors, which is similar to the traditional steering system in our model. There have also been proposals for formalizing image specification, including spatial model checking [2] and attempts to formalize vision ontology [13, 14].

### 3 Specification and Ontology

The selected sign classes for verification are Sh0, Sh1, and Wn7 as depicted in Figure 1. While these look similar, the semantic content is different. Sh0 means stop and the others signal safe passage. This makes properly distinguishing them a safety-critical issue. To ensure that the train comes to a stop when encountering a Sh0 sign on the current track, it is crucial to accurately detect and locate it. To achieve this, we employ an AI object detection system in the controller. Subsequently, the monitor verifies if the bounding box image aligns with the expected ontology. This provides additional confidence in the accuracy of the result.

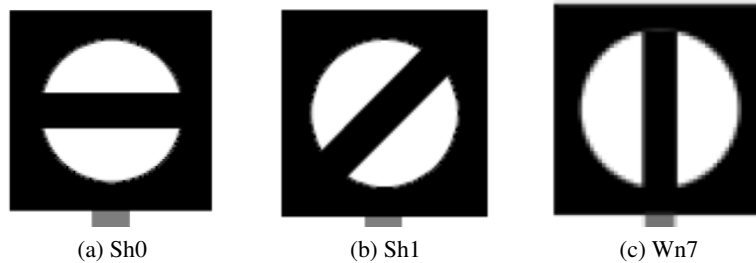


Figure 1: Train Control Shunting Signs

It is often challenging to provide a precise formal definition of an image class based solely on its features. Instead, we focus directly on detectable image characteristics. In this context, we can observe that the images include two semi-circles with only orientation as the distinguishing feature. This characteristic feature allows us to define the sign using the contours and orientation angles of the feature.

For a given image tensor  $I$  with height  $h$  and width  $w$ , consider the set of contours (sets of points) denoted as  $C(I)$ , which are identified by a contour detection algorithm. Let  $S_0$  be the set of images belonging to the Sh0 class. Also define  $A : C(I) \rightarrow \mathbb{R}^+$  as the area function, which calculates the area of a given contour. Similarly, let  $\sigma : C(I) \rightarrow \mathbb{Z}^+$  be an orientation function that determines the angle between the contour and the horizontal axis. We can then express membership of an image to one of the classes by considering an image a member of the set  $S_0$  if it contains a pair  $(c_1, c_2) \in C(I) \times C(I)$ , which fullfills all the following conditions, given some pre-determined error tolerances  $\delta_i, i \in \{1, 2, 3, 4, 5\}$ <sup>2</sup> and an expected angle  $a$  that depends on the class in question.

1.  $A(c_1)(1 - \delta_1) \leq A(c_2) \leq (1 + \delta_1)A(c_1)$
2.  $(1 - \delta_2)\sigma(c_1) \leq \sigma(c_2) \leq (1 + \delta_2)\sigma(c_1)$
3.  $\delta_3 h \leq A(c_i) \leq \delta_4 h, i \in 1, 2$

<sup>2</sup>In the prototype implementation the tolerance values used were  $\delta_{1,2,5} = 0.2$ ,  $\delta_3 = 0.1$  and  $\delta_4 = 0.3$

4.  $\delta_3 w \leq A(c_i) \leq \delta_4 w, i \in 1, 2$
5.  $c_1 \cap c_2 = \emptyset$
6.  $|\sigma(c_i) - a| \leq 90\delta_5, a = 0$

For the remaining two classes, the expected angle  $a$  in the final condition varies to 45 for Sh1 and 90 for Wn7. Otherwise, the definitions are identical. The conditions one to six define an Sh0 sign as an image with two contours that have similar angles and orientations. The orientation should be within a certain error threshold. Also, the definition expects, that the areas do not overlap. While ideally, we expect an orientation of zero, variations can occur due to different photo angles. Thus, the inclusion of an error term accounts for this discrepancy in measurement accuracy.

This definition is not flawless and permits the possibility of false positives. This implies that there may be instances where images that do not depict the intended sign could potentially be accepted (see Figure 2a). However, incorporating this check reduces the likelihood of such occurrences compared to those without it. The stringency of the monitoring process needs to be weighed against the decrease in true positives to strike a suitable balance. Adjustments can be made by selecting appropriate  $\delta$  values within certain limits. Now we can define a requirement for a correct implementation.

*REC*: The implementation accurately verifies whether an image meets the ontology requirements of a specific class.

## 4 Implementation and Experiments

While the following implementation is not yet verified in terms of *REC*, we aim to do so in future work. Here we provide a prototype, which is developed enough to indicate the potential usefulness of such an implementation. Given an image and an expected class, it either validates or rejects the image. We then integrated it with a YOLOv8 object detection model and measured the influence on common performance metrics (see. Table 2b). In the following sections, we present details on the implementation and the performed experiments.

### 4.1 Implementation

The controller component is a simple wrapper for the YOLOv8<sup>3</sup> implementation of an object detection model known as YOLO [15]. The outcomes obtained from this model are packaged into a certificate and transmitted to the monitor. To have the model detect the signs in question, we created and labeled a custom sign-detection dataset [9], on which we trained three model variants. These were the nano, small and medium versions of the model with 3.2M, 11.2M and 25.9M parameters respectively. The training was done for 200 epochs with a batch size of 16. They achieved mAP50 values of 0.827, 0.90 and 0.93 on the test set.

From the model results the controller generates a *certificate* – in the sense of certified control (see Section 2) – consisting of the following components:

1. The original image.
2. The assigned class result.
3. The bounding box, represented as a tuple in the format  $(x, y, w, h)$ , with values normalized to fit the dimensions of the image.

---

<sup>3</sup><https://github.com/ultralytics/ultralytics>

This Python object is then given to the monitor. In a production implementation, it would be preferable to serialize and send this data to a statically typed version of the monitor for optimal security.

The monitor implementation utilizes Python’s OpenCV [7] library to apply simple and well-tested computer vision algorithms to the given images. To begin, the bounding box image is resized to 206x206 and converted to grayscale to facilitate contour detection. Subsequently, a filtering process is applied to the contours to ensure their area falls within the specified size boundaries (refer to Section 3). We then need to calculate the area and orientation of the detected contours to determine if some of them fit the requirements for the ontology. The area of each contour is extracted using an available function within OpenCV. In addition, we utilize OpenCV once again by fitting a line through each contour as a means of determining its orientation. With that, we can calculate the orientation using the following equation.

$$\sigma(c) = \frac{180 \arccos(\vec{e}_1 \cdot \vec{v})}{\pi |\vec{v}|}$$

Next, we evaluate the remaining contours in pairs to determine if they satisfy the similarity conditions for area and orientation (refer to Section 3 for details). If a pair is found that meets these conditions, we then verify if its orientation aligns with the expected orientation for the corresponding class. If it does, the monitoring system considers this as a valid certificate. However, if any of these criteria are not met, the certificate will be rejected.

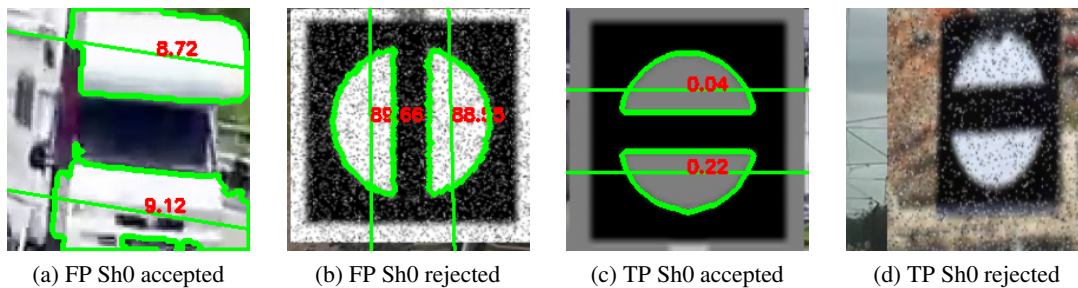


Figure 2: Visual Examples of Successful and Failing Monitor Checks

## 4.2 Experiments

In contrast to the automotive field, which benefits from large-scale image datasets like KITTI [3] for efficient object detection model evolution using road scene images, the railway industry faces limitations in terms of relevant datasets. Recently, interesting multi-sensor benchmark datasets [20] have started to emerge, but do not fit our particular use case. This lack of labeled, high-quality data poses a challenge when it comes to training and validating AI-based systems for this particular case. When evaluating the performance of the prototype, we have to confront this lack of data in the field. Since the relevant publicly available datasets do not cover the classes in question, we resort to custom labeling for training and a data generation approach for the evaluation of the system. For the generation, we chose a small number of base images of the signs in question, which are put through different random perturbation combinations and then pasted in random amounts – one to four – onto images from train footplate rides, gathered from the web. By this method we generated 28283 unique images containing 43638 signs. There are up to four signal per picture, which is typical of a shunting yard. For this work we ignore the selection of relevant signals and only focus on detection. The following perturbations were applied:

1. Horizontal Flip

Model	Detected	TP	FP	Model	Detected	TP	FP
n	30111	25514	4597	n	21716	21714	2
s	30335	26790	3545	s	22834	22831	3
m	28672	22728	5944	m	20460	20460	0

(a) Results without Monitor

(b) Results with Monitor

Table 1: Raw numbers for Models on Generated Data

Model	Precision	Recall	$F_1$ score	Model	Precision	Recall	$F_1$ score
n	0.85	0.58	0.69	n	1.00	0.50	0.67
s	0.88	0.61	0.72	s	1.00	0.52	0.68
m	0.79	0.52	0.63	m	1.00	0.47	0.64

(a) Results without Monitor

(b) Results with Monitor

Table 2: Model Metrics on Generated Data (values rounded to two decimal places)

2. Gaussian Noise (Salt and Pepper with Levels of 0.05 and 0.075)
3. Scaling (Up and back down to square images of 50,100,213,416 and 832 px)
4. Blur (normalized box filter with kernel sizes 3, 5, 7)
5. Brightness change (levels 0.5,1.5)

In Figure 2 we see examples of these images with monitor visualizations applied. It shows cut YOLO bounding boxes with the contours, lines and corresponding orientations detected by the monitor. Figure 2a shows one of the few remaining false positives. The image fits all the defined criteria of the  $S_0$  ontology for these  $\delta$  values but is not actually of that class. Given stricter tolerances (e.g.  $90\delta_5 < 8$ ) this mistake would not occur. Overall the results seen in Table 2b show a slight reduction in model performance in terms of recall and an evenly weighted F-score compared to the prior results in Table 2a. The concrete detection numbers can be found in Table 1. The drop in recall and F-score is expected due to the reduction in true positives. However, almost all false positives have been recognized and can thus be prevented. The tolerances can be adjusted to further reduce false positives, at the cost of more recall and F-score, or to allow more leeway to the perception system. In terms of runtime performance, the monitor checks a certificate in approximately 0.7 ms on an Intel i5-12600K processor. In comparison, the inference of the YOLOv8 model will range from 2 ms – for the nano model variant – to 8 ms for the m version. This means that the performance overhead is likely not a major concern in a production environment.

## 5 Conclusion and Future Work

In conclusion, this study demonstrates the potential utility of certified control runtime monitoring for object detection of formally definable and safety critical classes. The resulting trade-off in our tests is promising enough to warrant further investigation into different application possibilities. However, further research is necessary to fully validate its implementation in a type-safe language following the *REC* guidelines. The obtained results should be verified in appropriate field test for any real world application. Additionally, it should be noted that a significant portion of the perception system remains unverified. Moving forward, our future work will involve evaluating the applicability of a similar architecture for other components of the perception system such as obstacle detection. This evaluation will include examining different sensor types such as LIDAR and radar on benchmark datasets.

## References

- [1] J.R. Abrial & A. Hoare (2005): *The B-Book: Assigning Programs to Meanings*. Cambridge University Press, doi:10.1017/CBO9780511624162.
- [2] Vincenzo Ciancia, Diego Latella, Michele Loreti & Mieke Massink (2016): *Model Checking Spatial Logics for Closure Spaces*. *Log. Methods Comput. Sci.* 12(4), doi:10.2168/LMCS-12(4:2)2016.
- [3] Andreas Geiger, Philip Lenz & Raquel Urtasun (2012): *Are we ready for autonomous driving? The KITTI vision benchmark suite*. In: *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3354–3361, doi:10.1109/CVPR.2012.6248074.
- [4] Jürgen Grossmann, Nicolas Grube, Sami Kharma, Dorian Knoblauch, Roman Krajewski, Mariia Kucheiko & Hans-Werner Wiesbrock (2023): *Test and Training Data Generation for Object Recognition in the Railway Domain*. In Paolo Masci, Cinzia Bernardeschi, Pierluigi Graziani, Mario Koddenbrock & Maurizio Palmieri, editors: *Software Engineering and Formal Methods. SEFM 2022 Collocated Workshops*, Springer International Publishing, Cham, pp. 5–16, doi:10.1007/978-3-031-26236-4\_1.
- [5] Jan Gruteser, David Geleßus, Michael Leuschel, Jan Roßbach & Fabian Vu (2023): *A Formal Model of Train Control with AI-based Obstacle Detection*. In Birgit Milius, Simon Collart-Dutilleul & Thierry Lecomte, editors: *Reliability, Safety, and Security of Railway Systems. Modelling, Analysis, Verification, and Certification*, Springer Nature Switzerland, pp. 128–145, doi:10.1007/978-3-031-43366-5\_8.
- [6] Underwriters Laboratories Inc (2020): *4600 Standard for Evaluation of Autonomous Products*. Technical Report, Underwriters Laboratories Inc.
- [7] Itseez (2015): *Open Source Computer Vision Library*. <https://github.com/itseez/opencv>.
- [8] Daniel Jackson, Valerie Richmond, Mike Wang, Jeff Chow, Uriel Guajardo, Soonho Kong, Sergio Campos, Geoffrey Litt & Nikos Aréchiga (2021): *Certified Control: An Architecture for Verifiable Safety of Autonomous Vehicles*. *CoRR* abs/2104.06178, doi:10.48550/arXiv.2104.06178. arXiv:2104.06178.
- [9] KILOK (2023): *Sign Detection Dataset*. <https://universe.roboflow.com/kilok/sign-detection-4oqe4>. Visited on 2023-08-09.
- [10] Bettina Könighofer, Florian Lorber, Nils Jansen & Roderick Bloem (2020): *Shield Synthesis for Reinforcement Learning*. In Tiziana Margaria & Bernhard Steffen, editors: *Leveraging Applications of Formal Methods, Verification and Validation: Verification Principles*, Springer International Publishing, Cham, pp. 290–306, doi:10.1007/978-3-030-61362-4\_16.
- [11] Michael Leuschel & Michael Butler (2003): *ProB: A Model Checker for B*. In: *Proceedings FME, LNCS 2805*, pp. 855–874, doi:10.1007/978-3-540-45236-2\_46.
- [12] Jan Peleska, Anne E. Haxthausen & Thierry Lecomte (2022): *Standardisation Considerations for Autonomous Train Control*. In Tiziana Margaria & Bernhard Steffen, editors: *Leveraging Applications of Formal Methods, Verification and Validation. Practice*, Springer Nature Switzerland, pp. 286–307, doi:10.1007/978-3-031-19762-8\_22.
- [13] Daniele Porello, Marco Cristani & Roberta Ferrario (2013): *Integrating ontologies and computer vision for classification of objects in images*. In: *Proceedings of the Workshop on Neural-Cognitive Integration in German Conference on Artificial Intelligence*, pp. 1–15.
- [14] "K. K. Thyagarajan R. I. Minu" (2014): *Semantic Rule Based Image Visual Feature Ontology Creation*. *International Journal of Automation and Computing* 11(20140504), doi:10.1007/s11633-014-0832-3.
- [15] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick & Ali Farhadi (2016): *You Only Look Once: Unified, Real-Time Object Detection*. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE Computer Society, Los Alamitos, CA, USA, pp. 779–788, doi:10.1109/CVPR.2016.91.
- [16] Danijela Ristić-Durrant, Marten Franke & Kai Michels (2021): *A Review of Vision-Based On-Board Obstacle Detection and Distance Estimation in Railways*. *Sensors (Basel, Switzerland)*, doi:10.3390/s21103452.

- [17] Claudio Filipi Gonçalves dos Santos & João Paulo Papa (2022): *Avoiding Overfitting: A Survey on Regularization Methods for Convolutional Neural Networks*. CoRR abs/2201.03299, doi:10.1145/3510413, arXiv:2201.03299.
- [18] Satoshi Suzuki & Keichi Abe (1985): *Topological structural analysis of digitized binary images by border following*. *Computer Vision, Graphics, and Image Processing* 30(1), pp. 32–46, doi:10.1016/0734-189X(85)90016-7.
- [19] Ruifan Tang, Lorenzo De Donato, Nikola Besinovic, Francesco Flammini, Rob M.P. Goverde, Zhiyuan Lin, Ronghui Liu, Tianli Tang, Valeria Vittorini & Ziyulong Wang (2022): *A literature review of Artificial Intelligence applications in railway systems*. *Transportation Research Part C: Emerging Technologies* 140, p. 103679, doi:10.1016/j.trc.2022.103679.
- [20] Roman Tilly, Philipp Neumaier, Karsten Schwalbe, Pavel Klasek, Rustam Tagiew, Patrick Denzler, Tobias Klockau, Martin Boekhoff & Martin Köppel (2023): *Open Sensor Data for Rail 2023*, doi:10.57806/9MV146R0.

# Model Checking for Closed-Loop Robot Reactive Planning

Christopher Chandler

School of Computing Science  
University of Glasgow

`christopher.chandler@glasgow.ac.uk`

Alice Miller

School of Computing Science  
University of Glasgow

`alice.miller@glasgow.ac.uk`

Bernd Porr

School of Biomedical Engineering  
University of Glasgow

`bernd.porr@glasgow.ac.uk`

Giulia Lafratta

School of Engineering  
University of Glasgow

`giulia.lafratta@glasgow.ac.uk`

In this paper, we show how model checking can be used to create multi-step plans for a differential drive wheeled robot so that it can avoid immediate danger. Using a small, purpose built model checking algorithm *in situ* we generate plans in real-time in a way that reflects the egocentric reactive response of simple biological agents. Our approach is based on chaining temporary control systems which are spawned to eliminate disturbances in the local environment that disrupt an autonomous agent from its preferred action (or *resting state*). The method involves a novel discretization of 2D LiDAR data which is sensitive to bounded stochastic variations in the immediate environment. We operationalise multi-step planning using invariant checking by forward depth-first search, using a cul-de-sac scenario as a first test case. Our results demonstrate that model checking can be used to plan efficient trajectories for local obstacle avoidance, improving on the performance of a reactive agent which can only plan one step. We achieve this in near real-time using no pre-computed data. While our method has limitations, we believe our approach shows promise as an avenue for the development of safe, reliable and transparent trajectory planning in the context of autonomous vehicles.

## 1 Introduction

Simple biological systems (or agents) can safely navigate through a previously unseen environment by responding in real-time to sensory inputs. On sensing an unexpected input (e.g., an obstacle) the agent responds by performing an action to change its state. This action takes the form of a motor output which results in a change to the environment, which is in turn sensed by the agent and the loop repeats [4]. The behaviour is egocentric and reactive—the agent is only concerned with its immediate environment and only deviates from its course (or *resting state*) when necessary. Naturally, complex agents are capable of more sophisticated behaviours, such as the prediction of disturbances and the generation of plans to counteract them. This requires distal sensor information and spatial understanding of the wider environment. Indeed, there is evidence that in biological systems an innate “core” understanding of world physics and causality allow organisms to organise their behaviours in accordance with predicted outcomes [29, 19].

Model checking [1] is a widely used technique for automatically verifying reactive systems. It is based on a precise mathematical and unambiguous model of the possible system behaviour. To verify that the model meets specified requirements (usually expressed in temporal logic), all possible executions are checked systematically. If an execution failing the specification is found, the execution path which caused the violation is returned. Model checking has previously been successfully used in a variety of different systems. It helped to ensure the safety and reliability of safety-critical systems like flight control [31], space-craft controllers [13] and satellite positioning systems [24]. It has also been successfully applied to many aspects of software verification, e.g., for industrial systems and operating systems [2, 30].

In autonomous robotic systems, model checking has been used for static and runtime verification [23, 10, 5, 33] and has been proposed for strategy synthesis [11, 17]. It has been used in many contexts, for example: to generate real-time action plans with formal guarantees for steerable needles [21]); industrial robots [32]; and for assistive-care robots which dynamically re-calibrate their path in real-time [12].

In recent work [26], we combined the Spin [14] model checker’s ability to identify paths violating temporal properties with sensor information from a 3D Unity simulation of an autonomous vehicle, to plan and perform consecutive overtaking manoeuvres on a traffic-heavy road. The model checker received information from the (simulated) autonomous vehicle, updated its current model, derived a safe path, then communicated the path back to the autonomous vehicle. Although a useful proof-of-concept, the time delay due to model compilation (approx. 3 secs—even though verification of the model itself only took around 20 ms) and the communication between the model checker and game engine was unacceptable. In addition, the requirement to divide the underlying action space into discrete sections made the approach feasible only for less congested environments, such as rural roads.

In this paper we investigate adapting the approach of [26] to a real autonomous agent. Our main objective is to demonstrate the effectiveness of two measures which should address the time lag and accuracy problems described above: (i) a simplified model checking algorithm with faster compilation time, and (ii) the use of model checking situated *on board* the autonomous agent.

In Section 2 we give an overview of our method. Specifically we show how a robot uses multi-step plans derived using model checking to move through a domain, avoiding immediate danger. We present the underlying formal model and describe both how it is used to generate *solution paths* to eliminate disturbances and how the model is updated in real-time. In Sections 3 and 4 we describe our implementation and present results. In Section 5 we discuss the implications of our approach—how it compares with previous work, an alternative physics modelling approach, and its current limitations.

## 2 Method

A typical scenario in which a wheeled robot is driving in an environment avoiding walls and obstacles is shown in Figure 1 A. The robot agent is initially in a resting state executing a preferred task  $T_0$  which makes the robot drive in a straight line. However, from sensory input, the robot can infer that if it continues to follow this path, it will crash into a wall, an unwanted and unexpected event. The robot can predict what is going to happen as it knows its own direction and velocity and can reason about possible courses of action. For example, it can switch to task  $T_R$  (turn right) which is executed until the disturbance  $D_1$  has been avoided. The task  $T_R$  can be viewed as a temporary control system (see Figure 1B) with the goal of counteracting the disturbance until it has been avoided, turning states (i.e., sensor data) into actions. Once the disturbance is eliminated,  $T_R$  is not needed and the robot returns to task  $T_0$ .

While the robot has turned right in Figure 1A, the control goal of counteracting the disturbance  $D_1$  could have equally been achieved by turning left, executing task  $T_L$ . However, as shown in Figure 1C, the robot would have soon encountered a second disturbance  $D_2$ , a situation which cannot be solved by the temporary control system in Figure 1B as it can only reason one step ahead. In Figure 1C, if the robot turns left, it can go straight for a while but then has to turn left or right again. Furthermore, if the robot then turns right, it is possible that it might get trapped in a corner, which is undesirable behaviour. To overcome this limitation, the robot needs some ability to reason about the outcome of chained sequences of tasks for a given number of steps we call the *horizon*. This reasoning process translates into the tree structure of tasks and sensory inputs shown in Figure 1D which forms the foundation of our method.

Our basic approach is illustrated in Figure 2. In Figure 2A, a wheeled robot drives towards the far



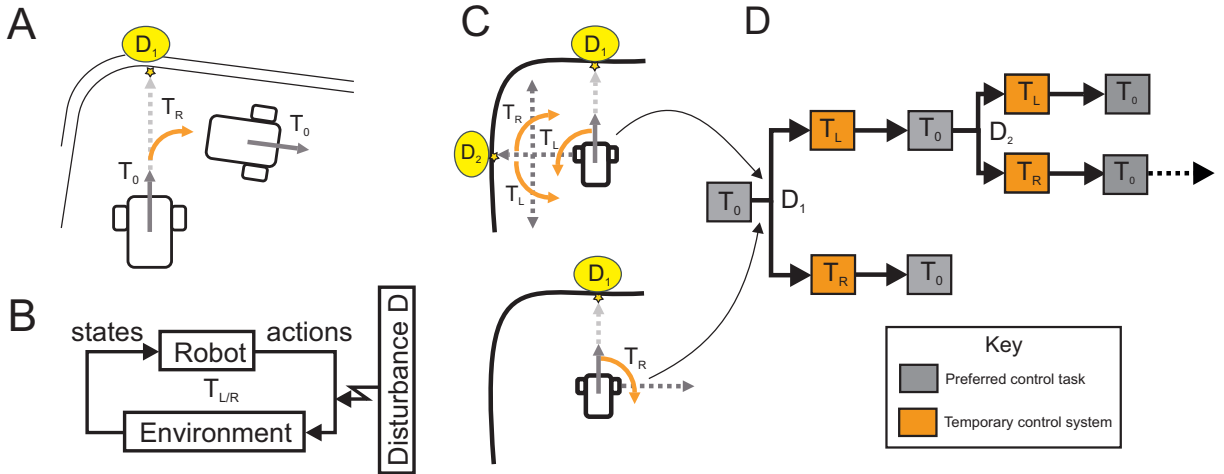


Figure 1: Overview of the general concept. In A the robot is executing its preferred task  $T_0$  which makes the robot drive in a straight line. The robot approaches a wall and senses a disturbance  $D_1$  which disrupts the robot agent from its preferred course. The robot spawns a control task  $T_R$  (turn right) to counteract  $D_1$  which can be viewed as the temporary control system shown in B. The task  $T_R$  only exists until the disturbance is has been counteracted. However, the disturbance  $D_1$  could have been eliminated by turning left  $T_L$ . C shows that this would have immediately put the robot in a complex situation which cannot be solved in a stimulus-response fashion by spawning a control task. In D a reasoning tree is shown for chaining the spawning of temporary control tasks in response to disturbances in the environment.

wall of a cul-de-sac in its preferred task  $T_0$  and senses a disturbance  $D$ , operationalised in our case as the nearest sensed point within a distance  $d = vt_{look}$  where  $v$  is the robot velocity and  $t_{look}$  is a set lookahead time. The width of the visual field is determined by the width of the robot plus some tolerance and a check is made for disturbances each iteration of the control loop (approx. every 200 ms). In this scenario, if the robot turns left or right to avoid the disturbance  $D$ , the robot will immediately encounter another disturbance. Furthermore, in either case the robot could get trapped in a corner which is undesirable.

To address this problem, our model is first updated with state information from the environment using a lightweight procedure explained in Section 2.2. In essence, the procedure involves checking whether a proper subset of the robot workspace is empty (i.e., free of obstacles) based on a novel abstraction of the point cloud data. We utilise symmetry on the axes of a 2D vector space and perform simple filtering to determine whether a given subset is disturbance free in which case the corresponding horizon state is determined *safe*. Figure 2B shows a graphical representation of the procedure outcome, which in this case has updated the model to reflect that turning left twice then returning to  $T_0$  is the safe option.

A valid path for the model is generated utilising a bespoke implementation of model checking (see Figure 2C). We extract transitions from the path to recover the trajectory, which in our case forms a sequence of control tasks from the set  $Act = \{T_0, T_L, T_R\}$ .  $T_0$  denotes the preferred task of driving in a straight line, and both  $T_L$  and  $T_R$  are temporary control systems (see Figure 1B) which rotate the robot left or right by 90 degrees to eliminate a disturbance. Real execution is never exact, however error is permissible so long as the control goal of eliminating the disturbance  $D$  is achieved. The resulting plan is executed when the disturbance is some distance  $d < d_{safe}$  from the robot. Once the plan has been executed, the robot returns to  $T_0$  until it encounters another unwanted disturbance in the environment.

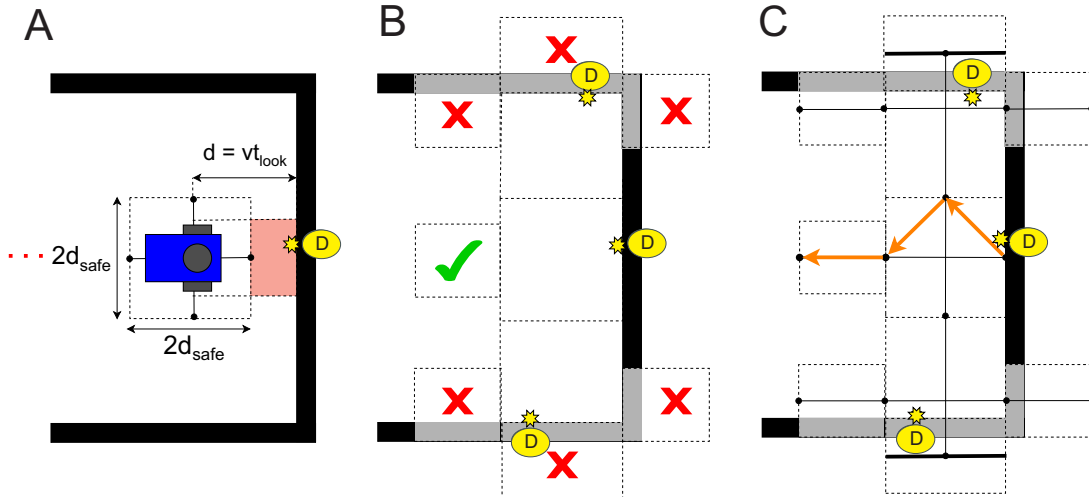


Figure 2: Representation of planning sequence for a cul-de-sac. In A the robot detects a disturbance  $D$  (indicated by yellow star), the nearest detectable point associated with the far wall. A distance  $d_{safe}$  defines a safe zone so the robot can rotate and acts as a threshold for the robot to start executing plans.  $d = vt_{look}$  is the sensing range derived from the velocity and a lookahead time  $t_{look}$ . B shows a graphical representation of the model update procedure. The robot cannot turn right/left then straight nor can it plan an extra step, the only empty set is the one behind. In C a path is generated using in situ model checking. Transitions are extracted and the resulting sequence of tasks (i.e., plan) is executed by control.

## 2.1 Preliminary model checking

### 2.1.1 Trajectory specification in LTL

We define trajectory specification as the desired sequence of discrete control tasks for an obstacle avoidance scenario, limited in this initial case to static environments. The planning problem normally consists of two conditions: (i) do not hit any obstacles and (ii) make progress towards a goal [27]. We address (i) by model checking of a regular safety property. Usually when model checking, counterexamples constitute violation of the property under scrutiny. We however denote any violation a *solution path*.

In this paper we focus on trajectory specification given as a Linear Temporal Logic (LTL) formula due to its power for discrete sequential planning. LTL formulae are built from a finite set of atomic propositions  $AP$ , Boolean connectors such as conjunction  $\wedge$  and negation  $\neg$ , and two temporal modalities  $\bigcirc$  (“next”) and  $\bigcup$  (“until”) [1]. The atomic proposition  $a \in AP$  stands for the state label  $a$  in a transition system. For example, in this initial work  $AP = \{safe, horizon\}$  where  $horizon \in AP$  is true in states of the transition system defined as valid planning steps and  $safe \in AP$  is true in states that can be reached without encountering a disturbance. The states where  $horizon$  is true are known a priori and therefore fixed for our transition system while the states where  $safe$  is true is decided at runtime. A complete description of our transition system is provided in Section 2.1.2.

LTL formulae over the set  $AP$  of atomic propositions are formed according to the following grammar [1]:  $\varphi ::= \text{true} \mid a \mid \varphi_1 \wedge \varphi_2 \mid \neg\varphi \mid \bigcirc\varphi \mid \varphi_1 \bigcup \varphi_2$  where a  $\varphi$  with an index denotes some arbitrary but distinct formula in LTL. Hence in our case, the atomic propositions  $\varphi_1 = safe$  and  $\varphi_2 = horizon$  are both LTL formulae, so by the grammar the conjunction  $\varphi_1 \wedge \varphi_2$  is also a formula as is its negation  $\neg(\varphi_1 \wedge \varphi_2)$ . From this basic grammar, other operators can be derived, such as  $\square$  (“always”) and  $\diamond$  (“eventually”), however the derivation is omitted here for brevity (see [1] for a detailed discussion).

Our approach utilises a single regular safety property  $\Box\varphi$  where

$$\varphi = \neg(\text{safe} \wedge \text{horizon}) \quad (1)$$

is an invariant expected to hold in each state of the system. Intuitively, the property  $\Box\varphi$  says that for any state  $s$  in the transition system at least one  $a \in AP$  is always false. As our interest is in *solution paths* not error paths, the invariant  $\varphi$  negates the desired outcome, so that what would normally be the set of counterexamples for an infinite run of the system, becomes a set of solutions. Consequently, the set of solutions paths for our model is the set of paths with a state satisfying the negation of the invariant:

$$\neg\varphi = \text{safe} \wedge \text{horizon} \quad (2)$$

The set of counterexamples of a regular safety property constitute a language of finite words which can be recognised by a nondeterministic finite automaton (NFA) [1]. We therefore construct the NFA  $\mathcal{A}_{\Box\varphi} = (Q, \Sigma, \delta, Q_0, F)$  where  $Q$  is a finite set of states,  $\Sigma = 2^{AP}$  is a finite alphabet defined as the power set of the  $AP$ ,  $\delta : Q \rightarrow 2^Q$  is a transition relation,  $Q_0 \subseteq Q$  is a set of initial states, and  $F \subseteq Q$  is a set of accept states [1]. In fact, for any invariant  $\varphi$ , the language of all counterexamples (i.e., solutions) can be represented by an NFA with two states. In our specific case, the NFA  $\mathcal{A}_{\Box\varphi}$  progresses to the accepting state and terminates if and only if for some state in the transition system the conjunction in (2) is true.

### 2.1.2 Task-driven transition system

A finite transition system is used as a model to describe the behaviour of the robot and provides semantics for trajectory specification in LTL. The discretized workspace consists of  $n$  states  $S = \{s_0, s_1, \dots, s_n\}$  and control tasks are interpreted as labelled state transitions to reflect the reasoning tree in Figure 1D.

**Definition 2.1** (Finite transition system). A finite transition system  $TS$  is a tuple  $(S, Act, \rightarrow, I, AP, L)$  where

- $S$  is a finite set of states,
- $Act$  is a finite set of actions,
- $\rightarrow \subseteq S \times Act \times S$  is a transition relation,
- $I \subseteq S$  is a set of initial states,
- $AP$  is a finite set of atomic propositions, and
- $L : S \rightarrow 2^{AP}$  is a labelling function.

The labelling function  $L$  relates a set  $L(s) \in 2^{AP}$  of atomic propositions to a state  $s$ , where  $2^{AP}$  denotes the powerset of  $AP$  (i.e., the set of all  $AP$  subsets including itself and the empty set) [1]. Hence the labelling function  $L$  assigns truth-values, determining which atomic propositions are satisfied for some state  $s$ .

The intuitive behaviour of a finite transition system is as follows. The finite transition system  $TS$  starts in some initial state  $s_0 \in I$  and evolves according to the transition relation  $\rightarrow$ . For convenience, we represent a transition between states with  $s \xrightarrow{\alpha} s'$  instead of  $(s, \alpha, s') \in \rightarrow$ . If  $s$  is the current state, then a transition  $s \xrightarrow{\alpha} s'$  originating from  $s$  is selected and executed, i.e., the action  $\alpha$  associated with the transition is performed, evolving the transition system from state  $s$  to  $s'$ . In cases where the current state  $s$  has more than one outgoing transition, the action  $\alpha$  is chosen in a nondeterministic fashion.

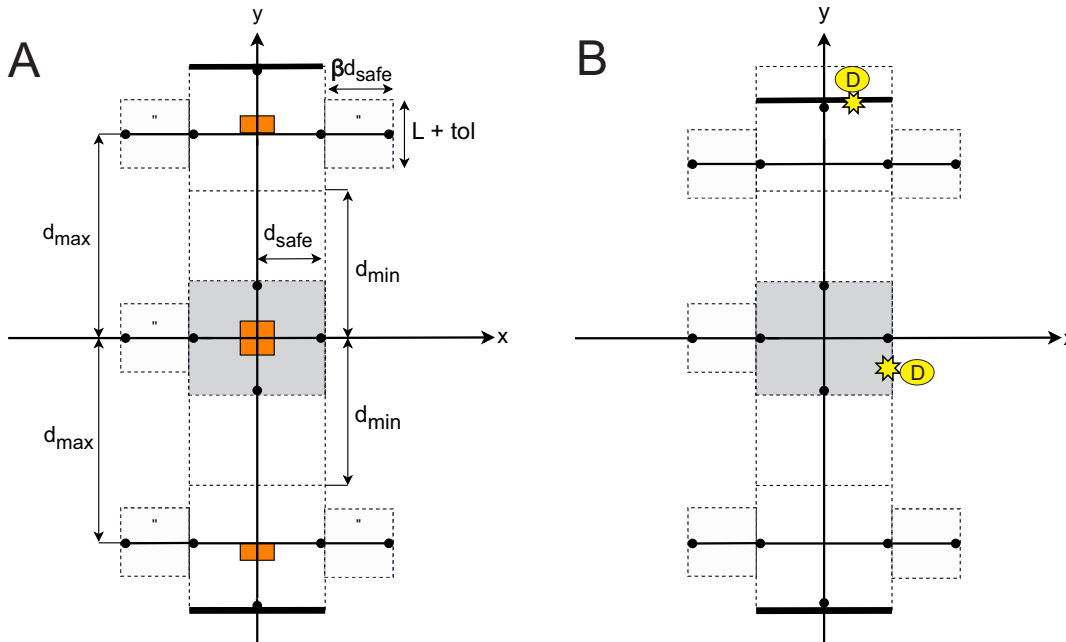


Figure 3: Depiction of point cloud abstraction and states. Here A shows the structure of the abstraction where  $d_{max}$ ,  $d_{min}$  and  $d_{safe}$  are tunable parameters constrained by the dimensions of the robot.  $L + tol$  represents the wheelbase of the robot plus some tolerance, which does not need to be as wide as the main driving corridor because the robot can only go straight. The states (black dots) represent a fixed point distance  $d_{safe}$  in front of the robot for future robot configurations as it navigates the abstraction, facing towards positive  $x$  by convention with the initial state  $s_0$  directly in front. In B an example scenario shows five states in the positive lateral direction adjusting within the tolerance  $d_{max} - d_{min}$  to the location of the nearest disturbance. Hence lateral states adapt to bounded stochastic variations in the local environment.

Our concrete model of robot behaviour is based upon abstraction of point cloud data from a 2D LiDAR with 360 degrees field of view (see Figure 3A for an illustration). States  $S = \{s_0, s_1, \dots, s_{14}\}$  represent points on a 2D vector space distance  $d_{safe}$  in front of possible future robot configurations; the robot faces towards positive  $x$  by convention when entering the initial state  $s_0$  to generate a plan. Hence path generation can be seen as the robot reasoning about where it will end up after executing a given sequence of tasks whilst respecting the robot safe zone. As shown in Figure 3A, the four states on the edge of the safe zone (grey box) represent a fixed point distance  $d_{safe}$  from the origin, defined to ensure the robot always has adequate clearance for rotational movements when initiating a plan. In addition, the furthestmost state behind the robot on the  $x$ -axis is also a fixed distance from the origin; this state provides the robot with an option to turn around and go back the way it came if relevant for the scenario. The  $y$ -coordinate of states on the lateral extremes of our discretization is variable, however, determined at runtime by the location of disturbances in the positive and negative lateral directions (see Figure 3B for an illustration of the abstraction adjusting to a disturbance in the positive lateral direction).

As mentioned above, control tasks are interpreted as labelled transitions between states to reflect the desired reasoning tree in Figure 1D. Hence we call our model a *task-driven* finite transition system where

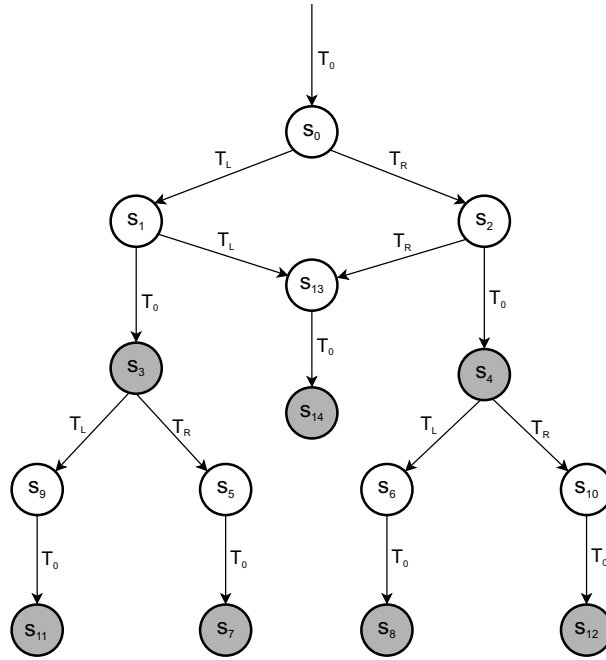


Figure 4: Task-driven finite transition system. Grey states indicate that *horizon* is true. These states are known a priori and fixed for the model.  $S_{left} = \{s_3, s_5, s_7, s_9, s_{11}\}$  and  $S_{right} = \{s_4, s_6, s_8, s_{10}, s_{12}\}$  are determined by the locations of disturbances in the lateral direction. The remaining states represent fixed points. Transitions reflect the reasoning tree shown in Figure 1D for the spawning of control tasks.

- $S = \{s_0, s_1, \dots, s_{14}\}$  is a set of states representing a point distance  $d_{safe}$  in front of the robot for possible future configurations as it navigates the abstraction shown in Figure 3,
- $Act = \{T_0, T_L, T_R\}$  is the set of discrete control tasks defined in Section 2,
- $\rightarrow \subseteq S \times Act \times S$  is a transition relation where  $s \rightarrow s'$  is admissible if and only if there exists a control task which can evolve the model from state  $s$  to  $s'$ ,
- $I = \{s_0\}$  is the initial state distance  $d_{safe}$  in front of the robot,
- $AP = \{safe, horizon\}$  is the set of atomic propositions defined in Section 2.1.1, and
- $L : S \rightarrow 2^{AP}$  is a labelling function such that  $L(s)$  determines if property (2) is true at state  $s$ .

Figure 4 shows the structure of our transition system. In our model, the states where *horizon* is true are known a priori and therefore fixed, indicated in Figure 4 by grey states. This is a modelling choice to reflect the final step of any plan, which is to return the robot to its resting state, the preferred task  $T_0$ . The transitions to horizon states provide some assurance that the robot has time to replan if another disturbance is encountered soon after plan execution, e.g., due to an error in the real task execution.

Our model generates a plan of one, two or three steps, as any transition to a safe horizon state (i.e., a state where the property in (2) is true) represents a return to the preferred task  $T_0$  and is thus excluded. One step plans reflect scenarios where the robot can spawn a single temporary control system  $T_{L/R}$  to counteract a disturbance then return immediately to the preferred task  $T_0$ ; two step plans are reserved for scenarios where the robot is boxed in and needs to spawn a sequence of two  $T_{L/R}$  to about turn and evade the situation; and three step plans add an extra  $T_{L/R}$  to counteract disturbances in the lateral direction.

### 2.1.3 Product transition system and NFA

In model checking, we are normally interested in establishing whether  $\Box\varphi$  is true for all possible runs of a system. This is equivalent to checking if  $Traces_{fin}(TS) \cap \mathcal{A}_{\Box\varphi} = \emptyset$  where  $Traces_{fin}(TS)$  is the set of finite traces for the transition system [1]. To check this, we first construct the product transition system  $TS \otimes \mathcal{A}_{\Box\varphi}$ , then derive an invariant  $\varphi$  for the product from the accept states of  $\mathcal{A}_{\Box\varphi}$  such that  $Traces_{fin}(TS) \cap \mathcal{A}_{\Box\varphi} = \emptyset$  if and only if  $TS \otimes \mathcal{A}_{\Box\varphi} \models \Box\varphi$  (i.e., the property  $\Box\varphi$  is satisfied in the product transition system). Verification of a regular safety property can therefore be reduced to invariant checking on the product. However, as we are interested in counterexamples of the property as solution paths, our focus is instead on generating paths in the system for states where the property  $\Box\varphi$  is false.

**Definition 2.2** (Product transition system). The product transition system  $TS \otimes \mathcal{A}_{\Box\varphi}$  is a tuple  $(S', Act', \rightarrow', I', AP', L')$  where

- $S' = S \times Q$ .  $s' = \langle s, q \rangle \in S'$ ,  $\forall s \in S$  and  $\forall q \in Q$ ,
- $\rightarrow'$  is the smallest relation defined by the rule  $\frac{s_i \xrightarrow{\alpha} s_j \wedge p \xrightarrow{L(s_j)} q}{\langle s_i, p \rangle \xrightarrow{\alpha'} \langle s_j, q \rangle}$ ,
- $I' = \{ \langle s_0, q \rangle \mid s_0 \in I \wedge \exists q_0 \in Q_0. q_0 \xrightarrow{L(s_0)} q \}$ .
- $AP' = Q$ , and
- $L' : S \times Q \rightarrow 2^Q$  is given by  $L'(\langle s, q \rangle) = \{q\}$ .

It suffices to perform a reachability analysis on  $TS \otimes \mathcal{A}_{\Box\varphi}$  to check the invariant  $\varphi$ . In this paper, we implement and perform invariant checking by forward depth-first-search (f-DFS) (see Algorithm 4 in [1] for details). Our model is updated with state information from sensors at runtime to determine which horizon states are safe and generate a sequence of states from which the associated control tasks can be extracted. If  $\Box\varphi$  is false for some state, then the execution path in which the state is reached (normally referred to as a counterexample) is called a *solution path* for the task-driven finite transition system, which in our case reflects a sequence of discrete robot configurations in the workspace. We then extract the tasks associated with each transition to recover the trajectory for the control layer to execute.

## 2.2 Model update procedure

As mentioned above, states where *horizon* is true are known a priori (see Figure 4 for details), so for a solution path to be generated it remains for us to decide which of these states is also safe. We utilise longitudinal and lateral offsets of the point cloud to simulate respective displacements and take advantage of symmetry on the axes of a 2D vector space to determine if a subset of our abstraction, specified to represent a practical over-approximation of the task execution workspace, contains no disturbance. If a given subset is empty, the corresponding state is reachable and determined safe, indicated in Figure 5.

Our abstraction in Figure 3 assumes the robot is at the origin of an underlying 2D vector space used as a model for the point cloud data. The robot is facing towards positive  $x$  by convention with the initial state  $s_0$  distance  $d_{safe}$  directly in front of the robot on the  $x$ -axis. In Figure 5, subset  $o_1$  is a reflection of  $o_2$  with the  $x$ -axis forming a line of symmetry (grey box indicates the robot safe zone, however both sets extend to the  $x$ -axis). For each of these subsets, the  $y$ -axis also forms a line of symmetry which splits each subset in half, such that each half is a reflection of the other. Symmetry on the  $x$ -axis means we can use a simple inequality to decide whether an observation has a qualifying  $y$ -coordinate (within the lateral bounds of the abstraction, positive for set  $o_1$  and negative for set  $o_2$ ). Symmetry on the  $y$ -axis means we

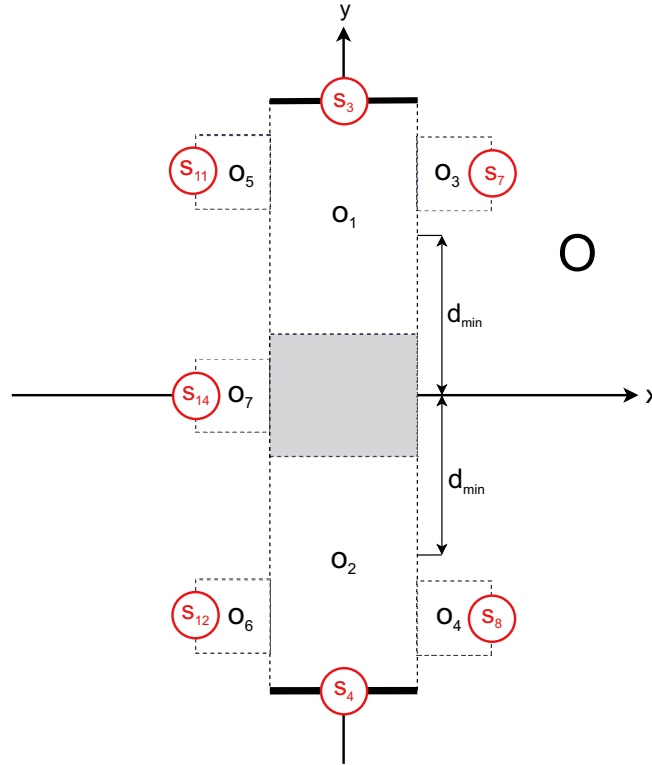


Figure 5: The set of point cloud observations  $O$  and the disjoint subsets  $o_1, o_2, \dots, o_7$  which form the abstraction. If some  $o_i = \emptyset$ , the atomic proposition *safe* is true for the horizon state on the edge of the set (indicated with red circles), otherwise the proposition is false. In addition, if the robot cannot travel a lateral distance  $d_{min}$  in either direction, then it is assumed  $o_7 = \emptyset$  to address trap situations, e.g., getting stuck in a corner. In static environments this is valid as the direction the robot came from should be safe.

can use absolute values to ignore the sign and ensure that any included observation has an  $x$ -coordinate which does not exceed some maximum distance from the  $y$ -axis. For sets  $o_1$  and  $o_2$ , this distance is  $d_{safe}$ , hence the longitudinal dimension of these sets respects the robot safe zone.

When a disturbance  $D$  is sensed, the procedure is initiated. First a longitudinal offset is calculated by subtracting  $d_{safe}$  from the  $x$ -coordinate so we can forward simulate our abstraction to its location:

$$\Delta_x = D_x - d_{safe} \quad (3)$$

However, if  $D_x \leq d_{safe}$  then  $\Delta_x = 0$ , as the abstraction is already at the desired displacement from the disturbance in the longitudinal direction, so no forward simulation is necessary for reasoning.

We then iterate the observations  $O$ , forward simulate each observation by subtracting  $\Delta_x$  from the  $x$ -coordinate, and use symmetry on the axes of the vector space to sort observations into relevant subsets:

$$o_1 = \{o \in O \mid 0 < o_y < d_{max} + d_{safe} \wedge |o_x| \leq d_{safe}\} \quad (4)$$

$$o_2 = \{o \in O \mid -(d_{max} + d_{safe}) < o_y < 0 \wedge |o_x| \leq d_{safe}\} \quad (5)$$

where  $0$  is a constant to distinguish between positive and negative  $y$ -coordinates, the expression  $d_{max} + d_{safe}$  represents the absolute distance to lateral extremes of our abstraction, and  $|o_x| \leq d_{safe}$  ensures the

width of the subset respects the robot safe zone. Subsequent reasoning can be seen as the robot predicting what will happen if it executes a one step plan to avoid the disturbance and return to its resting state.

If  $o_1 = \emptyset$  and  $o_2 = \emptyset$ , then in either case we can infer that the subset is free of disturbances, meaning that the robot can execute  $T_L$  or  $T_R$  then return to the preferred task  $T_0$  at least for distance  $d_{max}$  (a set parameter for the furthestmost possible lateral configuration of the robot from the origin of the vector space, as per the point cloud abstraction shown in Figure 3A). From the robot agent perspective, this means that both states  $s_3$  and  $s_4$  can be reached without encountering a disturbance (insofar as it knows) and a path is generated nondeterministically. If  $o_1 = \emptyset$  and  $o_2 \neq \emptyset$ , then the robot can only execute  $T_L$  before returning to  $T_0$ , so a path is generated for  $s_3$  not  $s_4$ . If  $o_1 \neq \emptyset$  and  $o_2 = \emptyset$ , then the robot can only execute  $T_R$  before returning to  $T_0$ , so a path is generated for  $s_4$  not  $s_3$ . Where any of these conditions hold, a one step plan is generated without progressing the procedure, as the robot would like to return to its resting state as soon as possible. However, if both  $o_1$  and  $o_2$  are not empty, we can conclude that the robot will soon encounter another disturbance once it returns to  $T_0$ , so a one step plan is not possible.

Next a two step plan is considered, i.e., whether the robot is boxed in and should turn 90 degrees twice to go back the way it came, or has enough room in the lateral directions to execute a three step plan. The nearest positive lateral disturbance  $D^+ = \min |o_y|$  for  $o \in o_1$  and nearest negative lateral disturbance  $D^- = \min |o_y|$  for  $o \in o_2$  is acquired. As long as  $|D_y^+|$  or  $|D_y^-|$  is greater than  $d_{min}$ , it is concluded that the robot can travel at least distance  $d_{min} - d_{safe}$  in the associated lateral direction. It is therefore considered a safe initial direction and a valid three step plan exists for the scenario. Otherwise, the robot infers that it is boxed in, so the only empty set is  $o_7$  (immediately behind the robot), leading to the generation of a two step plan in which the robot turns around to go back the way it came (assumed to be safe). In this case, a two step plan is generated and the procedure terminates, as further reasoning is unnecessary.

However, if at least one direction is determined initially safe for a three step plan, positive and negative lateral offsets are calculated so that we can reason about counteracting any lateral disturbances:

$$\Delta_y^+ = D_y^+ - d_{safe} \quad (6)$$

$$\Delta_y^- = D_y^- + d_{safe} \quad (7)$$

where  $\Delta_y^+$  estimates the maximum lateral displacement of the robot in the positive direction while respecting the robot safe zone, and  $\Delta_y^-$  represents the same for the negative direction. At this point in the procedure, we have no more use for sets  $o_1$  and  $o_2$  so they do not participate in any further reasoning. Instead we build sets  $\{o_3, o_5, o_4, o_6\}$  to reason about counteracting any lateral disturbances. Progressing this far means that we have already reasoned about the initial two steps the robot can execute.

Lateral offsets translate the point cloud data, such that the robot remains at the origin and the axes of the vector space again form lines of symmetry on the relevant subsets. For example, when  $\Delta_y^+$  is applied, the robot is at the origin distance  $d_{safe}$  from the nearest disturbance on the left of the robot, representing an egocentric perspective of its future location if it first executed the sequence  $\langle T_L, T_0 \rangle$  (the orientation of the robot is of course different, however for our purposes this can be ignored). As a result, states  $s_{11}$  and  $s_7$ , which are by design a fixed lateral distance  $d_{safe}$  from any sensed lateral disturbances, sit on the  $x$ -axis, such that  $s_{11}$  is directly behind the robot and  $s_7$  is directly in front. The states are a midpoint in the lateral dimension for the associated subsets, hence  $o_5$  is behind the robot,  $o_3$  is in front, and the  $x$ -axis forms a line of symmetry which splits each subset in half. This means we can use absolute values to ignore the sign and ensure that any included observation has a  $y$ -coordinate that does not exceed some maximum distance from the  $x$ -axis. Symmetry on the  $y$ -axis means we can use a simple bounded inequality to decide whether an observation has a qualifying  $x$ -coordinate (whilst excluding any of subsets  $o_1$  and  $o_2$ ).



We therefore iterate  $O$ , forward simulating each observation subtracting  $\Delta_x$  from the  $x$ -coordinate as before, but adjusting for the lateral displacement by subtracting  $\Delta_y^+$  or  $\Delta_y^-$  from the  $y$ -coordinate.  $O$  is iterated once for each direction. In either case, the lateral offset places the relevant three step horizon states on the  $x$ -axis so we can use symmetry on axes of the vector space for specifying relevant subsets:

$$o_3, o_4 = \{o \in O \mid d_{safe} < o_x \leq \beta d_{safe} \wedge |o_y| \leq \frac{1}{2}(L + tol)\} \quad (8)$$

$$o_5, o_6 = \{o \in O \mid -\beta d_{safe} \leq o_x < -d_{safe} \wedge |o_y| \leq \frac{1}{2}(L + tol)\} \quad (9)$$

where  $d_{safe}$  is a constant to distinguish between positive and negative  $x$ -coordinates (while excluding subsets  $o_1$  and  $o_2$ ),  $\beta$  is a coefficient for tuning the length of the subsets and  $\frac{1}{2}(L + tol)$  defines the width (see Figure 3A for an illustration). From a cognition perspective, subsequent reasoning can be seen as the robot predicting what will happen if it returns to its resting state after eliminating a second disturbance.

Suppose the positive and negative lateral directions have both been calculated as safe for a three step plan, i.e., both  $|D_y^+|$  and  $|D_y^-|$  are greater than  $d_{min}$ . If subsets  $\{o_3, o_5, o_4, o_6\}$  are empty, then states  $\{s_7, s_{11}, s_8, s_{12}\}$  are deemed safe and a path is generated for one of the states in a nondeterministic way. However, if any of the subsets are non-empty, the corresponding state is unsafe and excluded from path generation via model checking. If  $|D_y^+| \leq d_{min}$ , then states  $s_7$  and  $s_{11}$  are automatically considered unsafe (i.e., the positive lateral direction is invalid for a three step plan). If  $|D_y^-| \leq d_{min}$ , then states  $s_8$  and  $s_{12}$  are automatically considered unsafe (i.e., the negative lateral direction is invalid for a three step plan). As mentioned above, prior to generating a three step plan, if  $|D_y^+|$  and  $|D_y^-|$  are less than or equal to  $d_{min}$ , neither direction is considered safe, so a two step plan to turn around and evade the situation is generated.

### 3 Implementation

As a case study, we implemented our method on a differential drive robot shown in Figure 6A. Our robot was adapted from a widely available mobile robot development platform, AlphaBot by Waveshare<sup>1</sup>. For sensing the environment, we equipped the robot with a low cost 360 degree 2D laser scanner, RPLiDAR A1M8 by Slamtec<sup>2</sup>, and for actuation we used two continuous rotation servos by Parallax<sup>3</sup>. The hardware programming interface for the robot was a Raspberry Pi 3 Model B<sup>4</sup> included with the AlphaBot development kit running a Quad Core 1.2GHz Broadcom 64bit CPU with 1GB RAM and wireless LAN.

Our method was implemented<sup>5</sup> in C++ using the closed-loop agent architecture shown in Figure 6B. At runtime, LiDAR scans generate a callback from a dedicated thread which passes observations to the agent event handler (approx. every 200 ms). The agent then initiates the task execution step which sends a control signal to the actuator thread. If the robot is in the preferred task  $T_0$  after the control signal is sent to the servos, a check is made for new disturbances in the environment. If a temporary control system  $T_{L/R}$  is the current task, the progress of counteracting the disturbance is checked. In either case, a task result is returned to the agent indicating whether the task has been a success or has failed. When in the preferred task  $T_0$ , if there is no plan available, a new disturbance initiates the model update procedure described in the previous section and generates a plan using invariant checking by f-DFS. The plan is

<sup>1</sup><https://www.waveshare.com/alphabot-robot.htm>

<sup>2</sup><https://www.slamtec.com/en/LiDAR/A1/>

<sup>3</sup><https://www.parallax.com/product/parallax-continuous-rotation-servo/>

<sup>4</sup><https://www.raspberrypi.com/products/raspberry-pi-3-model-b/>

<sup>5</sup><https://github.com/possibilia/mc-avoid>

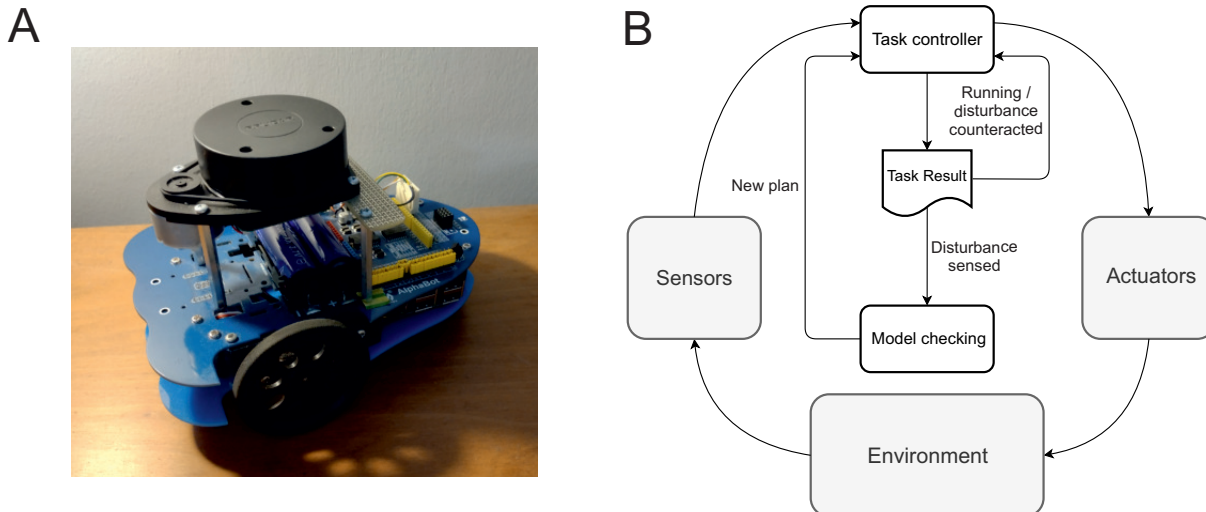


Figure 6: A our robot. B the agent architecture.

executed when the disturbance is distance  $d < d_{safe}$  from the origin of the point cloud data. Once the plan has been executed, the robot defaults to the preferred task  $T_0$  until a disturbance repeats the process.

## 4 Results

Our goal was to develop an egocentric method for chaining temporary control systems in response to disturbances in the environment using in situ model checking. Specifically, we were interested in improving on the case where an agent can only spawn a single control task in response to disturbances (see Section 2 for details). In this paper, we restricted our attention to static environments and local tactical planning for avoiding obstacles, focusing on a cul-de-sac scenario as an initial test case for our method. Our results show an improvement on one step planning yielding efficient trajectories for avoiding a cul-de-sac. In addition, for both comparisons our model checking procedure was executed in less than 11 ms.

Trajectories for the the first comparison are shown in Figure 7. Using our method, the robot approaches the cul-de-sac in Figure 7A and upon sensing a disturbance in the environment (i.e., a point on the wall on the left) generates the three step plan  $\langle T_R, T_0, T_R \rangle$  in Figure 7B. Here the model update procedure determines that sets  $o_1$  and  $o_2$  (see Figure 5 above) are non-empty and that it can drive at least  $d_{min} = 0.5\text{m}$  in either direction. However, as sets  $o_3$  and  $o_5$  are also non-empty, it infers that it will meet a disturbance if it turns either left or right after initially going in the left direction. In this particular case, both sets  $o_4$  and  $o_6$  are empty and so the corresponding horizon states  $s_8$  and  $s_{12}$  are determined safe. A path to one of the states is chosen nondeterministically, in this case  $s_{12}$ . Execution time for the combined model update procedure and path generation using in situ model checking was 9.22 ms. Figure 7C shows the one step comparison, which in this case manages to navigate out of the cul-de-sac by chance, however not without entering it first, making the trajectory less efficient.

The second comparison is shown in Figure 8. In this case, the robot approaches the bottom right corner of the cul-de-sac in Figure 8A and infers that it is boxed in using our method. Consequently, it generates the two step plan  $\langle T_L, T_L \rangle$  in Figure 8B then returns to the preferred task  $T_0$ . As in the previous comparison, the model update procedure has determined that sets  $o_1$  and  $o_2$  are non-empty, however in this case the robot cannot travel at least  $d_{min} = 0.5\text{m}$  in any direction, so a three step plan is invalid; it is

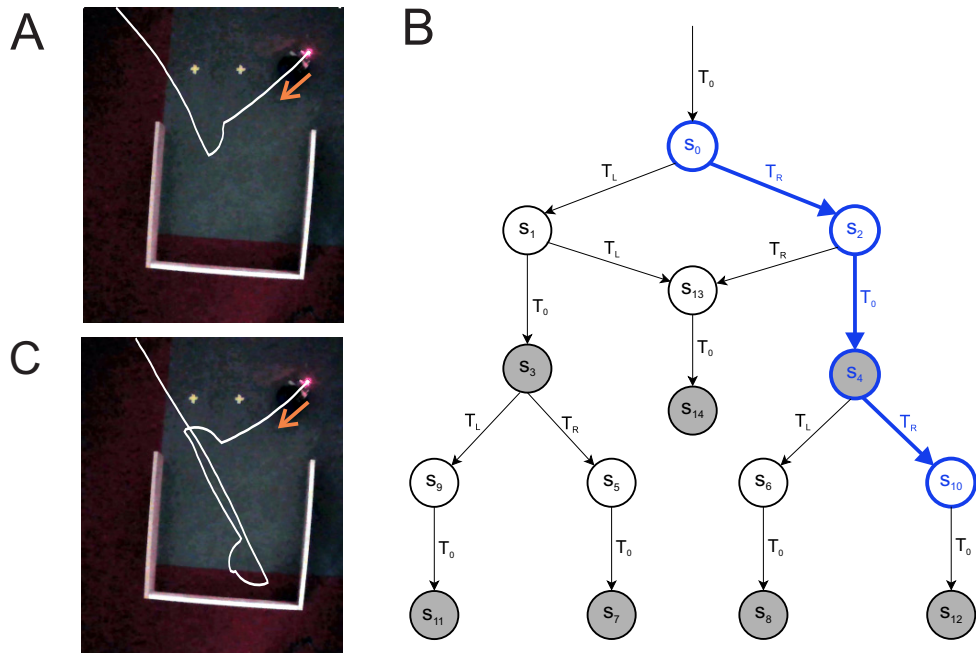


Figure 7: Comparison 1. In A the robot recognises that it would get boxed in turning left so it makes a three step plan for the right direction (last step cropped). In B the path and transitions in the model which generated the behaviour are shown. C shows trajectory followed by agent which can only plan one step.

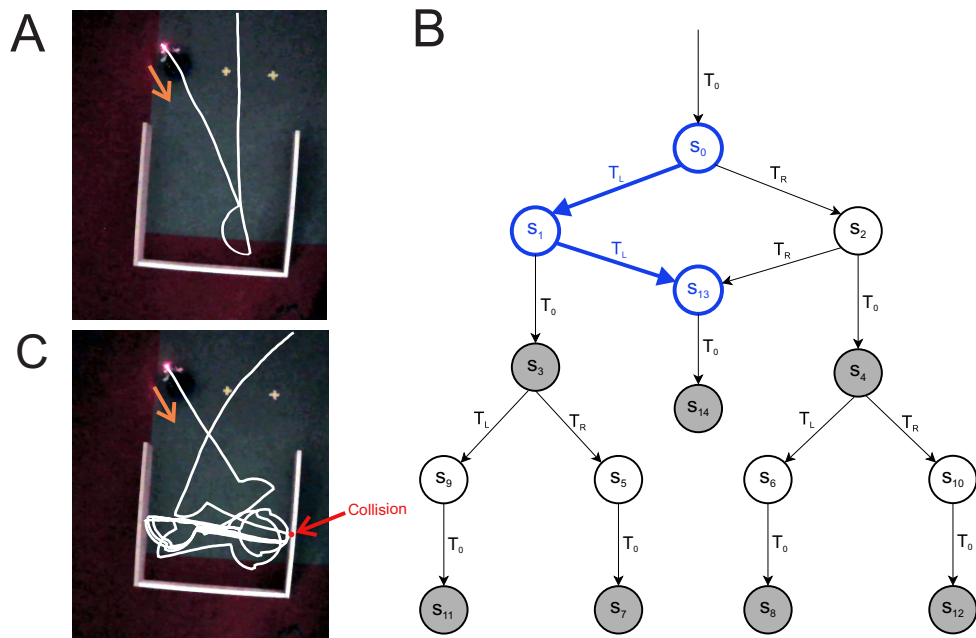


Figure 8: Comparison 2. In A the robot approaches the bottom right corner of the cul-de-sac and infers that it is boxed in so generates the two step plan shown in B to evade the situation. C shows that the robot gets trapped between two walls for the one step planning case, eventually colliding with a wall.

therefore assumed that only set  $o_7$  is empty and the corresponding horizon state  $s_{14}$  is safe. As there is only one safe horizon state in this situation, nondeterminism is resolved in the algorithm so the path is fully determined. Similar to the previous comparison, execution time for the combined update procedure and path generation was 10.87 ms. In contrast, the one step planning trial causes the robot to get trapped between two walls for a significant period of time before eventually colliding with one of them, as shown in Figure 8C. Hence our method generates a more efficient trajectory and can avoid a trap situation.

## 5 Discussion

In this paper, we have shown that it is possible to use live model checking to plan a safe sequence of discrete control tasks for a planning horizon of more than one step. In previous work [26], it was demonstrated that model checking could be used to plan overtaking manoeuvres for an autonomous vehicle (AV) as a proof of concept. As mentioned in the introduction, however, one of the main issues was compilation time in Spin [14] (approx. 3 secs), making the real-world application of model checking for trajectory planning impractical, even though verification of the model itself only took around 20 ms. We have overcome the compilation time problem by creating a stripped down model checking algorithm situated on board an autonomous agent. While our results are preliminary, we consider this work a successful first step towards real-time model checking for reliable and safe AV trajectory planning.

As finite state model checking relies on a discrete action space, another limitation in [26] was division of the underlying continuous system into 21 meter long segments, which is sufficient to represent rural roads or empty motorways but not congested urban environments. While using a fine-grained discretization would allow for more accurate modelling of the environment and vehicle speed, the state-space explosion problem imposes hard practical limits. If the state-space is too large, the additional computation would make real-time application of model checking infeasible, especially in high speed environments such as autonomous driving. In [22] it was argued that even a lag of 100 ms is unacceptable.

We have therefore introduced a novel discretization which moves away from static grids commonly applied in model checking, using abstraction to represent bounded stochastic variations in the continuous system. This has the benefit of keeping the state-space small but the model sensitive to fine-grained variations in the local environment. In this initial work, the upshot of our discretization is not obvious, however we believe that for any real-world application of model checking for AV trajectory planning, a discrete representation of the environment with adaptive characteristics will be necessary.

One limitation, however, is that our discretization of the LiDAR data models points on a 2D vector space. While sufficient for initial research, realistic driving scenarios will require richer information about the geometry of the local environment for reliable decision-making. If the dimensionality of objects in the environment is not known, then sophisticated manoeuvres will be impossible. For example, overtaking can be broken down into three distinct sub-manoevres [8]: (i) lane change to overtaking lane, (ii) pass leading vehicle(s), and (iii) lane change back to original lane, so for a successful overtake, the AV needs to know the necessary lateral and longitudinal displacements relative to the leading vehicle(s). Points on a plane have no dimensionality, so it is expected that some LiDAR preprocessing yielding richer geometry information about objects will be required for performant AV decision-making.

In a real-world driving scenario, AVs will need to cope with unpredictable traffic and changing weather conditions, so trajectory planning should in addition be sensitive to uncertainty in the driving environment. For this initial work we have chosen to keep the complexity of the model low to simplify the problem; one benefit is that the restricted modelling palette forces basic assumptions to be questioned in pursuit of a solution, as opposed to relying on the modelling apparatus alone. In our case, we have at-

tempted to stretch assumptions about relevant discretizations, moving away from typical fixed structures. In future work, however, a probabilistic model checker such as PRISM [16] could be stripped down and implemented on an autonomous agent to handle uncertainty in the local environment.

The main limitation of our approach is that it is not goal-directed. While our method is capable of local tactical planning for obstacle avoidance, any realistic scenario involving AVs (e.g., overtaking) requires the ability to make progress towards a goal relevant to the task. For example, in our previous work [26] a simulated LiDAR array provided sensory input sufficient for determining whether an overtake was possible. Subsequent research will seek to improve on our solution by adding goal-directed behaviour to the autonomous agent as a next step towards transparent and ecologically valid trajectory planning.

The major benefit of our approach is that it is tailored for individual scenarios and takes place in real-time using no pre-computed data [26]. Furthermore, it makes hard decisions which are in principle transparent. Popular machine learning approaches for trajectory planning, such as deep reinforcement learning [15], are either trained offline on datasets irrelevant to the immediate context, or trained online within a simulated environment. While predictions are fast in offline methods, they can result in unexpected or risky behaviour for unseen cases, which in autonomous driving can be dangerous, as evidenced in recent high profile accidents like the Tesla crash while in autopilot mode [20] and the Uber autonomous taxi crash [7], both of which were fatal. Online training of reinforcement learning in a simulated environment might be able to generate a larger variability of situations for training (and a much larger dataset than any real world data), but it would ultimately generate a rigid black box system which would not be transparent or guaranteed to react safely in all situations. Our solution relies on data from the local environment and generates explainable trajectories in real-time. Unlike machine learning methods, our basic approach makes clear decisions which are transparent by design.

Safety is ensured by our abstraction through over-approximation of the robot workspace and strict adherence to the robot safe zone. Similar concepts restricting the local behaviour of robots have been used elsewhere. In [18], for example, the notion of a *safe maneuvering zone* (SMZ) was used with a kinematic model for obstacle avoidance. The SMZ defines a circular boundary around the closest detected obstacle, creating a temporary sub-goal which minimally deforms the original robot path. The contour of any encountered obstacle is navigated until the robot can return to its original path. However, as the SMZ places safety bounds on obstacles, efficient strategies for obstacle avoidance may be ignored. Our safe zone is egocentric, so paths are not constrained by the contours of obstacles in the environment.

Previous work [3] combining lattice-based planning with optimal control has informed the development of our method. Here optimal path planning algorithms were used to generate motion primitives which can then be chained, producing locally optimal solutions to the path planning problem. However, it is a non-trivial task for a robot to precisely determine its position and follow a trajectory, often the approach favoured by safe navigation methods in the literature, such as control barrier functions [28]. Instead we focus attention on possible collisions which are tracked until out of reach whilst respecting the robot safe zone. Hence the trajectories produced by sequences of motion primitives (i.e., control tasks) are flexible in our method, as long as the control goal of eliminating disturbances is achieved.

## 5.1 Comparison with physics modelling

In a similar vein to model checking, we have approached closed-loop navigation from a lower level of abstraction; in brief, this consists of representing the robot, its behaviour and the external environment in the physics engine Box2D [6] (unpublished data). In this framework, we replace the model checking step by simulation of the execution of a task (or sequence of tasks) in the physics engine. The outcome of simulation can then be used to determine a sequence that best satisfies a goal. Modelling through a widely

utilised and validated physics library allows for accurate simulation of complex, dynamic environments. On the other hand, the speed of this method is strongly correlated with the number of bodies used in the simulation. In simple scenarios, optimal path selection is possible within the LiDAR sampling rate. However, cluttered environments and/or the creation of large task trees result in highly variable performance to a point where it no longer suits real-world applications in robots with low CPU clock speed (our data was collected on a 1.4GHz CPU). This limitation is not in the model checking approach, making it more robust, reliable and suited to a wider range of scenarios.

## **Acknowledgements**

This work was supported by a grant from the UKRI Strategic Priorities Fund to the UKRI Research Node on Trustworthy Autonomous Systems Governance and Regulation [EP/V026607/1, 2020-2024]; the UKRI Centre for Doctoral Training in Socially Intelligent Artificial Agents [EP/S02266X/1]; and the UKRI Engineering and Physical Sciences Research Council Doctoral Training Partnership award [EP/T517896/1-312561-05].

## References

- [1] Christel Baier & Joost-Pieter Katoen (2008): *Principles Of Model Checking*. 950, The MIT Press, Cambridge, Mass, doi:10.1093/comjnl/bxp025.  
Publication Title: MIT Press ISSN: 00155713.
- [2] Davide Basile, Alessandro Fantechi & Irene Rosadi (2021): *Formal Analysis of the UNISIG Safety Application Intermediate Sub-layer*. In Alberto Lluch Lafuente & Anastasia Mavridou, editors: *Formal Methods for Industrial Critical Systems*, Springer International Publishing, Cham, pp. 174–190, doi:10.1007/978-3-030-85248-1\_11.
- [3] Kristoffer Bergman, Oskar Ljungqvist & Daniel Axehill (2021): *Improved Path Planning by Tightly Combining Lattice-Based Path Planning and Optimal Control*. *IEEE Transactions on Intelligent Vehicles* 6(1), pp. 57–66, doi:10.1109/TIV.2020.2991951.
- [4] V. Braitenberg (1986): *Vehicles: Experiments in Synthetic Psychology*. MIT Press, Cambridge, Massachusetts. Available at [https://books.google.co.uk/books/about/Vehicles.html?id=7KkUAT\\_q\\_sQC&redir\\_esc=y](https://books.google.co.uk/books/about/Vehicles.html?id=7KkUAT_q_sQC&redir_esc=y).
- [5] R. C. Cardoso, G. Kourtis, L. A. Dennis, C. Dixon, M. Farrell, M. Fisher & M. Webster (2021): *A Review of Verification and Validation for Space Autonomous Systems*. *Current Robotics Reports* 2(3), pp. 273—283, doi:10.1007/s43154-021-00058-1.
- [6] E. Catto: *erincatto/box2d: Box2D is a 2D physics engine for games*. Available at <https://github.com/erincatto/box2d>.
- [7] K. Conger (2020): *Driver Charged in Uber’s Fatal 2018 Autonomous Car Crash*. *The New York Times*. Available at <https://www.nytimes.com/2020/09/15/technology/uber-autonomous-crash-driver-charged.html>
- [8] Shilp Dixit, Saber Fallah, Umberto Montanaro, Mehrdad Dianati, Alan Stevens, Francis Mccullough & Alexandros Mouzakitis (2018): *Trajectory planning and tracking for autonomous overtaking: State-of-the-art and future prospects*. *Annual Reviews in Control* 45, pp. 76–86, doi:10.1016/j.arcontrol.2018.02.001.
- [9] Marie Farrell & Matt Luckcuck, editors (2021): *Proceedings of the Third Workshop on Formal Methods for Autonomous Systems*. 348, Open Publishing Association, doi:10.4204/eptcs.348.
- [10] Angelo Ferrando, Louise A. Dennis, Davide Ancona, Michael Fisher & Viviana Mascardi (2018): *Verifying and Validating Autonomous Systems: Towards an Integrated Approach*. In Christian Colombo & Martin Leucker, editors: *Runtime Verification*, Springer International Publishing, Cham, pp. 263–281. doi:10.1007/978-3-030-03769-7\_15
- [11] D. Fraser, R. Giaquinta, Hoffmann, M. Ireland, A. Miller & G. Norman (2020): *Collaborative models for autonomous systems controller synthesis*. *Form Aspects of Computing* 32, pp. 157—186, doi:10.1109/TCST.2006.872519.
- [12] J. Hamilton, I. Stefanakos, R. Calinescu & J. Cámara (2022): *Towards Adaptive Planning of Assistive-care Robot Tasks*. In Luckcuck & Farrell [25], pp. 175–183, doi:10.4204/eptcs.371.
- [13] K. Havelund, M. Lowry & J. Penix (2001): *Formal Analysis of a Space-Craft Controller Using SPIN*. *Software Engineering, IEEE Transactions on* 27, pp. 749–765, doi:10.1109/32.940728.
- [14] G. Holzmann (2011): *The SPIN Model Checker: Primer and Reference Manual*, 1st edition. Addison-Wesley Professional.
- [15] S. Josef & A. Degani (2020): *Deep Reinforcement Learning for Safe Local Planning of a Ground Vehicle in Unknown Rough Terrain*. *IEEE Robotics and Automation Letters* 5(4), pp. 6748–6755, doi:10.1109/LRA.2020.3011912.
- [16] M. Kwiatkowska, G. Norman & D. Parker (2011): *PRISM 4.0: Verification of Probabilistic Real-Time Systems*. In Ganesh Gopalakrishnan & Shaz Qadeer, editors: *Computer Aided Verification*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 585–591, doi:10.1007/978-3-642-22110-1\_47.

- [17] M. Kwiatkowska, G. Norman & D. Parker (2022): *Probabilistic Model checking and autonomy*. *Annual review of control, robotics, and autonomous systems* 5(1), pp. 385–410. doi:10.1146/annurev-control-042820-010947
- [18] Lionel Lapierre & Rene Zapata (2012): *A guaranteed obstacle avoidance guidance system: The safe maneuvering zone*. *Autonomous Robots* 32(3), pp. 177–187, doi:10.1007/s10514-011-9269-5.
- [19] Yann LeCun (2022): *A Path Towards Autonomous Machine Intelligence — OpenReview*. Available at <https://openreview.net/forum?id=BZ5a1r-kVsf>.
- [20] David Lee (2016): *US opens investigation into Tesla after fatal crash*. Available at <https://www.bbc.co.uk/news/technology-36680043>.
- [21] S. Lehmann, A. Rogalla, M. Neidhardt, A. Schlaefer & S. Schupp (2021): *Online Strategy Synthesis for Safe and Optimized Control of Steerable Needles*. In Farrell & Luckcuck [9], pp. 128–135, doi:10.4204/EPTCS.348.9.
- [22] X. Li, Z. Sun, D. Cao, Z. He & Q. Zhu (2016): *Real-time trajectory planning for autonomous urban driving: Framework, algorithms, and verifications*. *IEEE/ASME Transactions on Mechatronics* 21(2), pp. 740–753, doi:10.1109/TMECH.2015.2493980.
- [23] Dennis Louise, Michael Fisher, Nicholas Lincoln, Alexei Lisitsa & Sandor Veres (2016): *Practical verification of decision-making in agent-based autonomous systems*. *Automated Software Engineering* 23(3), pp. 305–359, doi:10.1007/s10515-014-0168-9.
- [24] Y. Lu, A. Miller, C. Johnson, Z. Peng & T. Zhao (2014): *Availability analysis of satellite positioning systems for aviation using the Prism model checker*. In: *Proceedings of the 17th International Conference on Computational Science and Engineering (CSE 2014)*, pp. 704–713, doi:10.1109/CSE.2014.148.
- [25] Matt Luckcuck & Marie Farrell, editors (2022): *Proceedings of the Fourth International Workshop on Formal Methods for Autonomous Systems (FMAS) and Fourth International Workshop on Automated and verifiable Software sYstem DEvelopment (ASYDE)*. 371, Open Publishing Association, doi:10.4204/eptcs.371.
- [26] Daumantas Pagojus, Alice Miller, Bernd Porr & Ivaylo Valkov (2021): *Simulation and Model Checking for Close to Realtime Overtaking Planning*. In Farrell & Luckcuck [9], pp. 20–37, doi:10.4204/EPTCS.348.2.
- [27] Bruno Siciliano & Oussama Khatib, editors (2016): *Springer Handbook of Robotics*. Springer International Publishing, Cham, doi:10.1007/978-3-319-32552-1.
- [28] Andrew Singletary, Karl Klingebiel, Joseph Bourne, Andrew Browning, Phil Tokumaru & Aaron Ames (2020): *Comparative Analysis of Control Barrier Functions and Artificial Potential Fields for Obstacle Avoidance*. Available at <http://arxiv.org/abs/2010.09819>. ArXiv:2010.09819 [cs, eess].
- [29] Elizabeth S. Spelke & Katherine D. Kinzler (2007): *Core knowledge*. *Developmental Science* 10(1), pp. 89–96, doi:10.1111/J.1467-7687.2007.00569.X.
- [30] Güliz Tuncay, Soteris Demetriou, Karan Ganju & Carl A. Gunter (2018): *Resolving the Predicament of Android Custom Permissions*. In: *Network and Distributed System Security Symposium*, pp. 1–15, doi:10.14722/ndss.2018.23221.
- [31] L. Wang & F. Cai (2017): *Reliability analysis for flight control systems using probabilistic model checking*. *Proceedings of the IEEE International Conference on Software Engineering and Service Sciences, ICSESS 2017-Novem*, pp. 161–164, doi:10.1109/RAM.2017.7889773.
- [32] M. Weißmann, S. Bedenk, C. Buckl & A. Knoll (2011): *Model Checking Industrial Robot Systems*. In: *Model checking software (SPIN 2011)*, 6823, Springer Berlin Heidelberg, pp. 161–176, doi:10.1007/978-3-642-22306-8\_11.
- [33] Yi Yang & Tom Holvoet (2022): *Generating Safe Autonomous Decision-Making in ROS*. In Luckcuck & Farrell [25], pp. 184–192, doi:10.4204/eptcs.371.13.



# Online Reachability Analysis and Space Convexification for Autonomous Racing

Sergiy Bogomolov

Newcastle University,  
Newcastle upon Tyne, United Kingdom  
sergiy.bogomolov@ncl.ac.uk

Taylor T. Johnson

Vanderbilt University,  
Nashville, USA  
taylor.johnson@vanderbilt.edu

Diego Manzananas Lopez

Vanderbilt University,  
Nashville, USA  
diego.manzanas.lopez@vanderbilt.edu

Patrick Musau

Vanderbilt University,  
Nashville, USA  
patrick.musau@vanderbilt.edu

Paulius Stankaitis

Newcastle University,  
Newcastle upon Tyne, United Kingdom  
paulius.stankaitis@ncl.ac.uk

This paper presents an optimisation-based approach for an obstacle avoidance problem within an autonomous vehicle racing context. Our control regime leverages online reachability analysis and sensor data to compute the maximal safe traversable region that an agent can traverse within the environment. The idea is to first compute a non-convex safe region, which then can be convexified via a novel coupled separating hyperplane algorithm. This derived safe area is then used to formulate a nonlinear model-predictive control problem that seeks to find an optimal and safe driving trajectory. We evaluate the proposed approach through a series of diverse experiments and assess the runtime requirements of our proposed approach through an analysis of the effects of a set of varying optimisation objectives for generating these coupled hyperplanes.

## 1 Introduction

Over the last several years, autonomous racing has actively been pursued as a strategy to explore edge-case scenarios in autonomous driving [16]. Racing scenarios present unique challenges with respect to navigating high speeds and multi-agent interactions. In these contexts, vehicles must be able to operate at the edge of their operating envelopes in close proximity to static and dynamic obstacles. Several competitions have emerged over the last couple of years, such as the Indy Autonomous Challenge (IAC) [16], and the F1TENTH International Autonomous racing competition [34]. Although numerous racing strategies have been proposed over the last several years, head-to-head racing at high speeds remains a challenge. Unlike the time trials that are frequently used as qualification rounds in these competitions [33], head-to-head racing requires designing a regime to be able to predict the future trajectories that reflect the intentions of the other opponents and drive through the track as quickly as possible.

Within the autonomous racing space, one of the most popular frameworks for tackling the racing problem has been formulating and solving an optimisation problem that balances obstacle avoidance and travelling at high velocities [40, 20]. Specifically, the model-predictive control framework (MPC), which finds optimal control commands based on a model of the underlying system, while satisfying a set of constraints is the most widely used approach [18]. Although MPC approaches have enjoyed success

in these settings [40], one of the main limitations exhibited by many approaches is a lack of robust online risk assessment in often dynamic and uncertain environments, particularly around vehicle-to-vehicle interactions. While a lot of progress has been made in this area, collisions still occur due to misplaced estimations of the set of all possible trajectories that the vehicle could pursue [20]. Furthermore, as Katrakazas et al. note “exhaustively calculating and predicting the trajectories of other traffic participants at each epoch incurs a huge computational cost”. Currently, many existing approaches treat the vehicle as an isolated entity, and the behavioural models of other participants within the environment have not yet been widely incorporated into the MPC regime [20].

One of the ways that this challenge has been addressed has been through the use of reachability analysis approaches [1]. The idea is to compute the set of states that the other racing agents could occupy in the future, for a fixed time horizon, and plan trajectories for the ego vehicle that avoids this unsafe set [23, 26, 29]. This unsafe set allows for modelling the inherent uncertainty in the behaviour of other agents and for the synthesis of safe racing trajectories [1]. There are two main challenges that arise in these contexts. The first is that over long time horizons, reachability approaches will result in overly conservative behaviours as the set of avoidable states grows. The second is that reachability approaches are typically computationally challenging endeavours, thus leveraging them online is quite challenging. In light of these challenges, the following paper presents a model-predictive control framework leveraging real-time reachability for a 1/10 scale autonomous vehicle test-bed in a multi-agent racing setting modelled after the F1TENTH International Autonomous Racing Competition.

Finally, obtaining a solution to the MPC problem generally entails solving a convex optimisation problem, which guarantees convergence to a globally optimum solution. However, due to the presence of static and dynamic obstacles, the optimal control problem of obstacle avoidance is inherently a non-convex problem. Therefore, to solve this problem efficiently many approaches leverage state-space convexification. In the past, several state-space convexification approaches have been proposed, including region partitioning [28], computing separating hyperplanes [25, 31], and constructing approximations using stored data points [38] (further discussed in Section 2). In our framework, we propose a novel optimisation-based approach for convexifying non-convex state spaces by computing coupled separating hyperplanes. The coupling of separating hyperplanes makes it possible to compute optimal safe and convex regions. However, it comes at the cost of increased computation time. Therefore, in this paper, we investigate the feasibility (e.g., timing constraints) of computing coupled separating hyperplanes in a real-time autonomous racing scenario.

In summary, the contributions of this paper are: (1) we introduce a novel closed-loop model-predictive obstacle avoidance controller that integrates online reachability analysis and an optimisation-based state-space convexification approach, (2) we evaluate this approach across a diverse set of simulation experiments using the F1TENTH simulation platform. These experiments include varying the number of dynamic agents, the number of static and dynamic obstacles, and the racing environment. (3) We present a timing analysis of the state-space convexification approach. (4) Finally, we evaluate our approach against the well-known model-predictive contouring control approach, which has shown great success in obstacle avoidance tasks.

## 2 Related Work

Researchers have approached the obstacle avoidance problem from two major perspectives. The first strategy involved formulating and solving an optimisation problem. The second regime has typically involved a hierarchical decomposition of path planning and reference tracking. A variety of algorithms

such as artificial potential fields [47], genetic algorithms [45], rapidly-exploring random trees (RRT) [19], fuzzy logic algorithms [32], elastic band theory [12], and rolling window methods [46] have demonstrated success in numerous arenas. A key limitation of many path planning approaches is that they are incapable of respecting kinodynamic constraints, such as bounds on the acceleration, and often the trajectories must be passed to a low-level controller that utilises a higher fidelity dynamics model and respects control constraints [40]. Furthermore, in highly dynamic and uncertain environments, planners must be able to replan sufficiently fast to react appropriately to split-second environmental threats [19]. However, most planners typically do not replan sufficiently rapidly to ensure split-second reactivity to threats [23].

As mentioned previously, MPC approaches have demonstrated great success in generating optimal trajectories that respect kinodynamic constraints and recently researchers have combined these approaches with reachability analysis to generate provably collision-free paths [1, 5, 23, 35, 23]. Within this regime, [1, 35, 26] utilise forward reachability methods in order to eliminate areas of the state space that would result in collisions. While these methods are extremely effective, these approaches must be implemented carefully in order to ensure that the resulting trajectories do not result in overly conservative behaviours [23]. The alternative to these approaches is backward reachability approaches [5, 23] which utilise a target set representing a set of undesirable states, in order to design controllers that can guarantee dynamic and static obstacle avoidance with minimal intervention. However, these approaches are computationally demanding and typically the safety-ensuring control constraints, derived from these methods, are computed and cached offline before being incorporated into an MPC problem [23].

Beyond reachability methods, over the last several years, several space convexification approaches for the obstacle avoidance problem have been proposed. In [39] a feasible convex set for model-predictive control is obtained by computing two parallel time-varying hyperplanes on racetrack borders. However, the resulting hyperplanes do not consider static obstacles or dynamic agents. The works of Mercy et al. [31, 30] and Scholte et al. [41] utilise the concept of separating hyperplanes to compute hyperplanes which separate autonomous systems from convex obstacles. The paper [22] combines the model-predictive control and dynamic agent reachability analysis, and uses IRIS (Iterative Regional Inflation by Semi-definite programming) [11] for a state-space convexification. A similar approach has been introduced in [27] for motion planning. Finally, in [28] two (polar and convex) different types of convexification methods based on region partitioning for obstacle avoidance were proposed. Their convex partitioning regime utilises a convex partitioning algorithm [21] to compute the minimum number of convex regions that are needed to capture non-convex obstacles, whereas the polar partitioning approach derives a safe set by using a minimum number of triangles.

The state-space convexification approaches described above have two main limitations for the racing scenario: they generally aim to compute the largest convex region in the non-convex space (e.g., not necessarily in the travelling direction of the ego vehicle) or are not able to handle non-convex obstacles. In our approach, we also express the problem of computing separating hyperplanes as an optimisation problem, but we are interested in computing a correct set of separating hyperplanes that provide the largest safe convex region in the direction of the ego vehicle. Furthermore, our proposed approach is able to handle non-convex obstacles.

### 3 Preliminaries

#### 3.1 Model-Predictive Control

Let us suppose we have the following (1) discrete-time system where  $x \in X \subseteq \mathbb{R}^n$ ,  $u \in U \subseteq \mathbb{R}^m$  and  $t \in \mathbb{N}$ .

$$x_{t+1} = f(x_t, u_t) \quad (1)$$

The MPC problem can then be expressed as a finite horizon optimisation problem (1) where a cost function  $J$  is being minimised over a finite time horizon  $N$  subject to constraints (2.1 - 2.4).

$$J_{t \rightarrow t+N}(x_t) = \min_{u_0, \dots, u_{N-1}} p(x_{t+N}) + \sum_{k=t}^{t+N-1} q(x_k, u_k) \quad (2)$$

$$x_{k+1} = f(x_k, u_k), \forall k \in \{t, \dots, t+N-1\} \quad (2.1)$$

$$x_0 = x_s \quad (2.2)$$

$$x_k \in X, \forall k \in \{t, \dots, t+N-1\} \quad (2.3)$$

$$u_k \in U, \forall k \in \{t, \dots, t+N-1\} \quad (2.4)$$

The cost function  $J$  is made up of a stage cost function  $q$  and a terminal cost function  $p$  which determine the cost of being at the interim state  $x_k$  after applying an input  $u_k$ , and the cost of being at the final state  $x_{t+N}$ . The constraints (2.1 - 2.4) assert that the optimisation problem, given by equation (1), begins from an initial state  $x_s$  and that the interim state and control inputs must respect the constraint sets  $X$  and  $U$ .

If the dynamics and constraints can be formulated as linear expressions, then the MPC problem can be solved efficiently using standard convex optimisation techniques. However, if the dynamics or constraints are nonlinear, then the problem becomes a nonlinear optimisation problem that is much more computationally challenging to solve. On the other hand, allowing for nonlinear dynamics and constraints may permit one to track complex systems with a higher level of fidelity than using linear expressions. Thus, the computational cost must be evaluated against overall system performance [42].

#### 3.2 Reachability Analysis

Reachability analysis is a technique for computing the set of all reachable states of a dynamical system from a set of initial states. The reachable set of  $R_{t+1}$  can be defined formally as:

$$R_{t+1}(X_0) = f(X_0, U) \quad (3)$$

where  $X_0 \subseteq \mathbb{R}^n$  represents the set of initial states,  $U \subseteq \mathbb{R}^m$  represents the input set. More generally, reachability analysis methods aim to construct a *conservative* flowpipe (4) which encompasses all the possible reachable sets of a dynamical system over a time-horizon  $[0, T]$ . This can be formalised as follows (in practice the union is computed over a discretised interval):

$$R_{[0, T]}(X_0) = \bigcup_{t \in [0, T]} R_t(X_0) \quad (4)$$

Reachability analysis has been widely used in applications that range from the formal verification of systems to problems relating to the safe synthesis of complex systems [2]. The majority of reachability

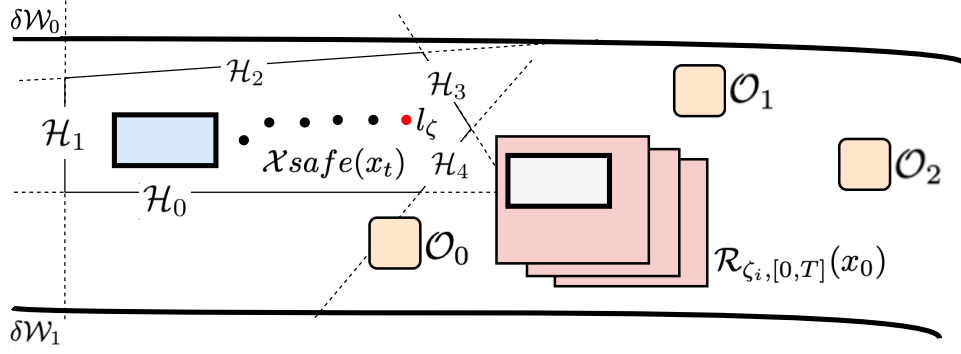


Figure 1: Visualisation of the autonomous racing problem with track boundaries,  $\{\delta W_0, \delta W_1\}$ , a dynamic opponent described by its reachable set  $R_{\zeta_i, [0, T]}(x_0)$  and static obstacles  $\{O_0, O_1, O_2\}$ . In this figure, the blue rectangle corresponds to the ego vehicle, and the white rectangle corresponds to a dynamic opponent. The main sub-problem is computing an n-number of separating hyperplanes ( $H_0 \dots H_4$ ) which jointly create a polyhedron  $X_{safe}$ . The computed  $X_{safe}$  must contain an ego vehicle and its target location  $l_\zeta$  as well as not overlap with observable obstacles.

analysis approaches leverage a combination of numerical analysis techniques, graph algorithms, and computational geometry [4, 3], and while in some cases it is possible to derive the exact reachable set of states, for many classes of systems computing the exact reachable set is infeasible. Thus, deriving the reachable set for these classes of systems involves obtaining a sound approximation (i.e., guarantee to contain a complete reachable set) of this set using a variety of set representations. Consequently, there is an inherent trade-off between the accuracy of the approximation and the time it takes to construct this set. We refer interested readers to the following papers [4, 3] for an in-depth discussion of these techniques.

## 4 Problem Statement and Space Convexification

### 4.1 Problem Formulation

In this paper, we consider the general autonomous racing problem (5) where a model predictive controller (2) is tasked with generating a sequence of control inputs  $u_{0 \dots T}$  that control a vehicle (1) such that it reaches the terminal state  $x_T \in X_f$  starting from an initial state  $x_s$  and steering through safe states, where  $X_{safe}$  and  $X_f$  are the safe states and terminal sets respectively. The goal is to steer the vehicle into the terminal set with the shortest time horizon  $T$ .

$$\begin{aligned}
 \min_{T, u_0, u_1, \dots, u_{T-1}} \quad & p(x_T) + \sum_{k=0}^T q(x_k, u_k) \quad s.t. \\
 x_{t+1} = & f(x_t, u_t), \quad x_0 = x_s \\
 x_t \in & X_{safe}, \quad u_t \in U, \quad x_T \in X_F
 \end{aligned} \tag{5}$$

In our formulation, the autonomous vehicle operates within a two-dimensional environment  $W \subset \mathbb{R}^2$  enclosed by boundaries  $\{\delta W_0, \delta W_1, \dots, \delta W_i\}$  as  $\delta W \subset W$ , among a set of dynamic agents  $\zeta = \{\zeta_0, \dots, \zeta_i\}$  ( $\zeta$  could be either ego  $\zeta_e$  or opponent vehicle  $\zeta_o$ ) and static obstacles  $\{O_0, O_1, \dots, O_i\}$  with  $O \subset W$ . The

region of space occupied by a dynamic agent  $\zeta_i(x_t) \in W$  in the environment over a time interval  $[t, t']$  from its current state  $x_t$  is given by its reachable set  $R_{\zeta_i, [t, t']}(x_t) \subset W$ . Our assumption is that the static and dynamic obstacles are contained within the two-dimensional environment. Furthermore, we refer to opponent vehicles within the racing environment as dynamic agents and refer to all other dynamic entities as dynamic obstacles.

To obtain a globally optimal solution to problem (5), as opposed to a locally optimal solution, the model-predictive control problem requires the state-space  $X$  to be convex. However, because of environment borders, static obstacles and dynamic agents,  $X$  is generally a non-convex entity. Therefore, the main sub-problem we are addressing in this paper is the computation of the safe, convex and *optimal* state-space  $X_{safe}$  in which a safe trajectory starting from  $x_0$  to a target location  $l_\zeta \in X_{safe}$  could be generated using an optimisation-based controller for the autonomous system (see Figure 1). The safe region of the state-space  $X_{safe}$  can be defined as follows:

$$X_{safe} = \{x \mid x \notin (\delta W \cup O \cup \bigcup_{i=0}^{N-1} R_{\zeta_i, [0, T]}(x_0))\} \quad (6)$$

where  $N$  is the number of observable dynamic agents. The computation of  $X_{safe}$  requires considering only *observable* obstacles, agents and borders. To define observable points we first introduce a notion of the LiDAR sensor which is mounted on the autonomous system and makes it possible to determine the distance to obstacles. The sensor sends  $M$  light pulses in an anti-clockwise direction around the autonomous system defined by  $\delta\theta$  increments and returns a set of observational points  $\{r_0(x_t), \dots, r_M(x_t)\}$  where a LiDAR observational point  $r_i(x_t) \in \mathbb{R}$  in the direction  $\theta_i$  can be formally defined in the following way (7):

$$r_i(x_t) = \min_{O_i \in O} \min_{z \in O_i} \|z - \zeta(x_t)\|_2 \quad s.t. \quad \text{atan2}(z - \zeta(x_t)) = \theta_i \quad (7)$$

Ranges of the observable LiDAR signals  $r_i(x_t)$  can be converted into a two-dimensional point cloud of the  $W$  where a single point  $p_i(x_t)$  of an agent  $\zeta(x_t)$  can be defined as a tuple (8):

$$p_i(x_t) = (\zeta(x_t) + r_i(x_t) \cos \theta_i, \zeta(x_t) + r_i(x_t) \sin \theta_i) \quad (8)$$

Now, we can define observable static obstacles of  $\zeta(x_t)$  as a set  $Q_{ob}$  of LiDAR points within a constant radius distance  $d$  from the agent's state  $\zeta(x_t)$ :

$$Q_{ob} = \{q \mid q \in \{p_0, \dots, p_{M-1}\} \wedge \|q - \zeta(x_t)\|_2 \leq d\} \\ d \in \mathbb{R}, 0 < d \leq \max(r_0(x_t), \dots, r_{M-1}(x_t)) \quad (9)$$

Furthermore, the observable unsafe space  $Q_{ob}$  should include reachable sets of other dynamic agents  $\{\zeta_0, \dots, \zeta_i\}$ . However, we are only interested in other dynamic agents which are within some distance  $d \in \mathbb{R}^+$  and so we update our definition  $Q_{ob}$  to include reachable regions of other *close* dynamic agents (10):

$$Q_{ob}^+ = Q_{ob} \cup \{q \mid q \in \bigcup_{i=0}^{N-1} R_{\zeta_i, [t, t']}(x_t) \wedge \|q - \zeta_e(x_t)\|_2 \leq d\} \quad (10)$$

## 4.2 Space Convexification via Separating Hyperplanes

This paper proposes a solution for the computation of  $X_{safe}$  which is based on the convexification of non-convex state space via separating coupled hyperplanes. A hyperplane  $H = \{x \mid a^\top x = b\}$ , where  $a \in \mathbb{R}^n, b \in \mathbb{R}, a \neq 0$ , is a set which splits set  $\mathbb{R}^n$  into two halfspaces. Let us also denote  $H^*$  (11) as one of the halfspaces of the hyperplane  $H$ . A separating hyperplane  $H$  is then said to separate two disjoint convex sets  $A, B$  such that  $A \subseteq H^+$  and  $B \subseteq H^-$  [7].

$$\begin{aligned} H^* &\in \{H^+, H^-\} & H^+ \cap H^- &= H \\ H^+ &= \{x \mid a^\top x \geq b\} & H^- &= \{x \mid a^\top x \leq b\} \end{aligned} \quad (11)$$

An intersection of finite halfspaces is a polyhedron  $P$  (12):

$$P_H = \{x \mid x \in \bigcap_{i=0}^{N-1} H_i^*\} \quad (12)$$

The idea behind a space convexification via separating coupled hyperplanes is to compute a set of hyperplanes  $HS = \{H_0, \dots, H_n\}$  such that together they *create* a polyhedron  $P_{HS}$  which (1) does not intersect with the set observable obstacles of the ego vehicle and (2) the ego vehicle  $\zeta_e(x_t)$  with its target location  $l_\zeta$  are within the polyhedron at the time  $t$  (13):

$$X_{safe} = \{x \mid x \in P_{HS} \wedge P_{HS} \cap Q_{ob}^+ = \emptyset \wedge \zeta_e(x_t) \in P_{HS} \wedge l_{\zeta(x_t)} \in P_{HS}\} \quad (13)$$

The problem of generating a set of separating coupled-hyperplanes  $HS$  can be defined as an optimisation or satisfiability problem (14) in which  $n$  number of hyperplanes are computed such that: 1) each hyperplane separates a part of observable obstacles from the ego vehicle and its target location and 2) all observable obstacles are separated by separating coupled hyperplanes

$$\begin{aligned} \text{compute } HS &= \{H_0, \dots, H_i\} \quad \text{s.t.} \\ \forall q_i^{ob+} \in Q_{ob}^+ &\Rightarrow \exists H_i \in HS \wedge q_i^{ob+} \in H_i^* \wedge \zeta(x_t), l_{\zeta(x_t)} \in \mathbb{R}^n \setminus H_i^* \end{aligned} \quad (14)$$

The convex and safe polyhedron  $X_{safe}$  is the intersection of halfspaces  $H_i^*$  of each hyperplane  $H_i \in HS$  for which  $\zeta_e(x_t) \in H_i^*$  and  $l_{\zeta(x_t)} \in H_i^*$  hold. The problem (14) can be expressed as an optimisation problem on the set of hyperplanes  $HS$  or polyhedron  $P_{HS}$ . One possible *performance* metric could be finding the largest  $P_{HS}$  [13, 8].

## 5 Autonomous Vehicle Control System

### 5.1 Overview of the Closed-Loop Control System

The closed-loop control system for obstacle avoidance which we propose in this paper combines online reachability analysis and non-linear model-predictive control (visualised in Fig. 2). The control cycle can be divided into four main procedures: sensing, environment data processing and local planning, state-space convexification and solving an optimal control problem.

The control system relies on the LiDAR sensor to obtain and identify the set of observable obstacles and safe regions. We then leverage the Ramer-Douglas-Peucker algorithm [14] to simplify the observed

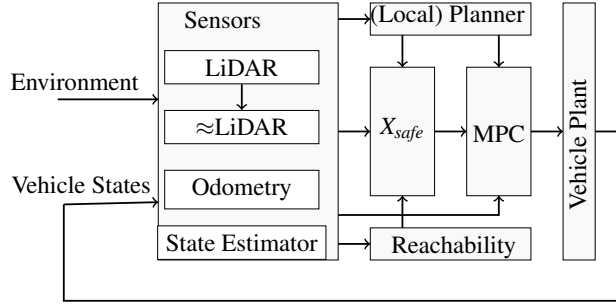


Figure 2: The architecture of the closed-loop control system for obstacle avoidance

LiDAR data and reduce the noisiness of its measurements. Doing so allows us to reduce the computation time needed to produce a set of coupled separating hyperplanes. The other sensors, namely, odometry measurements and the results of state-estimators, are used to determine the state of the ego vehicle and other agents respectively. In this work, we assume that the state of the ego vehicle and opponent agents are estimated perfectly. Therefore, we use the ground truth data provided by the simulator. This data is then passed to a (local) planner (e.g., Follow-the-Gap [43]) to select a target position. We then use reachability analysis to compute the set of reachable states for all agents within the environment.

The computation of separating coupled hyperplanes, which produces a safe and convex  $X_{safe}$ , involves using sensor information, the target location obtained from the local planner, and the set of reachable states of the dynamic agents within the environment. The hyperplanes are then passed to the model-predictive controller, together with the target location and odometry data, which solves an online optimal control problem (5) to determine the optimal inputs for the vehicle.

## 5.2 Computing Separating Coupled-Hyperplanes

The problem of computing separating coupled hyperplanes, which establishes a convex and safe  $X_{safe}$ , can be formulated as an optimisation (or satisfiability) problem. Thus, we present an optimisation-based method for solving (14) in order to separate the observable obstacle set  $Q_{ob}^+$  from the autonomous system  $\zeta_e(x_t)$  and its target location  $l_{\zeta_e(x_t)}$  at the state  $x_t$ .

In Algorithm 1, we describe the computation of our separating coupled-hyperplanes  $H_{0..n}$ . First, the set of unsafe states is included in the set  $Q_{ob}$ . The set of unsafe states consists of the set of observable obstacles from the LiDAR sensor and the reachable states of the dynamic agents. Using this set, we then make use of the state of the ego vehicle, the target location obtained from the local planner, and in the case case of the constrained optimisation method a predefined number of hyper-planes to formulate an optimisation problem. Furthermore, only obstacles (Algorithm 1 ln. 6) and reachable sets (Algorithm 1 ln. 7-8) within a distance  $d$  from the ego vehicle  $\zeta_e(x_t)$  are considered in the hyperplane computation.

**Constrained Optimisation Method** The first approach uses a derivative-free constrained optimisation formulation which utilises a linear approximation of the objective function and optimisation constraints to solve the aforementioned optimisation problem [37]. In the optimisation problem, an individual separating hyperplane  $H_n \in \{H_0, \dots, H_n\}$  is only *responsible* for separating a subset of  $Q_{ob}^+$  from  $\zeta_e(x_t)$  and  $l(x_t)$ , while the set of all hyperplanes considered should separate the vehicle from  $Q_{ob}^+$  as a whole.

For each  $q_{ob} \in Q_{ob}^+$  a separate constraint in the optimisation problem is defined which checks if  $q_{ob}$  is separated from the target location and autonomous system  $\zeta_e(x_t)$  with some hyperplane  $H_n$ . The constrained optimisation method can use different objective functions which characterise how the set of



---

**ALGORITHM 1** The overall algorithm for the computation of separating coupled hyperplanes

---

- 1: **Inputs:** observable radius distance  $d \in \mathbb{R}^+$
  - 2: **Inputs:** states of the ego vehicle  $\zeta_e(x_t) = \{x_e, y_e\}$  and other dynamic agents  $\{\zeta_0, \dots, \zeta_N\}$
  - 3: **Inputs:** static obstacle data  $P = \{p_0, \dots, p_{N-1}\}$  from the LiDAR (ranges Equation (7))
  - 4: Compute target states of the ego vehicle with the local planner  $l(x_t) = \{x_t, y_t\}$
  - 5: Compute reachable states  $R = \bigcup_{i=0}^N R_i$  of observable dynamic agents  $\{\zeta_0, \dots, \zeta_N\}$
  - 6: Compute  $Q_{ob}(x_t)$  by using static obstacle LiDAR data  $\{q \mid q \in P \wedge \|q - \zeta_e(x_t)\|_2 \leq d\}$
  - 7: Compute  $Q_{ob}^+(x_t)$  by combining static and dynamic obstacles  $Q_{ob}(x_t) \cup R$
  - 8: Encode  $q \in Q_{ob}^+, x \in \zeta_e(x_t), x \in l(x_t)$  as constraints of the optimisation problem and solve by using the constrained or bi-level optimisation method
  - 9: **Output:**  $\{H_0, \dots, H_n\}$
- 

hyperplanes is derived. For example, the optimisation problem could try minimising the distance between each  $H_n$  and its associated set of  $q_{ob}$ , or simply be expressed as a satisfiability problem with a constant objective function. We present an analysis of different optimisation objective functions for this purpose in the evaluation section.

**Bi-level Optimisation Method** The problem defined in (14) can also be encoded as a bi-level optimisation problem in (15). The problem is similar to one solved by Deits and Tedrake [11] except we are interested in computing a polygon defined by a minimum number of hyperplanes which contains the largest ellipsoids in the direction of the target location. The [11] maximises ellipsoid in any possible direction, which is not suitable for the racing context, as the most optimal trajectories produced by the MPC will most likely be along the ego vehicle to the target corridor.

The outer part of the problem computes the minimum set of separating hyperplanes (the size of the  $A$  matrix's diagonal) that separate obstacle points  $Q_{ob}^+$  from the ego vehicle and its target location. The inner part of the bi-level optimisation solves the Chebyshev centre [7] problem<sup>1</sup> by finding the centre  $q$  of the largest inscribable ellipsoid with radius  $R$ .

$$\begin{aligned}
 & \arg \min_{A,b} \|\mathbf{diag}(A)\| \text{ s.t.} \\
 & \quad Aq \geq b, \forall q \in Q_{ob}^+ \\
 & \quad Ax \leq b, \forall x \in \zeta_e(x_t) \cup l(x_t) \\
 & \quad \arg \max_{q,R} R \text{ s.t.} \\
 & \quad \quad a_j q + \|A\| R \leq b_j \\
 & \quad \quad R \geq 0
 \end{aligned} \tag{15}$$

### 5.3 Reachability Analysis of Dynamic Obstacles

To perform reachability analysis, we first identify a dynamical model of the vehicle and assume models for the dynamic obstacles within its environment.

---

<sup>1</sup>The reason for maximising the largest inscribable ellipsoid in contrast to directly maximising the area of the safe polyhedron  $X_{safe}$  is efficiency. There are no efficient methods for computing the area of irregular polyhedrons, while the Chebyshev centre problem can be solved sufficiently fast.

### 5.3.1 Dynamic Obstacle Model

The obstacle-tracking problem is a well-studied and challenging topic within the autonomous vehicle, computer vision, and robotics literature [44]. Typically, some assumptions are required in order to constrain the tracking problem to suit the context of the application. In our framework, we assume that the dynamic obstacles are described by a two-dimensional kinematic model and a corresponding bounding box. The equations describing the kinematic model are given as follows:

$$\dot{x} = v_x, \dot{y} = v_y$$

where  $v_x$  and  $v_y$  are the velocities in the  $x$  and  $y$  direction, respectively. Additionally, we make the assumption that we have access to the position and velocity of the other race participants.

While it is possible to use more sophisticated models to describe the behaviour of the dynamic obstacles within the vehicle's environment, for simplicity we selected a two-dimensional kinematic model. However, it is worth noting that there has been a growth in approaches that perform online parameter estimation for dynamic obstacles within a robot's environment through online system identification [15].

### 5.3.2 Online Reachability Computation

Using the dynamics models obtained in the previous sections, the crux of the real-time reachability algorithm is computing the set of reachable states  $R_{[0,T]}(X_0)$  over a finite time horizon. The algorithm utilised within this work is based on mixed face-lifting, which is part of a class of methods that deal with *flow-pipe construction* or *reachtube computation* [17]. This is done using snapshots of the set of reachable states that are enumerated at successive points in time, as outlined in Equation (3).

In general, it is not possible to obtain the exact reachable set  $R_{[0,T]}(X_0)$ , so we compute an over-approximation such that the actual system behaviour is contained within the over-approximation [24]. The algorithm utilised in this work utilises  $n$ -dimensional hyper-rectangles ("boxes") as the set representation to generate reachtubes [17]. Over long reach-times, the over-approximation error resulting from the use of this representation can be problematic. However, for short reach-times it is ideal in terms of its simplicity and speed [6].

Traditionally, reachability approaches have been executed offline because they are computationally intensive endeavours. However, in [6, 17], Bak et al. and Johnson et al. presented a reachability algorithm, based on the seminal mixed face-lifting algorithm [10], capable of running in real-time on embedded processors. The algorithm is implemented as a standalone C-package that does not rely on sophisticated (non-portable) libraries, recursion, or dynamic data structures and is amenable to the anytime computation model in the real-time scheduling literature. In this regime, each task produces a partial result that is improved upon as more computation time is available, known as an anytime algorithm [17]. We refer readers to the following papers for an in-depth treatment of these procedures [10, 6, 17].

## 6 Evaluation

In this section, we present a runtime analysis of proposed algorithms for computing separating coupled hyperplanes and an evaluation of the overall control system by using the F1TENTH simulation platform. In the following section, we first describe an optimisation-free method (MPCC) for computing separating hyperplanes, which will be used to compare against our proposed approaches.

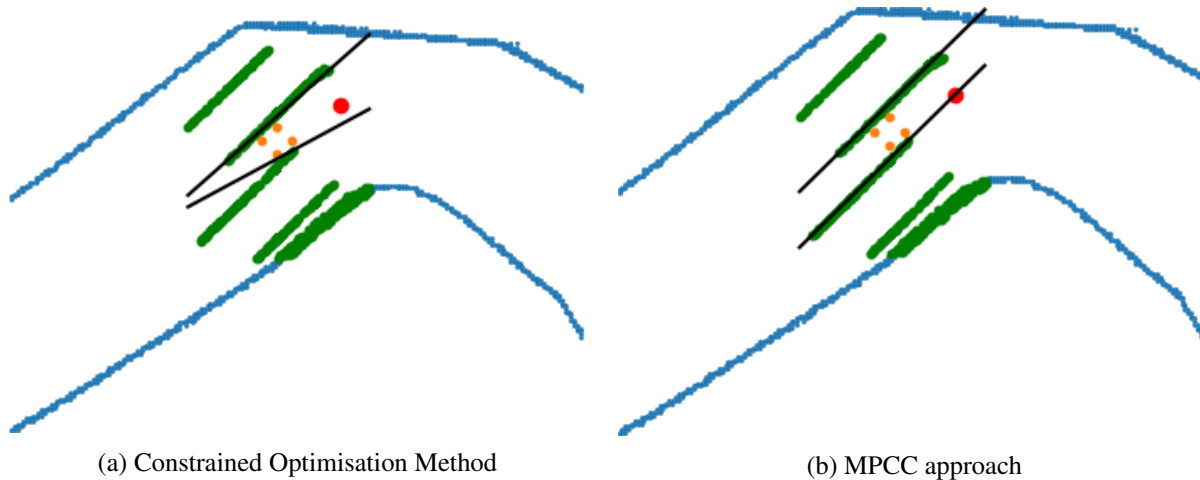


Figure 3: A snapshot of the artificial overtaking scenario with two opponents represented: combined observable static and dynamic obstacles  $Q_{ob}^+$  (green points), corners of the ego vehicle (orange points), target point (red point), computed hyperplanes (black) and the boundaries of the racetrack (blue).

### 6.1 MPCC Optimisation-free Hyperplane Approach

In [25] Liniger et al. tackled the autonomous racing problem via a nonlinear MPC problem that encoded the obstacle avoidance problem by means of a high-level corridor planner based on dynamic programming. The safe corridor that their framework utilised was constructed by projecting the points along the centre line of the track onto the racetrack borders (one for the left border, and one for the right border). Their regime demonstrated success in controlling 1/43 scale race cars, driven at speeds of more than 3 m/s using controllers executing at 50 Hz sampling rate on embedded computing platforms [25]. While their evaluation was limited to environments with static obstacles, we experimented with using such a scheme to obtain the separating hyperplanes framing our MPC problem. We refer readers to the following paper for an in-depth discussion of their approach [25].

### 6.2 Offline Analysis of Convexification Algorithms

Deploying optimisation-based methods into real-time autonomous control systems requires careful consideration of timing constraints issued by the optimisation method. The computation time of the proposed methods can be affected by the number of obstacle points being considered or in the case of the constrained optimisation method by the selected optimisation cost function. In these experiments, we aim to evaluate the *quality* of separating hyperplanes generated by different approaches and the computation time. The former is assessed by inscribing the largest circle with radius  $R$  between generated hyperplanes (the centre of the circle must be between the ego vehicle and its target location) and gives a reasonable size approximation of the  $X_{safe}$  in the travelling direction of ego vehicle. In the first set of experiments, we traversed the ego vehicle along a predefined path on one of the two racetracks: Porto (see Figure 5) and Walker without other dynamic agents, and then considered an artificially created overtaking scenario with one and two opponents (visualised in Figure 3). The results of this experiment are summarised in Table 1. The MPCC approach is clearly more time efficient compared to our proposed approaches as it is not an optimisation-based approach. However, it produces a smaller average inscribed circle radius (i.e., smaller  $X_{safe}$ ), particularly, in the overtaking scenarios with up to 16 per cent smaller  $R$

in comparison to the largest averaged  $R$ . The bi-level optimisation approach is around tenfold faster than the constrained optimisation approach, as its outer problem is a linear programming problem, which can be efficiently solved even with a larger number of obstacles (problem constraints). However, our experiments show that the bi-level optimisation method does not always produce the largest  $X_{safe}$  and in some cases generates more than two hyperplanes, which would negatively affect solving the MPC problem.

In the second experiment, we increased the number of generated (randomly positioned) obstacle points around an ego vehicle to evaluate our method’s scalability with respect to a larger number of obstacles. For the constrained optimisation method we considered three types of objective functions: Hausdorff, Euclidean distances and a satisfiability problem which only requires satisfying optimisation constraints. For each objective function and bi-level optimisation, the number of obstacle points varied from 10 to 2000. The evaluation results are visually shown in Figure 4.

Experiment Scenario	Approach	H	Time (s)	R
Porto (w/o obstacles)	MPCC	2	$6.46 \times 10^{-5}$	1.29
Porto (w/o obstacles)	Constrained Optimisation	2	0.132	<b>1.30</b>
Porto (w/o obstacles)	Bi-level Optimisation	2.16	0.019	1.072
Walker (w/o obstacles)	MPCC	2	$7.65 \times 10^{-5}$	0.702
Walker (w/o obstacles)	Constrained Optimisation	2	0.12	0.675
Walker (w/o obstacles)	Bi-level Optimisation	2.26	0.019	<b>0.715</b>
Overtaking (1 opponent)	MPCC	2	$8.679 \times 10^{-5}$	0.627
Overtaking (1 opponent)	Constrained Optimisation	2	0.18	0.661
Overtaking (1 opponent)	Bi-level Optimisation	2.13	0.022	<b>0.755</b>
Overtaking (2 opponents)	MPCC	2	$9.79 \times 10^{-5}$	0.465
Overtaking (2 opponents)	Constrained Optimisation	2	0.265	<b>0.488</b>
Overtaking (2 opponents)	Bi-level Optimisation	2.59	0.026	0.473

Table 1: Offline evaluation of different methods for computing separating hyperplanes with average computation time, an average inscribed radius  $R$  and an average number of hyperplanes  $H$ .

### 6.3 Real-Time Control System Evaluation

Our real-time evaluation of the overall control system (MPC Hype) includes a diverse set of experiments that include changing the number of racing agents present within the racetrack, including additional dynamic obstacles within the racetrack, adding static obstacles onto the racetrack, and changing the racing environment. We compare the performance of our approach against a set of controllers typically utilized within the F1TENTH racing competitions with respect to two metrics that we refer to as *efficiency*, and *safety*. *Efficiency* is the total distance that the F1TENTH vehicle traverses around the track divided by the amount of time it took to do so.<sup>2</sup> *Safety* corresponds to the controller’s ability to avoid collisions over a set of experimental runs (i.e., 10 collisions in 20 experiments corresponds to a safety score of 50%). The following controllers were utilised as a local planning mechanism for selecting the target point used in our MPC regime. Additionally, we utilised them as a baseline comparison for our approach.

**Pure Pursuit** The Pure Pursuit algorithm is a widely used path-tracking algorithm that was originally designed to calculate the arc needed to get a robot back onto a path [9]. It has shown great success in

<sup>2</sup>This is equivalent to the average speed attained during the experiment.

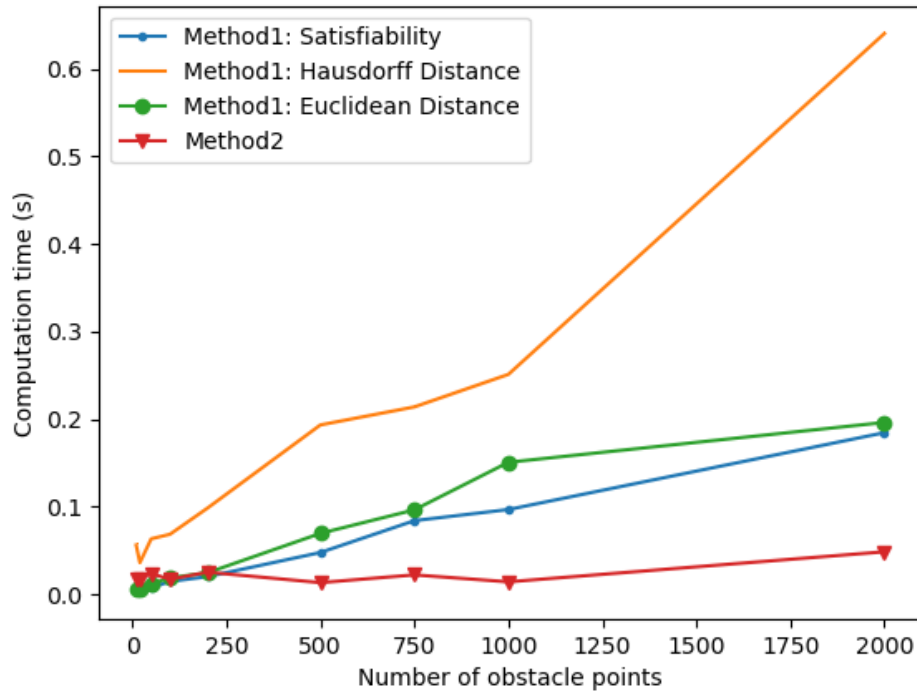


Figure 4: Offline evaluation of separating coupled hyperplane computation time against different numbers of obstacle points (optimisation constraints), different objective functions (method 1 - constrained optimisation approach, method 2 - bi-level optimisation).

being used in numerous contexts, and in this work, we utilise it to design a controller that allows the F1TENTH vehicle to follow a path along the centre of the racetrack.

**Gap Following** Obstacle avoidance is an essential component of a successful autonomous racing strategy. Gap following approaches have shown great promise in dealing with dynamic and static obstacles. They are based on the construction of a gap array around the vehicle used for calculating the best heading angle needed to move the vehicle into the centre of the maximum gap [34]. In this work, we utilise a gap following controller called the “disparity extender” by Otterness et al. that won the F1TENTH competition in April of 2019 [36].

Our evaluation included a sizeable diversity of experiments with respect to the number of vehicles present in the racing environment, the presence of static and dynamic obstacles, the racetrack used for the autonomous race, the local planner chosen to select goal points, and the method selected to obtain the separating hyperplanes. Each configuration was evaluated over 30 experimental runs of 60 seconds. Table 2 displays the results of experiments with two and three cars respectively (separated by horizontal line) on a single track without static obstacles (a screenshot of the two agent experiment is shown in Figure 5). In the table that follows, DE corresponds to the disparity extender, PP corresponds to pure pursuit, MPCC corresponds to the approach presented by Liniger et al. [25], and MPC Hype corresponds to the optimisation-based approach presented in this document. Finally, Race Duration corresponds to

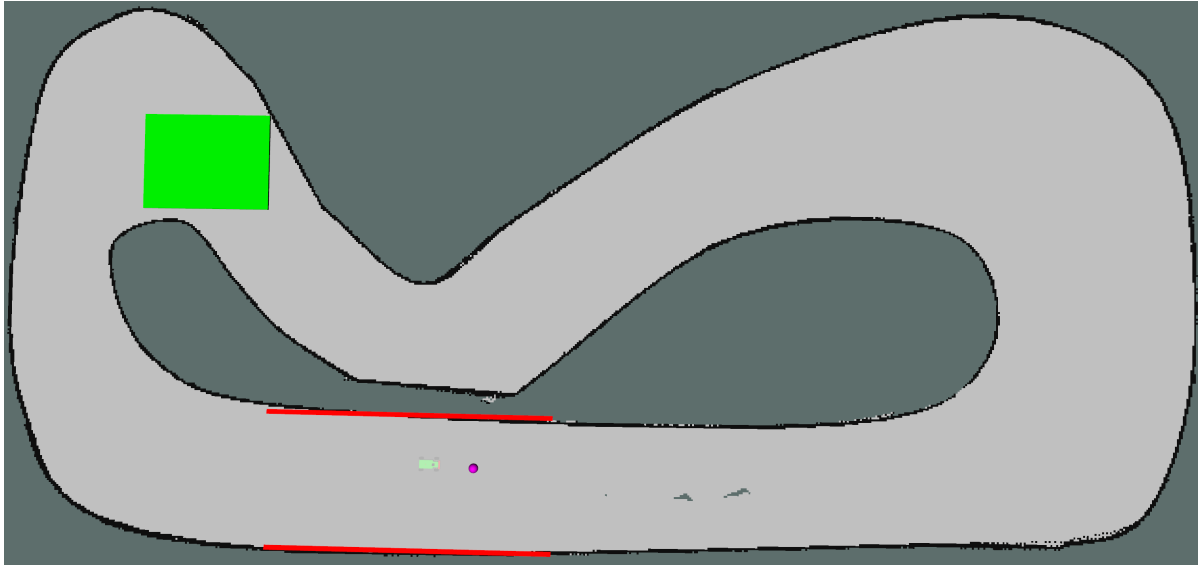


Figure 5: An example of a two-agent racing scenario. The bright green rectangle, represents the reachable set (convex hull) of the opponent vehicle over a  $t = 0.5$  second time horizon, while the faded green vehicle represents the ego vehicle. The purple dot corresponds to the target location obtained from the local planner. The red lines are the two parallel half spaces that approximate the traversable region within the racetrack.

Table 2: Performance summary of two-car experiments (without static obstacles): DE (Disparity Extender), PP (Pure Pursuit), MPCC (Model-Predictive Control with Contouring) and MPC Hype (our control system with the constrained optimisation hyperplane computation)

Track	Approach	Local Planner	Ego Efficiency (m/s)	Opponent Efficiency (m/s)	Race Duration (s)	Safety (%)
Porto	DE	DE	5.29	4.65	51.57	38.33
Porto	MPC Hype	DE	0.00	5.27	5.53	0.00
Porto	MPC Hype	PP	3.06	5.18	25.74	13.33
Porto	MPCC	DE	3.00	4.97	7.12	20.00
Porto	MPCC	PP	3.00	5.34	55.14	46.67
Porto	PP	PP	4.70	5.33	60.0	100.00
Porto	DE	DE	5.38	4.10	33.78	28.33
Porto	MPC Hype	DE	1.19	4.50	5.40	0.00
Porto	MPC Hype	PP	2.75	2.96	43.26	30.00
Porto	MPCC	DE	1.66	4.23	5.39	3.33
Porto	MPCC	PP	1.83	4.00	5.37	16.67
Porto	PP	PP	4.70	3.73	57.30	70.00

the amount of time the agents were able to race before a collision occurred.

The results from the experiments suggest that our proposed control system (MPC Hype) can increase autonomous vehicle safety without loss of efficiency (compared to MPCC), especially when the number of opponent vehicles is increased. However, results also suggest that the performance of our MPC implementation could be further improved, for example, by improving the MPC cost function to generate better speed profiles in corners. This would also provide us with more evidence of the hyperplane approach when ego velocity is increased. Furthermore, our experimentation setup did not differentiate

between different types of collisions, for example, collisions, where the opponent vehicle collided with the back of the ego vehicle and the reverse situation, were treated equally (i.e., counted the same in the safety metric). A more nuanced safety metric with a *blame* factor would provide a better understanding of our control system performance.

## 7 Conclusions and Future Work

This paper presented an optimisation-based approach for static and dynamic obstacle avoidance problems within an autonomous vehicle racing context. Our control regime leveraged online reachability analysis and sensor data to compute the maximal safe traversable region that an agent can traverse within the environment. We described a technique for computing a convex safe region via a novel coupled separating hyperplane algorithm. This derived safe area was then used to formulate a nonlinear model-predictive control problem that sought to find an optimal and safe driving trajectory with varying degrees of efficacy. Our experimental evaluation demonstrated that our approach was feasible as an obstacle avoidance strategy. Finally, we assessed the runtime requirements of our proposed approach by analysing the effects of a set of varying optimisation objectives for generating these coupled hyperplanes.

There are a number of future work directions we would like to explore. Firstly, our study did not consider uncertainty in sensors, our future work will seek to include uncertainties arising from the state estimation of opponent vehicles in their reachable set computation. Secondly, future studies would include an analysis against hierarchical control architectures that decompose the obstacle avoidance problem into planning and trajectory tracking. Lastly, we would like to evaluate the proposed approach on the physical F1TENTH platform in order to validate further that our approach admits low resource requirements.

## References

- [1] M. Althoff & J. M. Dolan (2014): *Online Verification of Automated Road Vehicles Using Reachability Analysis*. *IEEE Transactions on Robotics* 30(4), pp. 903–918, doi:10.1109/TRO.2014.2312453.
- [2] Matthias Althoff, Goran Frehse & Antoine Girard (2021): *Set Propagation Techniques for Reachability Analysis*. *Annual Review of Control, Robotics, and Autonomous Systems* 4(1), pp. 369–395, doi:10.1146/annurev-control-071420-081941.
- [3] Eugene Asarin, Thao Dang, Goran Frehse, Antoine Girard, Colas Le Guernic & Oded Maler (2007): *Recent progress in continuous and hybrid reachability analysis*. *Proceedings of the 2006 IEEE Conference on Computer Aided Control Systems Design, CACSD*, pp. 1582–1587, doi:10.1109/CACSD-CCA-ISIC.2006.4776877.
- [4] Eugene Asarin, Thao Dang & Antoine Girard (2003): *Reachability Analysis of Nonlinear Systems Using Conservative Approximation*. In Oded Maler & Amir Pnueli, editors: *Hybrid Systems: Computation and Control*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 20–35, doi:10.1007/3-540-36580-X\_5.
- [5] A. Bajcsy, S. Bansal, E. Bronstein, V. Tolani & C. J. Tomlin (2019): *An Efficient Reachability-Based Framework for Provably Safe Autonomous Navigation in Unknown Environments*. In: *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 1758–1765, doi:10.1109/CDC40024.2019.9030133.
- [6] S. Bak, T. T. Johnson, M. Caccamo & L. Sha (2014): *Real-Time Reachability for Verified Simplex Design*. In: *2014 IEEE Real-Time Systems Symposium*, pp. 138–148, doi:10.1109/RTSS.2014.21.
- [7] Stephen Boyd & Lieven Vandenbergh (2004): *Convex Optimization*. Cambridge University Press, doi:10.1017/CBO9780511804441.
- [8] J. S. Chang & C. K. Yap (1986): *A Polynomial Solution for the Potato-Peeling Problem*. *Discrete Comput. Geom.* 1(2), p. 155–182, doi:10.1007/BF02187692.

- [9] R. Craig Coulter (1992): *Implementation of the Pure Pursuit Path Tracking Algorithm*. Technical Report CMU-RI-TR-92-01, Carnegie Mellon University, Pittsburgh, PA.
- [10] Thi Xuan Thao Dang (2000): *Verification and Synthesis of Hybrid Systems*. Theses, Institut National Polytechnique de Grenoble - INPG.
- [11] Robin Deits & Russ Tedrake (2015): *Computing Large Convex Regions of Obstacle-Free Space Through Semidefinite Programming*, pp. 109–124. Springer International Publishing, Cham, doi:10.1007/978-3-319-16595-0\_7.
- [12] Huixu Dong, Ching-Yen Weng, Chuangqiang Guo, Haoyong Yu & I-Ming Chen (2021): *Real-Time Avoidance Strategy of Dynamic Obstacles via Half Model-Free Detection and Tracking With 2D Lidar for Mobile Robots*. *IEEE/ASME Transactions on Mechatronics* 26(4), pp. 2215–2225, doi:10.1109/TMECH.2020.3034982.
- [13] Reza Dorrigiv, Stephane Durocher, Arash Farzan, Robert Fraser, Alejandro López-Ortiz, J. Ian Munro, Alejandro Salinger & Matthew Skala (2009): *Finding a Hausdorff Core of a Polygon: On Convex Polygon Containment with Bounded Hausdorff Distance*. In Frank Dehne, Marina Gavrilova, Jörg-Rüdiger Sack & Csaba D. Tóth, editors: *Algorithms and Data Structures*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 218–229, doi:10.1007/978-3-642-03367-4\_20.
- [14] David H Douglas & Thomas K Peucker (1973): *Algorithms for the Reduction of the Number of Points Required to Represent a Digitized Line or its Caricature*. *Cartographica: The International Journal for Geographic Information and Geovisualization* 10(2), pp. 112–122, doi:10.3138/FM57-6770-U75U-7727.
- [15] Gowtham Garimella, Matthew Sheckells & Marin Kobilarov (2017): *Robust obstacle avoidance for aerial platforms using adaptive model predictive control*. In: *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5876–5882, doi:10.1109/ICRA.2017.7989692.
- [16] Gabriel Hartmann, Zvi Shiller & Amos Azaria (2021): *Autonomous Head-to-Head Racing in the Indy Autonomous Challenge Simulation Race*. *CoRR* abs/2109.05455. arXiv:2109.05455.
- [17] Taylor T. Johnson, Stanley Bak, Marco Caccamo & Lui Sha (2016): *Real-Time Reachability for Verified Simplex Design*. *ACM Trans. Embed. Comput. Syst.* 15(2), doi:10.1109/RTSS.2014.21.
- [18] Chanyoung Jung, Seungwook Lee, Hyunki Seong, Andrea Finazzi & David Hyunchul Shim (2021): *Game-Theoretic Model Predictive Control with Data-Driven Identification of Vehicle Model for Head-to-Head Autonomous Racing*. *CoRR* abs/2106.04094, doi:10.48550/arXiv.2106.04094.
- [19] Sertac Karaman, Matthew R. Walter, Alejandro Perez, Emilio Frazzoli & Seth Teller (2011): *Anytime Motion Planning using the RRT\**. In: *2011 IEEE International Conference on Robotics and Automation*, pp. 1478–1483, doi:10.1109/ICRA.2011.5980479.
- [20] Christos Katrakazas, Mohammed Quddus, Wen-Hua Chen & Lipika Deka (2015): *Real-time motion planning methods for autonomous on-road driving: State-of-the-art and future research directions*. *Transportation Research Part C: Emerging Technologies* 60, pp. 416–442, doi:10.1016/j.trc.2015.09.011.
- [21] Mark Keil & Jack Snoeyink (2002): *On The Time Bound For Convex Decomposition Of Simple Polygons*. *International Journal of Computational Geometry & Applications* 12(03), pp. 181–192, doi:10.1142/S0218195902000803.
- [22] Shivesh Khaitan, Qin Lin & John M. Dolan (2021): *Safe Planning and Control Under Uncertainty for Self-Driving*. *IEEE Transactions on Vehicular Technology* 70(10), pp. 9826–9837, doi:10.1109/TVT.2021.3108525.
- [23] Karen Leung, Edward Schmerling, Mengxuan Zhang, Mo Chen, John Talbot, J Christian Gerdes & Marco Pavone (2020): *On Infusing Reachability-Based Safety Assurance within Probabilistic Planning Frameworks for Human-Robot Vehicle Interactions*. *The International Journal of Robotics Research* 39(10-11), pp. 1326–1345, doi:10.1177/0278364920950795.
- [24] Qin Lin, Xin Chen, Aman Khurana & John Dolan (2020): *ReachFlow: An Online Safety Assurance Framework for Waypoint-Following of Self-driving Cars*. In: *International Conference on Intelligent Robots and Systems (IROS)*, IROS'2020, IEEE, doi:10.1109/IROS45743.2020.9341122.



- [25] Alexander Liniger, Alexander Domahidi & Manfred Morari (2014): *Optimization-based autonomous racing of 1:43 scale RC cars*. *Optimal Control Applications and Methods* 36(5), p. 628–647, doi:10.1002/oca.2123.
- [26] Changliu Liu, Tomer Arnon, Christopher Lazarus, Clark W. Barrett & Mykel J. Kochenderfer (2019): *Algorithms for Verifying Deep Neural Networks*. CoRR abs/1903.06758, doi:10.48550/arXiv.1903.06758, arXiv:1903.06758.
- [27] Changliu Liu, Chung-Yen Lin & Masayoshi Tomizuka (2017): *The Convex Feasible Set Algorithm for Real Time Optimization in Motion Planning*, doi:10.48550/ARXIV.1709.00627. Available at <https://arxiv.org/abs/1709.00627>.
- [28] Jiechao Liu, Paramsothy Jayakumar, Jeffrey L. Stein & Tulga Ersal (2018): *A nonlinear model predictive control formulation for obstacle avoidance in high-speed autonomous ground vehicles in unstructured environments*. *Vehicle System Dynamics* 56(6), pp. 853–882, doi:10.1080/00423114.2017.1399209.
- [29] Joseph Lorenzetti, Mo Chen, Benoit Landry & Marco Pavone (2018): *Reach-Avoid Games Via Mixed-Integer Second-Order Cone Programming*. In: *2018 IEEE Conference on Decision and Control (CDC)*, pp. 4409–4416, doi:10.1109/CDC.2018.8619382.
- [30] Tim Mercy, Wannes Van Loock & Goele Pipeleers (2016): *Real-time motion planning in the presence of moving obstacles*. In: *2016 European Control Conference (ECC)*, pp. 1586–1591, doi:10.1109/ECC.2016.7810517.
- [31] Tim Mercy, Ruben Van Parys & Goele Pipeleers (2018): *Spline-Based Motion Planning for Autonomous Guided Vehicles in a Dynamic Environment*. *IEEE Transactions on Control Systems Technology* 26(6), pp. 2182–2189, doi:10.1109/TCST.2017.2739706.
- [32] Siti Hajar Ashikin Mohammad, Muhammad Akmal Jeffril & Nohaidda Sariff (2013): *Mobile robot obstacle avoidance by using Fuzzy Logic technique*. In: *2013 IEEE 3rd International Conference on System Engineering and Technology*, pp. 331–335, doi:10.1109/ICSEngT.2013.6650194.
- [33] Matthew O’Kelly, Varundev Sukhil, Houssam Abbas, Jack Harkins, Chris Kao, Yash Vardhan Pant, Rahul Mangharam, Dipshil Agarwal, Madhur Behl, Paolo Burgio & Marko Bertogna (2019): *F1/10: An Open-Source Autonomous Cyber-Physical Platform*. CoRR abs/1901.08567, doi:10.48550/arXiv.1901.08567, arXiv:1901.08567.
- [34] Matthew O’Kelly, Hongrui Zheng, Dhurv Karthik & Rahul Mangharam (2020): *F1TENTH: An Open-source Evaluation Environment for Continuous Control and Reinforcement Learning*. In Hugo Jair Escalante & Raia Hadsell, editors: *Post Proceedings of the NeurIPS 2019 Demonstration and Competition Track*, Proceedings of Machine Learning Research, PMLR.
- [35] Michael Otte & Emilio Frazzoli (2016): *RRTX: Asymptotically optimal single-query sampling-based motion planning with quick replanning*. *The International Journal of Robotics Research* 35(7), pp. 797–822, doi:10.1177/0278364915594679.
- [36] Nathan Otterness (2019): *The "Disparity Extender" Algorithm, and F1/Tenth*. Available at <https://www.nathanotterness.com/2019/04/the-disparity-extender-algorithm-and.html>.
- [37] M. J. D. Powell (2007): *A View of Algorithms for Optimization without Derivatives*. Technical Report DAMTP 2007/NA03, University of Cambridge.
- [38] Ugo Rosolia & Francesco Borrelli (2019): *Learning How to Autonomously Race a Car: a Predictive Control Approach*. arXiv:1901.08184.
- [39] Ugo Rosolia, Stijn De Bruyne & Andrew G. Alleyne (2017): *Autonomous Vehicle Control: A Nonconvex Approach for Obstacle Avoidance*. *IEEE Transactions on Control Systems Technology* 25(2), pp. 469–484, doi:10.1109/TCST.2016.2569468.
- [40] Tobias Schoels, Luigi Palmieri, Kai Oliver Arras & Moritz Diehl (2019): *An NMPC Approach using Convex Inner Approximations for Online Motion Planning with Guaranteed Collision Freedom*. CoRR abs/1909.08267, doi:10.48550/arXiv.1909.08267.

- [41] Eelco Scholte & Mark E. Campbell (2008): *Robust Nonlinear Model Predictive Control With Partial State Information*. *IEEE Transactions on Control Systems Technology* 16(4), pp. 636–651, doi:10.1109/TCST.2007.912120.
- [42] Hiroya Seki, Satoshi Ooyama & Morimasa Ogawa (2002): *Nonlinear Model Predictive Control Using Successive Linearization*. *Transactions of the Society of Instrument and Control Engineers* 38, pp. 61–66, doi:10.1109/AIM.2017.8014275.
- [43] Volkan Sezer & Metin Gokasan (2012): *A novel obstacle avoidance algorithm: “Follow the Gap Method”*. *Robotics and Autonomous Systems* 60(9), pp. 1123–1134, doi:10.1016/j.robot.2012.05.021.
- [44] Alper Yilmaz, Omar Javed & Mubarak Shah (2006): *Object Tracking: A Survey*. *ACM Comput. Surv.* 38(4), p. 13–es, doi:10.1145/1177352.1177355.
- [45] Yang Zeqing, Liu Libing, Tan Zhihong & Liu Weiling (2008): *Application of Adaptive Genetic Algorithm in flexible inspection path planning*. In: *2008 27th Chinese Control Conference*, pp. 75–80, doi:10.1109/CHICC.2008.4605656.
- [46] Yalong Zhang, Zhenghua Liu & Le Chang (2017): *A new adaptive artificial potential field and rolling window method for mobile robot path planning*. In: *2017 29th Chinese Control And Decision Conference (CCDC)*, pp. 7144–7148, doi:10.1109/CCDC.2017.7978472.
- [47] Liu Zhiyang & Jiang Tao (2017): *Route planning based on improved artificial potential field method*. In: *2017 2nd Asia-Pacific Conference on Intelligent Robot Systems (ACIRS)*, pp. 196–199, doi:10.1109/CCDC52312.2021.9602174.

# Automatic Generation of Scenarios for System-level Simulation-based Verification of Autonomous Driving Systems\*

Srajan Goyal

Fondazione Bruno Kessler  
University of Trento  
Trento, Italy  
sgoyal@fbk.eu

Alberto Griggio

Fondazione Bruno Kessler  
Trento, Italy  
griggio@fbk.eu

Jacob Kimblad

Fondazione Bruno Kessler  
Trento, Italy  
jkimblad@fbk.eu

Stefano Tonetta

Fondazione Bruno Kessler  
Trento, Italy  
tonettas@fbk.eu

With increasing complexity of Automated Driving Systems (ADS), ensuring their safety and reliability has become a critical challenge. The Verification and Validation (V&V) of these systems are particularly demanding when AI components are employed to implement perception and/or control functions. In ESA-funded project VIVAS, we developed a generic framework for system-level simulation-based V&V of autonomous systems. The approach is based on a simulation model of the system, an abstract model that describes symbolically the system behavior, and formal methods to generate scenarios and verify the simulation executions. Various coverage criteria can be defined to guide the automated generation of the scenarios.

In this paper, we describe the instantiation of the VIVAS framework for an ADS case study. This is based on the integration of CARLA, a widely-used driving simulator, and its ScenarioRunner tool, which enables the creation of diverse and complex driving scenarios. This is also used in the CARLA Autonomous Driving Challenge to validate different ADS agents for perception and control based on AI, shared by the CARLA community. We describe the development of an abstract ADS model and the formulation of a coverage criterion that focuses on the behaviors of vehicles relative to the vehicle with ADS under verification. Leveraging the VIVAS framework, we generate and execute various driving scenarios, thus testing the capabilities of the AI components. The results show the effectiveness of VIVAS in automatically generating scenarios for system-level simulation-based V&V of an automated driving system using CARLA and ScenarioRunner. Therefore, they highlight the potential of the approach as a powerful tool in the future of ADS V&V methodologies.

## 1 Introduction

In the rapid evolution of Autonomous Driving Systems (ADS), the problem of ensuring their safety and reliability has become a paramount concern. The Verification and Validation (V&V) of these systems necessitate the assessment of their correctness in a multitude of dynamic and complex real-world scenarios. To address this challenge, the integration of powerful simulation tools with advanced verification methodologies has gained considerable attention [10, 18].

---

\*This work has been supported by: the “AI@TN” project funded by the Autonomous Province of Trento; the PNRR project FAIR - Future AI Research (PE00000013), under the NRRP MUR program funded by the NextGenerationEU; and the PNRR MUR project VITALITY (ECS00000041), Spoke 2 ASTRA - Advanced Space Technologies and Research Alliance.

Under the support of ESA funding, the VIVAS project [12] was dedicated to developing a generic framework tailored for system-level simulation-based V&V of autonomous systems. The approach is based on a simulation model of the system, an abstract model that describes symbolically the system behavior, and formal methods to generate scenarios and verify the simulation executions. It permits the specification of diverse coverage criteria, thereby directing the automated creation of scenarios, and formal properties to be verified on the simulation runs. The framework has been created for space applications and applied to two use cases employing AI for resource prediction and opportunistic science. The system under test is based on the robotic digital twin developed in ROBDT [2].

In the automotive context, the CARLA simulator [9] has established itself as a widely used platform for simulating intricate driving scenarios in a controlled virtual environment. Its ScenarioRunner tool [25] further enhances its capabilities by enabling the specification of diverse and complex scenarios based on the reuse of various predefined car behaviors. Various works proposed AI-based solutions (e.g., [4, 22, 29]) for the perception and control components of cars that are integrated with the CARLA simulator for their validation. The CARLA community also organized a competition to compare and rank such solutions [24]. These autonomous driving agents, grounded in AI methodologies, serve as crucial components in achieving the autonomy of cars. Their integration in CARLA allows evaluating the ADS behavior across different scenarios. However, such validation is so far based on a few manually crafted scenarios.

This paper explores the application of the VIVAS framework to automatically generate scenarios for a system-level simulation-based verification of autonomous driving systems. This integration facilitates a comprehensive assessment of ADS correctness under various conditions, contributing to the enhancement of their safety and reliability. The paper describes the abstract ADS model in the extended SMV language handled by the nuXmv [3] symbolic model checker, capturing essential aspects of its functionality and behavior in a simple highway traffic situation. It then details the formulation of a coverage criterion based on such an abstract model, focusing on the interactions between other vehicles and the vehicle of the ADS under verification (hereafter called *ego*). The VIVAS integration finally consists of a translation of the abstract traces generated from the abstract model to the ScenarioRunner specification and a mapping back of the simulation runs to abstract traces for runtime verification of formal properties. The experimental evaluation shows how VIVAS is able to generate interesting scenarios effectively evaluating the behavior of the AI-based agents.

The rest of the paper is organized as follows: in Section 2, we summarize the related works and compare them with our approach; Section 3 describes in more detail the VIVAS framework and its components; in Section 4, we detail the instantiation to the ADS application; Section 5 shows the results while Section 6 draws conclusions and some directions for future work.

## 2 Related Work

Over the last decade, we have witnessed significant efforts in the verification of AI-based autonomous systems using formal methods. Many works focus on formal verification of neural networks, for example encoding them into constraint solving (e.g., [14–16]) or using abstraction (e.g., [19, 23]), just to name a few approaches. Our approach instead is rooted in the line of research (e.g., [10, 26, 28]) that tackles the verification at the system level using a simulator. This integrates the AI components, potentially using machine learning (ML) models, for perception or control, in the context of a closed-loop cyber-physical system. As in VerifAI [10], the simulation traces are then formally analyzed with monitoring and runtime verification techniques.

Differently from the mentioned approaches, we exploit an abstract symbolic model to generate automatically the scenarios and define a coverage criterion for the generated test cases. While previous approaches focus on the automated synthesis of the simulation parameters for a specific scenario (e.g., different car movements to change lanes in front of the ego car), we concentrate on the generation of different functional scenarios (e.g., sequences of scenes with different change lanes of non-ego cars). Moreover, in this paper, we map such abstract symbolic scenarios to the scenario specification language of CARLA to verify ADS with different available AI solutions. So, the case study is based on available benchmarks for AI-based ADS taken as is.

There are in fact a variety of scenario specification languages that can be used in this context. VerifAI uses the Scenic language [13, 27] to model the abstract feature space defining the scenarios, which can be instantiated to test cases. Scenic is a probabilistic programming language for scenario generation specifically designed to test the robustness of systems containing AI and ML components by allowing the generation of rare events. It allows the specification of spatial and temporal relationships between objects of a scenario as well as composing several scenarios into more complex ones. By the use of distributions for encoding interesting parameters, Scenic will perform automatic test case generation through the use of sampling. Similarly, the Paracosm [18] framework is a programmatic interface that can be used to create various automotive driving simulation scenarios through the design of parameterized environments and test cases. The parameters control the environment in the scenario including the behavior of the actors and can include things such as pedestrians, lanes, and light conditions. Parameters are specified using either discrete or continuous domains and test cases are instantiated from the domain using random sampling and Halton sampling respectively. A coverage criterion is then defined over the coverage of the domains, where  $k$ -wise combinatorial coverage is used for the discrete domains and dispersion is used for continuous domains. Although Paracosm can provide output using the OpenDRIVE format, it is primarily coupled to be used with the Unity game engine, and as such scenarios are modeled using the C# programming language. The Measurable Scenario Description Language (M-SDL) [11] is another scenario description language similar to Scenic. In M-SDL, one captures the behavior of identified actors in scenarios. M-SDL makes use of pre-defined basic building blocks such as actors (including the AV) along with some pre-defined behaviors, sets of possible routes, and environmental conditions. Libraries then use the basic building blocks to implement more complex behaviors such as cars overtaking, running red lights, driving on a highway, etc. Since M-SDL scenarios are abstract and parameterized, a single scenario can map onto many concrete ones through the use of sampling. ScenarioRunner [25] is a module of CARLA that allows traffic scenario definition and execution for the CARLA simulator. The scenarios can be defined through a Python interface or using the OpenSCENARIO standard [1]. ScenarioRunner is used to validate AI solutions for ADS. These results can be validated and shared in the CARLA Autonomous Driving Leaderboard [24], an open platform for the community to fairly compare their progress, evaluating agents in realistic traffic situations.

For all the above languages, the scenario must be specified manually, to then derive the test cases automatically. VIVAS instead provides a model-based approach to generate the scenarios automatically based on a coverage criterion that defines the interesting combinations of situations. In this paper, we focus on the integration with CARLA, because it allows the verification of the solutions shared by the CARLA community. However, the approach can also work with different specification languages, and we have, for example, a prototype integration with Scenic interfaced with CARLA. We have not presented the results of this integration in this paper since the ego model is based on Newtonian physics, with no AI models involved in the autonomous driving pipeline.

Although not specifically focused on AI-based systems, another very relevant work is described in [17], which proposes an optimization-based approach to synthesize ADS scenarios from formal spec-

ifications and a given map. Their formal specification of scenarios corresponds to our abstract scenario and is also synthesized from a symbolic model. However, test case generation does not follow any coverage criteria but enumerates specifications starting from an initial scene. In principle, our coverage-driven generation of scenarios can be combined with various techniques to concretize the scenario with different trajectories and sampling of the different environment parameters.

TAF [20] is another tool for automated test case generation of autonomous systems. Their abstract model is defined in an XML-based domain-specific language. It includes semantic constraints on the initial conditions of the environment and its agents (unlike the additional state transition systems in our work), which are solved using SMT solvers to generate abstract test cases. Random sampling is combined with these solvers to diversify the test cases, with an expert given coverage of data values. Their coverage criteria is based on covering parameter values to instantiate the scenarios. Although constraints on time can be expressed, more generic temporal specification on the sequences of actions and the related coverage criteria are not supported as in our approach. On the other side, our framework can be extended to constraints with quantifiers and complex data structure as in [21], which are currently not supported in VIVAS.

### 3 The VIVAS Framework

VIVAS is a V&V framework for generating test cases for autonomous systems (possibly using AI/ML components) via a combination of system-level simulation and symbolic model checking. VIVAS makes use of formal, symbolic models of the environment and system components to generate *abstract test scenarios* for the autonomous system of interest using model checking techniques. The abstract test scenarios are then instantiated by the concretization of the abstract parameters to provide concrete scenarios to be executed on a system-level simulator encompassing AI/ML models, to obtain execution traces that are in turn analyzed by an automatically-generated monitor. The output of the framework is a V&V result consisting of coverage statistics of the executed traces with respect to the symbolic models and quantitative and qualitative information for each use case. The overview of the architecture can be seen in Fig. 1, which depicts the main parts of the VIVAS framework. These are the abstract scenario generator, the concrete scenario generator, the simulator, and the executor monitor.

**Abstract Scenario Generation.** Scenario generation is the first step of the approach. The starting point is a formal, symbolic model of the system, which provides an abstract view of both the environment and the components under test (including AI/ML parts). ML components are defined in a declarative manner, approximated in terms of input-output mapping. Abstract test scenarios are generated from the formal system model using symbolic model checking techniques by the abstract scenario generator. Abstract scenarios are defined as combinations of values of predicates describing interesting behaviors of the abstract system. From the technical point of view, each abstract scenario is encoded as a formal property that is expected to be violated by the system (i.e. a property specifying that “the scenario cannot occur in the abstract system”). For each such property defined by the abstract scenario generator, a model checker will be executed on the system model, with the goal of finding a counterexample to the property. By construction, each such counterexample corresponds to an execution trace witnessing the realization of the abstract scenario of interest.

**Concrete Scenario Generation.** Each of the traces produced by the model checker is then refined into a (set of) concrete scenarios that can be used to drive the system-level (concrete) simulator. Due

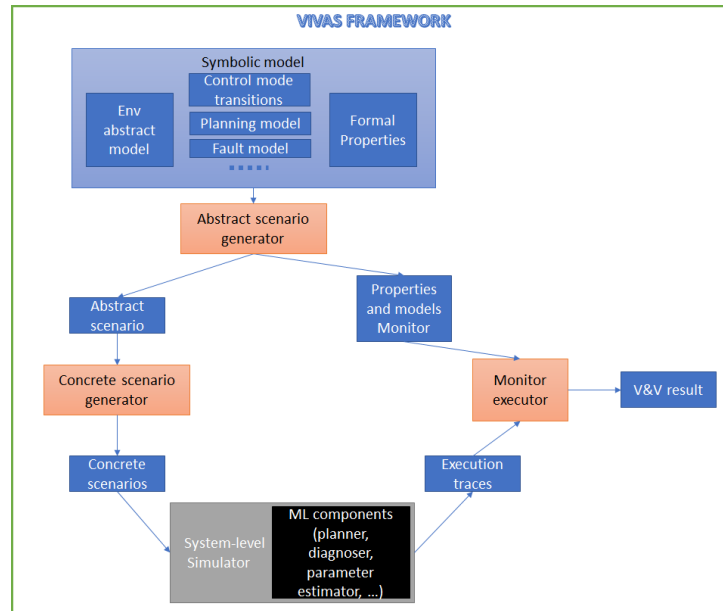


Figure 1: Top-level architecture of VIVAS framework. Blue boxes: artifacts, Orange boxes: code.

to uncertainties and abstractions in the abstract scenarios, a one-to-many mapping is defined where a single abstract scenario can be instantiated to many, possibly infinite, concrete scenarios. This is done by defining a mapping between the abstract values and the set of concrete values that they represent and then appropriately sampling from the sets. An example is that the abstract model might discretize the time of day into dusk, dawn, midnight, and midday. The concretizer then has to sample the actual time for the simulation.

**Simulation.** The task of the system-level simulator is to run a simulation of the target asset under the requested conditions, by configuring the system, its environment, and its inputs as specified in the concrete scenario produced by VIVAS. Upon completion, the simulator provides the corresponding execution trace of the system, containing all the necessary details to evaluate the properties of interest.

**Execution Monitor.** Each concrete scenario produced is executed by the simulator, which generates a corresponding concrete execution trace. This trace is then used to determine whether:

1. the concrete execution of the system satisfies the property of interest, and
2. the concrete execution of the system complies with the input abstract scenario (which defines the situation of interest for the current test).

This is done by formally evaluating the trace with a runtime monitor that is automatically generated from the formal specification of the property and the abstract system model. The trace evaluation can have four possible outcomes:

1. The trace complies with the abstract scenario (defining the situation under test), and it also satisfies the property: the test execution is relevant and the test passes.
2. The trace complies with the abstract scenario, but it does not satisfy the property: this corresponds to a test failure on a relevant scenario, and it should be reported to the user.

3. The trace satisfies the property, but it does not comply with the abstract scenario: this corresponds to a (good) execution in an unexpected situation, in which some of the assumptions defining the scenario might be violated. This might be due to imprecisions/abstractions in the symbolic model and in the concretizer, which might prevent the realization of the abstract scenario under analysis. This situation might be reported to the user, as it might suggest that a revision/refinement of the symbolic model might be needed.
4. The trace violates the property and it does not comply with the abstract scenario: this corresponds to a test failure in an unexpected situation. Similarly to the above, it might be a warning that the symbolic model of the system is not precise enough to capture the situations of interest defined by the abstract scenario.

**Abstract and concrete coverage.** Ensuring an adequate level of coverage is one of the primary goals of a good set of tests. In VIVAS, coverage is defined with respect to a domain-specific notion of “interesting situations”, which are those that are (implicitly) defined by the possible combinations of values of predicates that are used by the abstract scenario generator to produce abstract traces. By construction, therefore, VIVAS tries to enumerate abstract scenarios that ensure a 100% degree of coverage of the *abstract* situations of interest<sup>1</sup>. Each abstract scenario is then refined into one or more concrete simulation inputs, leading to corresponding concrete simulation traces. In order to determine the *concrete* coverage (i.e., the degree of coverage of interesting situations at the concrete level), the VIVAS monitor analyzes the execution traces. It checks for compliance with the property of interest and the corresponding abstract scenario’s specification (i.e., the “interesting situation”) from which the concrete executions originate.

Only executions that satisfy the abstract scenario specification contribute to the coverage at the concrete level: if an execution does not comply with its abstract specification, it represents an unexpected situation from which no coverage information can be drawn<sup>2</sup>.

## 4 Autonomous Driving Application

In order to apply the VIVAS methodology to ADS application, we instantiate various components of the VIVAS framework. We choose CARLA simulator as it is widely used in the automotive domain and it has a large community that provides various AI-based solutions for perception and control. We define an abstract model that focuses on highway scenarios where the ego is surrounded by other vehicles in various dynamic situations. In the following, we provide details about the different components.

### 4.1 CARLA Simulation Model and AI Components

CARLA [9] is a high-fidelity open-source simulator that provides a dynamic environment for the development, testing, and validation of AD systems. It is written in C++ as a plugin for Unreal Engine. As a standalone package, it provides pre-defined maps with 3D meshes ranging from city roads with intersections to highways, to mimic real-world landscapes for the agents to drive in. Various sensor models

---

<sup>1</sup>Note that a 100% degree of coverage might not be reached, either because some situations are not feasible already at the abstract level, or because the model checker cannot find a witness trace for the scenario specification within the given resource budget (time and/or memory).

<sup>2</sup>Note that in principle such a situation might still provide *some* information (e.g. it might still cover a different but still interesting situation); therefore, the test result is still reported to the user. However, determining this might not be obvious in general, and therefore we opted for the conservative choice of excluding the test from the computation of the degree of coverage in such cases.



(cameras, Lidar, radar, GPS, IMU) are provided to gather the data from the environment. The simulator includes many vehicle models, from small cars to large trucks, with different properties like mass, dynamics, and controls. A simulation is composed of (i) the CARLA Simulator that computes the physics and renders the scene and all actor properties, (ii) client scripts written using a Python API, that allows control of the actors, sensors, and environmental conditions.

**AI-based components.** The CARLA community through its leaderboard competition [24] provides various state-of-the-art AI solutions for end-to-end autonomous driving. However, only a few of them provide the necessary code and well-trained models for their methodologies to be evaluated and built upon. We specifically tested Interfuser [22], TCP [29], and LAV [4], all three currently in the top 5 of the leaderboard. Within a few test runs of the AI agent provided with TCP, we noticed that the ego vehicle brakes to a standstill as soon as any other vehicle arrives next to it in its adjacent lane. We consider it too conservative of an autonomous behavior to test our verification methodology. The LAV agent on the other hand behaved well autonomously (in accordance with its overall score on the leaderboard) in terms of route completion and collision avoidance. However, it had an erratic behavior of changing lanes non-deterministically at scenario instantiation. It would require us to make ad-hoc changes to relative positions of the non-egos with respect to ego in every concrete scenario we generate.

We therefore chose the Interfuser agent as the AI system under test for our V&V methodology. It is currently ranked 2 on the leaderboard (rank 1 among the open-source solutions). This solution primarily focuses on the safety of AD systems by generating interpretable semantic features of the environment through multi-model sensor fusion, for constraining the agent’s low-level control actions in real-time within safe sets. The perception system processes the data gathered by 3 RGB cameras and one Lidar sensor.

All three AI agents mentioned above share the following main characteristics:

- The maximum driving speed is limited to 5 *m/s*, which is quite conservative for highway driving;
- The ego always travels in its own lane: an external route for the ego to follow needs to be provided. It may change lanes only based on the waypoints of this route on the map. Hence, it never overtakes slow-moving cars in front of it in the same lane. Ego just follows them while maintaining a safe distance, or keeping a stand-still.
- Standard rules of the road for overtaking only on the left (or the right) are not applied.

Note that our V&V methodology is agnostic to the AI solution chosen for the simulator. Since the abstract test scenarios are generated from the symbolic model of the system, abstract coverage would be the same for different AI solutions, although the concrete coverage may vary. In future work, we will use our methodology to benchmark other AI solutions as well.

## 4.2 Abstract Model and Coverage Criterion

We specify our abstract model as a synchronous symbolic transition system written in the language of the nuXmv [3] model checker. The model consists of 3 vehicles (one “ego” car, representing the autonomous system under test, and two other cars) moving on a highway with 3 lanes. The vehicles all drive in the same direction. The ego is constrained to stay in the middle lane and tries to maintain a given cruise speed, braking when necessary to avoid collisions with other cars, and possibly accelerating to reach the target speed. The other two “non-ego” cars can move freely on the highway, with arbitrary accelerations, braking, and lane change maneuvers (subject to physical constraints about min/max acceleration rates

```

MODULE Car(id)
IVAR acceleration : real;
VAR pos : real;
    lane : 0 .. MAX_LANE;
    speed : real;
DEFINE changing_lane := next(lane) != lane;
TRANS
    changing_lane -> (speed <= max_lane_change_speed &
        next(speed) <= max_lane_change_speed);
TRANS
    changing_lane -> (acceleration <= max_lane_change_acceleration &
        acceleration >= (- max_lane_change_braking));
TRANS
    next(speed) = max(speed + acceleration * TIME_STEP, 0);
TRANS
    changing_lane ?
        (next(pos) = pos + (speed + next(speed)) / 2 * TIME_STEP * 0.95) :
        (next(pos) = pos + (speed + next(speed)) / 2 * TIME_STEP);
TRANS
    (next(lane) = lane) | (next(lane) = lane + 1) | (next(lane) = lane - 1);

MODULE Ego(car1, car2)
-- VAR declarations...
DEFINE
    time_to_stop := speed / (-MAX_BRAKING);
    collision_next := (car1.lane = lane & car1.pos >= pos & speed > 0 &
        (car1.pos - pos) / speed <= time_to_stop) |
        (car2.lane = lane & car2.pos >= pos & speed > 0 &
        (car2.pos - pos) / speed <= time_to_stop);
TRANS
    collision_next ?
        (acceleration = MAX_BRAKING & target_speed = 0) :
        (target_speed = EGO_CRUISE_SPEED &
        ((speed < target_speed) ->
            (next(speed) = min(target_speed,
                speed + MAX_ACCELERATION * TIME_STEP))));
INVAR -- the cars do not crash into each other on purpose
    ((abs(pos - car1.pos) > SAFE_DISTANCE) | (lane != car1.lane)) &
    ((abs(pos - car2.pos) > SAFE_DISTANCE) | (lane != car2.lane)) &
    ((abs(car1.pos - car2.pos) > SAFE_DISTANCE) | (car1.lane != car2.lane));

```

Figure 2: Excerpt of the nuXmv code for the abstract model.

and speed limits, taken from publicly available online car databases), but are not allowed to crash into each other or the ego. We use a discrete model of time, in which each transition of the system corresponds to a time-lapse of 1 second. We use the theory of real arithmetic to encode the transition relation of the system, using mostly linear constraints to compute the updates to the speed and locations of the vehicles (thanks to the discretization of time). An excerpt of the symbolic model is shown in Fig. 2. The module `Car` is shared by different non-ego vehicles. Different transition relations on speed, acceleration, and position need to hold when a non-ego changes lane (with `changing_lane`). For the `Ego` module, we define the collision condition (`collision_next`) with non-ego vehicles (`car1` & `car2`, in this case). If True, the ego brakes with the `max_braking` until it stops; else it continues with (or reach towards) its `target_speed`.

In order to enumerate abstract scenarios encoding potentially-interesting traffic situations, we define for each non-ego car a set of predicates specifying its position relative to the ego, in terms of occupation

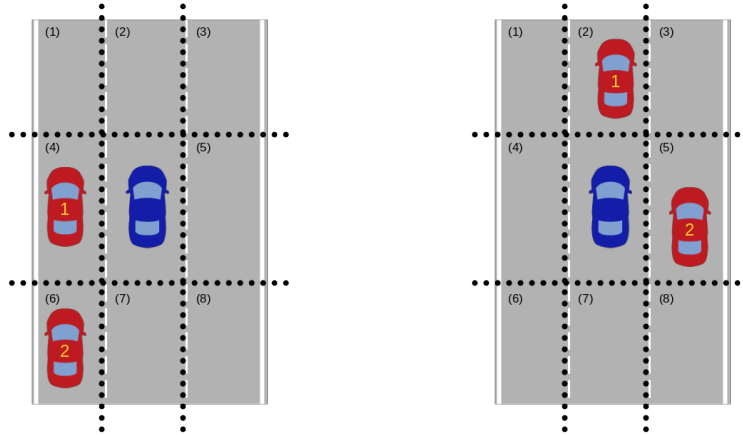


Figure 3: Example traffic situations for constructing abstract scenarios, specifying positions of non-ego cars (in red) in terms of the occupation of cells of an abstract 3x3 grid centered on the ego car (in blue).

of cells of an abstract 3x3 “grid” centered on the ego. Examples of the possible configurations that can be expressed in this way are shown in Fig. 3. We then define an *abstract scenario* as a combination of constraints about the different positions of the non-ego cars on the grid at different points in time. More specifically, each abstract scenario is specified as an LTL property of the following form:

$$\neg \mathbf{F}(\text{car1\_grid\_pos} = \text{CELL\_A1} \wedge \text{car2\_grid\_pos} = \text{CELL\_A2} \wedge \mathbf{X}(\mathbf{F}(\text{car1\_grid\_pos} = \text{CELL\_B1} \wedge \text{car2\_grid\_pos} = \text{CELL\_B2}))), \quad (1)$$

(where  $\text{cari\_grid\_pos}$  encodes the position of the  $i$ -th non-ego car in the grid and  $\text{CELL\_*}$  represent possible target positions for the cars.) By asking the model checker to find a counterexample to Eq. 1, we generate traces in which the non-ego cars first reach the configuration with  $\text{car1}$  in position A1 and  $\text{car2}$  in position A2, and then subsequently move to the configuration with  $\text{car1}$  in position B1 and  $\text{car2}$  in position B2, performing the necessary maneuvers while avoiding collisions with each other or with the ego.

The space of scenarios that is being explored therefore consists of all the possible combinations of transitions from configurations of the non-ego cars in terms of their position in the grid defined above. Enumerating all of them would give 4096 scenarios. We define our coverage criterion by selecting a subset of *abstract scenarios of interest*, consisting of various combinations of the traffic situations that can be modeled by positioning the non-ego cars in the grid around the ego. In total, for the experiments, we defined 144 such interesting scenarios.

### 4.3 VIVAS Interface with CARLA

In order to generate a scenario for the CARLA simulator, the abstract counterexample trace generated by the model checker is parsed for relevant information to be fed as input to the simulator. As an interface to the simulator, we used the CARLA module ScenarioRunner. This provides a Python interface to specify the routes for the ego as well as complex traffic scenarios by defining the behavior of the non-ego agent(s). ScenarioRunner also allows for running CARLA on a specified map at a particular location, while the user is allowed to implement their own AI-based ego agent. Every state of the abstract scenario trace is concretized into the corresponding behavior of every non-ego agent. Each behavior is then specified in

```

1 # instantiate car1's sequential behaviors
2 car1_behavior = py_trees.composites.Sequence("car1_behavior")
3
4 # first behavior in the sequence for car1: to drive forward.
5 ## Atomic behavior WaypointFollower runs in parallel with an atomic trigger condition DriveDistance
6 drive_car1 = py_trees.composites.Parallel("Drive forward",
7     policy=py_trees.common.ParallelPolicy.SUCCESS_ON_ONE)
8 drive_car1.add_child(WaypointFollower(car1, speed=3, avoid_collision=False))
9 ## drive_car1 terminates when car1 reaches the given distance
10 drive_car1.add_child(DriveDistance(car1, distance=2.6))
11 car1_behavior.add_child(drive_car1)
12
13 # second behavior in the sequence for car1: to change lane
14 lane_change_car1 = py_trees.composites.Parallel("lane change",
15     policy=py_trees.common.ParallelPolicy.SUCCESS_ON_ONE)
16 ## change lanes to the left with a speed of 2 m/s
17 lane_change_car1.add_child(LaneChange(car1, speed=2, direction='left',
18     distance_same_lane=1.3, distance_other_lane=60, distance_lane_change=9))
19 lane_change_car1.add_child(DriveDistance(car1, distance=12))
20 car1_behavior.add_child(lane_change_car1)
21
22 # third behavior in the sequence for car1: to stand still for 1 second.
23 still_car1 = py_trees.composites.Parallel("standstill",
24     policy=py_trees.common.ParallelPolicy.SUCCESS_ON_ONE)
25 still_car1.add_child(WaypointFollower(car1, speed=0.0, avoid_collision=False))
26 still_car1.add_child(StandStill(car1, name="StandStill", duration=1.0))
27 car1_behavior.add_child(still_car1)

```

Figure 4: Excerpt of automatically generated Scenario runner code (in Python3) for a non-ego (car1) behavior. Comments (in green) explain the behavior tree.

Python to generate a behavior tree for each corresponding non-ego vehicle. The behavior trees of all the non-egos present in the environment are then run in parallel during the simulation.

We first parse the initial coordinates and lanes of all the non-egos relative to the ego to instantiate them on the map. In each behavior of a behavior tree, the corresponding non-ego has to drive at a certain speed for a certain distance, following the waypoints given by the map on the same lane it is instantiated on. Although the duration of each state transition in the abstract trace is 1 second, the non-egos may take longer to drive that particular distance in the CARLA simulator, due to potential mismatches between the symbolic model and the simulator models. In case the vehicle stands still for  $n$  states in the abstract trace, it stands still for  $n$  seconds in the concrete scenario once it comes to a halt.

In the symbolic model, lane changes occur in one time step, with zero lateral distance traveled (since lanes have no width in the symbolic world). However, we constrain the successive lane changes of the same car to be  $N$  steps apart<sup>3</sup> to model the fact that a lane change is not instantaneous overall. To concretize this particular state transition, the non-ego transverses 9m while changing lanes, with this behavior terminating after traveling a total distance of 12m for the next behavior in the tree to be instantiated. Below these values, lane changes were not possible in CARLA at the speed ranges the vehicles drive in our scenario. Note that similar to the symbolic model, a non-ego can change only one lane at a time, with inputs {left, right} meaning change lane to the left or to the right.

We leverage the behavior library of ScenarioRunner to write these atomic behaviors and trigger conditions. An example behavior tree for one non-ego (car1) with all the above explained three behaviors is shown in Fig. 4. Here, lines 6-11: drive straight forward for 2.6 m, with a speed of 3 m/s; lines 14-20: perform a lane change to the left, with a speed of 2 m/s, driving a total of 12 m within which 9 m is the distance traveled while changing lanes; lines 23-37: stand still for 1 second. We do not need to

<sup>3</sup>We used  $N = 6$  in our experiments.

extract ego’s behavior from the abstract trace, since it is expected to make decisions autonomously in the simulator. Only initial spawn position and destination need to be extracted for the AI-based agent to follow the route.

The concrete simulation traces are then mapped back to the abstract trace to measure coverage, to check if the same sequence of scenes was encountered in the concrete scenarios or not. In particular, a predicate map is defined to map the absolute positions of the non-egos in the concrete simulation trace to the abstract 3x3 grid shown in Fig. 3. Here, we show examples of mapping the positions of non-egos to the cell locations 1,4 and 8 of the abstract grid:

$$\begin{aligned}
 \text{CELL\_1} &: (\text{car\_i.lane} < \text{ego.lane}) \wedge (4 \leq |\text{car\_i.pos} - \text{ego.pos}| \leq 24) \wedge (\text{car\_i.pos} > \text{ego.pos}) \\
 \text{CELL\_4} &: (\text{car\_i.lane} < \text{ego.lane}) \wedge (|\text{car\_i.pos} - \text{ego.pos}| \leq 10) \\
 \text{CELL\_8} &: (\text{car\_i.lane} > \text{ego.lane}) \wedge (4 \leq |\text{car\_i.pos} - \text{ego.pos}| \leq 24) \wedge (\text{car\_i.pos} < \text{ego.pos}),
 \end{aligned} \tag{2}$$

where  $\text{car\_i.pos}$  is the longitudinal position (in meters) of the  $i$ -th car in the simulation trace (and similarly for  $\text{ego.pos}$ ). In this way, we define the boundaries of the cells on the abstract grid.

#### 4.4 Monitoring of Properties

As described above, the monitor component of VIVAS is used to determine whether the concrete system (simulator) satisfies the system-level formal specification. For the automotive application, the simulation output traces include sequences of all states and actions executed by the ego vehicle, along with the time evolution of other observable parameters, which must be checked for property satisfaction/violation. In this study, we primarily need to check whether the ego vehicle crashes with another vehicle in the environment. Since the ego always travels in its own lane, we limit the check for the case when the ego crashes with any non-ego in front of it in its own lane. We do this by leveraging the continuous data stream from the collision sensor mounted on the ego. The monitor is currently hard-coded for monitoring specifically the output of this sensor, i.e., it checks whether the ego crashes or not at any time step in the simulation trace. Along with the satisfaction/violation of this property, the positions of non-egos in the simulation trace are mapped back to the abstract grid, to measure the degree of concrete coverage as described in §3.

In the future, we plan to use a runtime monitor based on NuRV [7, 8], to check standard LTL properties on ego behavior, e.g., if ego brakes within  $n$  time-steps as soon as any non-ego comes in front within its safe driving distance, or if the lane change of another vehicle is detected by the perception component of the ego within  $m$  time steps.

## 5 Results

In this section, we report on our experimental evaluation of our instantiation of VIVAS for the ADS application using the CARLA simulator. We first describe the experimental setup in §5.1, including the choice of parameters for the vehicles and environment in the symbolic model and in the CARLA simulator, necessary to generate meaningful scenarios. We then present the results of the evaluation in §5.2 and discuss them in §5.3.

We ran the experiments on an Intel i7 with NVIDIA GeForce RTX 2080 8GB GPU. These are the minimum hardware requirements to run the simulations on the CARLA simulator with AI models. All the experiments take roughly 22 hours to complete. We used a timeout of 200 seconds for each abstract scenario generation. This timeout was never reached by the model checker during our experiments: on

average, model checking took less than 10 seconds for each instance. Rather, the performance bottleneck turned out to be the time to instantiate a scenario in CARLA and perform the simulations, which took 3 minutes on average.

The code and data necessary for reproducing our experiments are available at <https://es-static.fbk.eu/people/sgoyal/fmas23>.

## 5.1 Experimental setup

At the level of the symbolic model, we define here some fixed parameters for simple, but meaningful scenario generation:

$$\begin{aligned}
 \text{ego\_cruise\_speed} &= 5 && (m/s) \\
 \text{non\_ego\_speed} &= [0, 12] && (m/s) \\
 \text{max\_acceleration} &= 5.6 && (m/s^2) \\
 \text{max\_braking} &= -4.6 && (m/s^2) \\
 \text{safe\_distance} &= 7 && (m) \\
 \text{lanes} &= \{\text{left, center, right}\}
 \end{aligned} \tag{3}$$

All the scenarios that we generate consist of 2 non-ego agents and one ego agent, all of which start from the same longitudinal position  $x = 0$ , with ego in the center lane and 2 non-egos on each side of it. Note that fixing the initial positions would not make a difference to the abstract scenario generation since the acceleration, braking, and speed for the non-egos are picked non-deterministically by the model checker for every time step, while respecting the above bounds. Since the AI-based agent we use in CARLA can only drive at around  $5 m/s$ , we limit the ego agent cruising speed to the same. All the agents start from  $0 m/s$ , with ego reaching its cruising speed with `max_acceleration`. To avoid collisions, it brakes with `max_braking` to maintain at minimum the `safe_distance` with all the non-egos.

To improve the robustness of abstract scenario generation, we reduce the size of each cell in the abstract  $3 \times 3$  grid in the symbolic model (see Fig. 3) by 3m in each direction, compared to the grid we use for evaluation of the coverage on the simulator. The lower values of the relative distances here are chosen to specifically create situations where non-egos stay close to ego and challenge its perception and control system with their braking and lane-changing maneuvers.

To mimic the symbolic environment model, we instantiate the CARLA simulator on a section of a highway of Town06, with 5 straight lanes, with ego positioned on the center lane and two non-egos on each lane next to it, corresponding to the symbolic model. Left-most and right-most lanes are not used. To compensate for the mismatch between the vehicle dynamics in the symbolic model and CARLA simulator, we concretize the initial positions of the non-egos at:

$$x = \{-3.5, 0, 3.5\} m \tag{4}$$

i.e., the non-egos start at  $3.5 m$  behind, same level, and  $3.5 m$  ahead of the ego in their respective lanes in different simulation runs. All the vehicles start from  $0 m/s$ , as parsed from the abstract trace. In future work, we could also concretize further for the simulations with one non-ego ahead and the other one behind the ego, to check if it extends the coverage results.

CARLA provides the possibility to change weather conditions (e.g., rainy, cloudy, night, etc.) at the beginning of simulation runs. We perform all the simulations in "clear noon" setting, for the ego's perception components to operate in the least challenging conditions. In future work, we will evaluate the AI solution in different weather conditions.

Table 1: Coverage Results.

Non-ego position	Total Scenarios	Coverage OK	Property FAIL	Coverage OK $\cap$ Property FAIL
3.5 m	144	68	3	2
0	144	73	17	4
- 3.5 m	144	45	54	19
<b>Set Union</b>		<b>81</b>	<b>61</b>	<b>25</b>

## 5.2 Evaluation

The model checker produces a total of 144 abstract scenarios based on the coverage criteria given in §4.2. Each abstract scenario is concretized into three concrete scenarios, by varying the initial positions of the non-egos according to Eq. 4, which gives us  $144 * 3 = 432$  concrete scenario outputs from the simulator. The evaluation results are shown in Table 1. The columns have the following meanings:

**Coverage OK:** Each point of the grid of coverage criteria represents a scenario with a fixed order of scenes. The concrete simulation run passes ("OK") if the abstract scenario generated by the model checker could indeed be generated on the simulator as well.

**Property FAIL:** The system-level property fails if the autonomous ego agent collides with at least one non-ego in front. We do not take into account the situations where non-egos crash into each other or hit the ego from behind.

**Coverage OK  $\cap$  Property FAIL:** the intersection of the above two conditions. These are the set of "interesting" cases (along with the other cases where property failed), where the coverage criteria passed, but the ego crashed with a non-ego in front.

**Set Union:** Combines the results for all concrete scenarios with respect to the abstract scenarios.

## 5.3 Analysis

All intended 144 abstract counterexamples could be generated by the model checker, meaning that the configurations we defined for the coverage did not violate any constraints in the symbolic model. As we see from the obtained results, not all the abstract scenarios generated by the model checker could be covered in the simulator. We could cover a total of only 81 out of 144 scenarios. This is primarily due to the mismatch in (a) behavior models and (b) vehicle dynamics, between the symbolic model and the simulator.

**Behavior model mismatch.** Since the ego is based on AI models, its non-deterministic behavior is not fully represented in the symbolic model. Ego speed is always varying within  $\pm 1$  m/s compared to the constant cruising speed in the symbolic model. The bounds of cells in the predicate map are defined for the relative position of the non-egos with respect to ego, as specified in Eq. 2. Hence, in some cases, the non-egos can not reach the required region within these bounds during the scenario, since the ego is traveling too fast or slow. In principle, we could overcome this by conditioning the non-egos' behavior to the ego's in terms of the distance traveled relative to the ego instead of the absolute distance on the lane, while translating the abstract scenario to the concrete one. However, no such atomic behaviors or conditions yet exist in ScenarioRunner. This could be included in one of the future works.

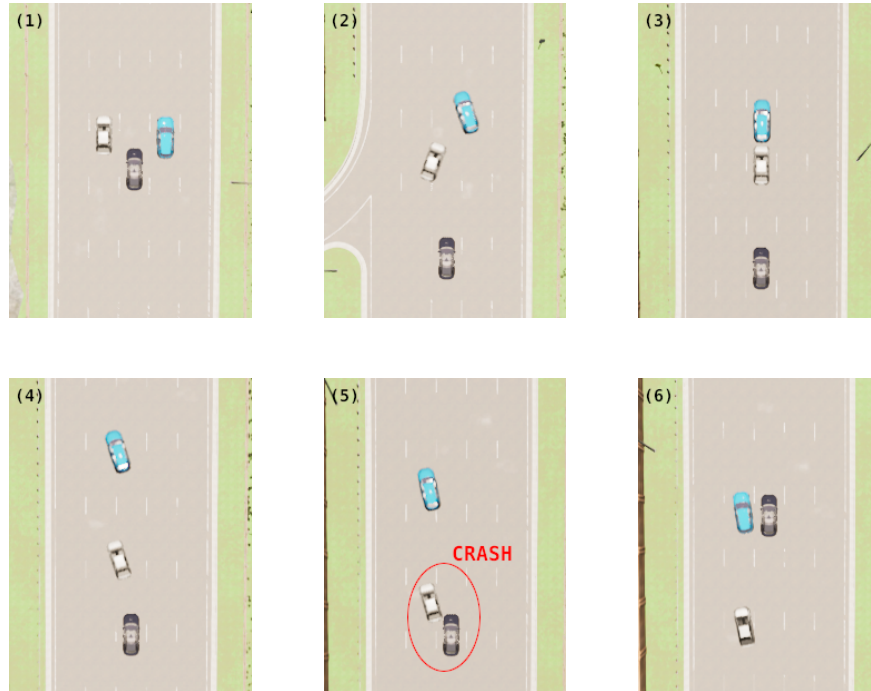


Figure 5: Scenario defined in Eq. 5, with Coverage OK  $\cap$  Property FAIL. Ego crashing with a non-ego in front (scene 5).

**Vehicle dynamics mismatch.** The vehicle dynamics models are based on OEM data, hard-coded in the simulator. The physical constraints for maximum acceleration and braking in the symbolic model are instead generic. However, we found a big mismatch during the simulations, with maximum acceleration values of vehicle models in CARLA reaching as high as  $12 \text{ m/s}^2$ , and maximum braking below  $-15 \text{ m/s}^2$ . Since these parameter values also vary from vehicle to vehicle in CARLA, we even came across situations where non-egos crashed into each other while changing lanes.

**Interesting scenarios.** Even though we could not cover all the abstract scenarios in the concrete simulator, there were 61 scenarios where AI-based ego collided with a non-ego agent in front. In particular, we obtained 25 interesting scenarios that met their abstract specification, but with ego crashing into a non-ego in front. Fig. 5 shows 6 scenes (in temporal order of 1-6) extracted from one such scenario. This corresponds to the LTL property specified below in Eq. 5 (with reference to Eq. 1):

$$\neg \mathbf{F}(\text{car1\_grid\_pos} = 2 \wedge \text{car2\_grid\_pos} = 2 \wedge \mathbf{X}(\mathbf{F}(\text{car1\_grid\_pos} = 6 \wedge \text{car2\_grid\_pos} = 4))) \quad (5)$$

Here, the grid positions correspond to the cell numbers mentioned in Fig. 3 for the abstract grid space. We now describe the scenario in Fig. 5.

- All the agents are initialized at  $0 \text{ m/s}$ , with car1 on the left of ego, and car2 to its right, with both cars starting  $3.5 \text{ m}$  ahead of ego (scene 1). Here, “car1\_grid\_pos = 4  $\wedge$  car2\_grid\_pos = 5” on the abstract grid.
- The non-egos travel faster than the ego to change lanes in front of it (scene 2), and end up in the configuration (scene 3) where the first predicate, “car1\_grid\_pos = 2  $\wedge$  car2\_grid\_pos = 2” is satisfied.



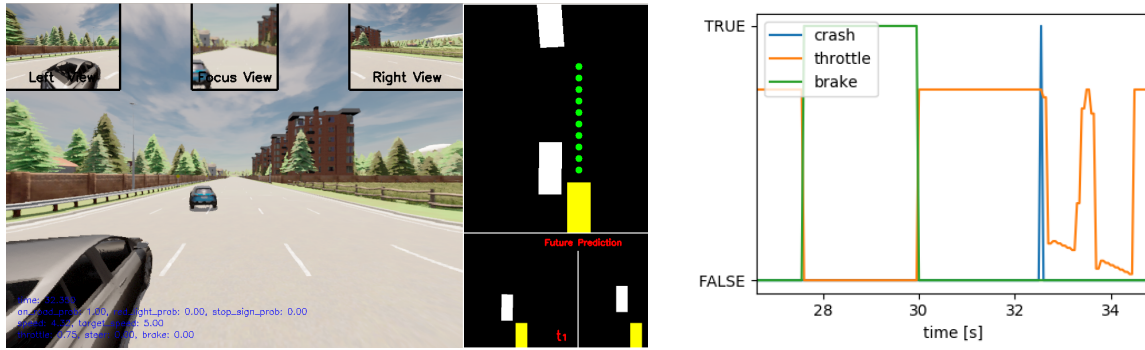


Figure 6: Ego collision; (left): Different camera views of AI-based ego, with the perception component's output; (right): Plot showing ego's throttle and braking values during the collision (in blue).

- The non-egos then change lanes to the left slowly (scene 4), when the ego crashes with car1 (scene 5), even when the car1 has still not completed the lane change.

Fig. 6 (left) shows the front view of the AI agent at the instant of crashing with non-ego. The real-time telemetry shows that ego's brake = 0 and throttle = 0.75. The perception component here seems to mis-detect the actual position of car1. The corresponding output from the simulation trace is shown in the right plot, with ego's throttle = True (and brake = False) for more than 2 seconds leading to the collision.

- The ego keeps driving forward when the second predicate, “ $\text{car1\_grid\_pos} = 6 \wedge \text{car2\_grid\_pos} = 4$ ” is satisfied (scene 6).

## 6 Conclusions

This paper showed the application of system-level simulation-based verification of ADS using formal methods to generate abstract scenarios. The verification toolchain includes nuXmv for model checking and generating abstract scenarios, CARLA for simulating the ego behaviors in concrete scenarios, and mappings from abstract to concrete scenarios and back. We presented an abstract model of the system and a coverage criterion that allows the automated generation of abstract scenarios with model checking. The generated abstract scenarios cover different sequences of traffic scenes that are relevant to test the reaction of the ego's behavior to see if it avoids crashing into other cars. The simulation with CARLA of the corresponding concrete scenarios showed various crashes caused by the ego, although not all simulations reproduce the expected abstract scenario. Inspecting some of the simulations reporting a crash in a covered abstract scenario confirms that the ego behavior is indeed buggy and this is probably due to the AI-based perception component.

During this study, we gained many insights that may lead to some interesting future research directions. These include more efficient techniques to generate abstract scenarios for minimizing the number of model checking runs needed to achieve a certain coverage level; the integration of effective sampling techniques that synthesize various simulation parameters for the same abstract scenario; extending the abstract model by incorporating uncertainty in the ego behavior or a more precise representation of the continuous-time behavior with timed or hybrid version of SMV [5, 6]; finally, enhancing the concrete scenario specification with conditional behaviors of non-ego vehicles that react to the choices of the ego.

## References

- [1] Association for Standardization of Automation & Measuring Systems (ASAM): *OpenSCENARIO*. <https://www.asam.net/standards/detail/openscenario/>. Accessed: 2023-08-30.
- [2] Marco Bozzano, Riccardo Bussola, Marco Cristoforetti, Srajan Goyal, Martin Jonáš, Konstantinos Kapellos, Andrea Micheli, Davide Soldà, Stefano Tonetta, Christos Tranoris & Alessandro Valentini (2023): *RobDT: AI-enhanced Digital Twin for Space Exploration Robotic Assets*. In: *The Use of Artificial Intelligence for Space Applications*, Springer Nature Switzerland, pp. 183–198, doi:10.1007/978-3-031-25755-1\_12.
- [3] Roberto Cavada, Alessandro Cimatti, Michele Dorigatti, Alberto Griggio, Alessandro Mariotti, Andrea Micheli, Sergio Mover, Marco Roveri & Stefano Tonetta (2014): *The nuXmv Symbolic Model Checker*. In: *CAV, Lecture Notes in Computer Science 8559*, Springer, pp. 334–342, doi:10.1007/978-3-319-08867-9\_22.
- [4] D. Chen & P. Krahenbuhl (2022): *Learning from All Vehicles*. In: *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE Computer Society, Los Alamitos, CA, USA, pp. 17201–17210, doi:10.1109/CVPR52688.2022.01671.
- [5] Alessandro Cimatti, Alberto Griggio, Enrico Magnago, Marco Roveri & Stefano Tonetta (2019): *Extending nuXmv with Timed Transition Systems and Timed Temporal Properties*. In Isil Dillig & Serdar Tasiran, editors: *Computer Aided Verification*, Springer International Publishing, Cham, pp. 376–386, doi:10.1007/978-3-030-25540-4\_21.
- [6] Alessandro Cimatti, Alberto Griggio, Sergio Mover & Stefano Tonetta (2015): *HyComp: An SMT-Based Model Checker for Hybrid Systems*. In: *TACAS, Lecture Notes in Computer Science 9035*, Springer, pp. 52–67, doi:10.1007/978-3-662-46681-0\_4.
- [7] Alessandro Cimatti, Chun Tian & Stefano Tonetta (2019): *Assumption-Based Runtime Verification with Partial Observability and Resets*. In: *RV, Lecture Notes in Computer Science 11757*, Springer, pp. 165–184, doi:10.1007/978-3-030-32079-9\_10.
- [8] Alessandro Cimatti, Chun Tian & Stefano Tonetta (2019): *NuRV: A nuXmv Extension for Runtime Verification*. In Bernd Finkbeiner & Leonardo Mariani, editors: *Runtime Verification*, Springer International Publishing, Cham, pp. 382–392, doi:10.1007/978-3-030-32079-9\_23.
- [9] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez & Vladlen Koltun (2017): *CARLA: An Open Urban Driving Simulator*. In Sergey Levine, Vincent Vanhoucke & Ken Goldberg, editors: *Proceedings of the 1st Annual Conference on Robot Learning, Proceedings of Machine Learning Research 78*, PMLR, pp. 1–16, doi:10.48550/arXiv.1711.03938.
- [10] Tommaso Dreossi, Daniel J. Fremont, Shromona Ghosh, Edward Kim, Hadi Ravanbakhsh, Marcell Vazquez-Chanlatte & Sanjit A. Seshia (2019): *VerifAI: A Toolkit for the Formal Design and Analysis of Artificial Intelligence-Based Systems*. In Isil Dillig & Serdar Tasiran, editors: *Computer Aided Verification*, Springer International Publishing, Cham, pp. 432–442, doi:10.1007/978-3-030-25540-4\_25.
- [11] O. foretellix: *Open M-SDL*. [https://releases.asam.net/OpenSCENARIO/2.0-concepts/M-SDL\\_LRM\\_OS.pdf](https://releases.asam.net/OpenSCENARIO/2.0-concepts/M-SDL_LRM_OS.pdf). Accessed: 2023-08-07.
- [12] Simone Fratini, Patrick Fleith, Nicola Policella, Alberto Griggio, Stefano Tonetta, Srajan Goyal, Thi Thieu Hoa Le, Jacob Kimblad, Chun Tian, Konstantinos Kapellos, Christos Tranoris & Quirien Wijnands (2023): *Verification and Validation of Autonomous Systems with Embedded AI: The VIVAS Approach*. In: *ASTRA*, p. To appear. Available at <https://az659834.vo.msecnd.net/eventsairwesteuprod/production-atpi-public/070740b67e5b4a32a9be94228c9ac40d>.
- [13] Daniel J. Fremont, Edward Kim, Tommaso Dreossi, Shromona Ghosh, Xiangyu Yue, Alberto L. Sangiovanni-Vincentelli & Sanjit A. Seshia (2022): *Scenic: a language for scenario specification and data generation*. *Machine Learning*, doi:10.1007/s10994-021-06120-5.
- [14] Xiaowei Huang, Marta Kwiatkowska, Sen Wang & Min Wu (2017): *Safety Verification of Deep Neural Networks*. In Rupak Majumdar & Viktor Kunčák, editors: *Computer Aided Verification*, Springer International Publishing, Cham, pp. 3–29, doi:10.1007/978-3-319-63387-9\_1.

- [15] Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian & Mykel J. Kochenderfer (2017): *Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks*. In: *CAV (1), Lecture Notes in Computer Science* 10426, Springer, pp. 97–117, doi:10.1007/978-3-319-63387-9\_5.
- [16] Guy Katz, Derek A. Huang, Duligur Ibeling, Kyle Julian, Christopher Lazarus, Rachel Lim, Parth Shah, Shantanu Thakoor, Haoze Wu, Aleksandar Zeljic, David L. Dill, Mykel J. Kochenderfer & Clark W. Barrett (2019): *The Marabou Framework for Verification and Analysis of Deep Neural Networks*. In: *CAV (1), Lecture Notes in Computer Science* 11561, Springer, pp. 443–452, doi:10.1007/978-3-030-25540-4\_26.
- [17] Moritz Klischat & Matthias Althoff (2020): *Synthesizing Traffic Scenarios from Formal Specifications for Testing Automated Vehicles*. In: *IV, IEEE*, pp. 2065–2072, doi:10.1109/IV47402.2020.9304617.
- [18] Rupak Majumdar, Aman Mathur, Marcus Pirron, Laura Stegner & Damien Zufferey (2021): *Paracosm: A Test Framework for Autonomous Driving Simulations*. In Esther Guerra & Mariëlle Stoelinga, editors: *Fundamental Approaches to Software Engineering*, Springer International Publishing, Cham, pp. 172–195, doi:10.1007/978-3-030-71500-7\_9.
- [19] Corina S. Păsăreanu, Ravi Mangal, Divya Gopinath, Sinem Getir Yaman, Calum Imrie, Radu Calinescu & Huafeng Yu (2023): *Closed-Loop Analysis of Vision-Based Autonomous Systems: A Case Study*. In Constantin Enea & Akash Lal, editors: *Computer Aided Verification*, Springer Nature Switzerland, Cham, pp. 289–303, doi:10.1007/978-3-031-37706-8\_15.
- [20] Clément Robert, Jérémie Guiochet, Héléne Waeselynck & Luca Vittorio Sartori (2021): *TAF: a Tool for Diverse and Constrained Test Case Generation*. In: *2021 IEEE 21st International Conference on Software Quality, Reliability and Security (QRS)*, pp. 311–321, doi:10.1109/QRS54544.2021.00042.
- [21] Luca Vittorio Sartori, Héléne Waeselynck & Jérémie Guiochet (2023): *Pairwise Testing Revisited for Structured Data With Constraints*. In: *ICST, IEEE*, pp. 199–209, doi:10.1109/ICST57152.2023.00027.
- [22] Hao Shao, Letian Wang, Ruobing Chen, Hongsheng Li & Yu Liu (2023): *Safety-Enhanced Autonomous Driving Using Interpretable Sensor Fusion Transformer*. In Karen Liu, Dana Kulic & Jeff Ichnowski, editors: *Proceedings of The 6th Conference on Robot Learning, Proceedings of Machine Learning Research* 205, PMLR, pp. 726–737, doi:10.48550/arXiv.2207.14024.
- [23] Gagandeep Singh, Timon Gehr, Markus Püschel & Martin Vechev (2019): *An Abstract Domain for Certifying Neural Networks*. *Proc. ACM Program. Lang.* 3(POPL), doi:10.1145/3290354.
- [24] CARLA Team: *CARLA Autonomous Driving Leaderboard*. <https://leaderboard.carla.org/leaderboard/>. Accessed: 2023-08-30.
- [25] CARLA Team: *CARLA ScenarioRunner*. <https://carla-scenariorunner.readthedocs.io>. Accessed: 2023-08-30.
- [26] Cumhuri Erkan Tuncali, Georgios Fainekos, Hisahiro Ito & James Kapinski (2018): *Sim-ATAV: Simulation-Based Adversarial Testing Framework for Autonomous Vehicles*. In: *Proceedings of the 21st International Conference on Hybrid Systems: Computation and Control (Part of CPS Week), HSCC '18*, Association for Computing Machinery, New York, NY, USA, p. 283–284, doi:10.1145/3178126.3187004.
- [27] Eric Vin, Shun Kashiwa, Matthew Rhea, Daniel J. Fremont, Edward Kim, Tommaso Dreossi, Shromona Ghosh, Xiangyu Yue, Alberto L. Sangiovanni-Vincentelli & Sanjit A. Seshia (2023): *3D Environment Modeling for Falsification and Beyond with Scenic 3.0*. In Constantin Enea & Akash Lal, editors: *Computer Aided Verification*, Springer Nature Switzerland, Cham, pp. 253–265, doi:10.1007/978-3-031-37706-8\_13.
- [28] Hermann Winner, Karsten Lemmer, Thomas Form & Jens Mazzega (2019): *PEGASUS—First Steps for the Safe Introduction of Automated Driving*. In Gereon Meyer & Sven Beiker, editors: *Road Vehicle Automation 5*, Springer International Publishing, Cham, pp. 185–195, doi:10.1007/978-3-319-94896-6\_16.
- [29] Penghao Wu, Xiaosong Jia, Li Chen, Junchi Yan, Hongyang Li & Yu Qiao (2022): *Trajectory-guided Control Prediction for End-to-end Autonomous Driving: A Simple yet Strong Baseline*. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho & A. Oh, editors: *Advances in Neural Information Processing Systems*, 35, Curran Associates, Inc., pp. 6119–6132, doi:10.48550/arXiv.2206.08129.

# Enforcing Timing Properties in Motorway Traffic

Christopher Bishopink

Carl von Ossietzky University Oldenburg, Oldenburg

bischopink@informatik.uni-oldenburg.de

In previous work [4], we proposed a *Runtime Enforcement* Approach to deal with timing properties in motorway traffic, which are present in form of *Timed Multi-Lane Spatial Logic (TMLSL)* formulae, a logic tailored to express both spatial and timing properties. Employing communication between the cars, we utilised a nondeterministic controller “guessing” which actions to execute next for each car, before asking the local monitors of the cars for permission to execute the announced actions. In this contribution, we consider a more reasonable controller that only considers sequences that satisfy its own properties. This is done utilising region automata that one can generate from the cars’ specifications. In the approach, we also came along a minor decidability result for TMLSL.

## 1 Introduction

With the number of (at least partially) autonomous cars increasing on the roads around the globe, challenges and advantages in the specification and verification of their behaviour occur. If one assumes that the cars are able to communicate with each other, a more detailed interplay between them is possible than with human drivers and allows finding solutions for complicated traffic situations that human drivers could easily miss.

The roads we consider here are motorways, formalised as *traffic snapshots* [9] with a logic to reason about them called *Multi-Lane Spatial Logic (MLSL)* [9]. We extended this spatial logic towards *Timed Multi-Lane Spatial Logic (TMLSL)* to also cover the timing aspect of a car’s specification in [3]. Based on TMLSL, we proposed a runtime-enforcement approach in [4], employing a nondeterministic controller that asked for the permission of other cars for the actions it wants to execute. Due to the nondeterminism, completely unreasonable sequences that even violate the own car’s specification could be announced. In the lack of a result that allow announcing/checking only reasonable sequences, the nondeterministic controller still allowed us to show that the approach is complete.

In this work, we propose a more reasonable approach, utilising the region automaton of the cars’ specifications. Still, all cars announce sequences they want to achieve, but this time all announced sequences would yield satisfying runs at least for the car that announces them. The announced sequences are then checked by a central entity, e.g. a road-side unit, for a run that is satisfying for all cars and informs the cars accordingly. We furthermore present a minor decidability result, eliminating one of the causes for the semi-decidability of TMLSL [3].

**Related Work** In the context of MLSL, different topologies have been explored in addition to motorway traffic, namely country roads [8] and urban traffic [17], as well as their satisfaction problems [13][6] and controllers for cars in these topologies with different desirable properties such as liveness and fairness [5, 18]. Other approaches in the context of autonomous or automated driving systems use e.g. *differential dynamic logic* [12] or a specification with extended types of timed automata [10]. These approaches mostly concentrate on a top-level view of the system under control. A more technical view of the evolution of a cars dynamics in an adaptive cruise control setting is e.g. given in [1]. Runtime

Enforcement [16] and runtime verification [7] are also well studied topics. To the best of our knowledge, runtime enforcement approaches however are more intensively considered in more restricted settings than motorway traffic, where the system evolves quiet dynamically to the input given. Another feature is that in our case, the input and output of the system under control are different from each other.

## 2 Preliminaries

In this Section, we introduce the formal concepts our approach is build on. We start with the model of motorway traffic, its logic and evolution in Sect. 2.1 and continue with the timing model used, Sect.2.2. The combination of them is called TMLSL and covered in Sect. 2.3.

### 2.1 Spatial Model of and Logic for Motorway Traffic

**Model** The spatial model we use in the setting of motorway traffic was introduced in [9]. It allows only traffic in one direction on a fixed set of lanes  $\mathbb{L} = \{1, \dots, N\}$  with an infinite extension each. On these lanes, cars from the set of car identifiers  $\mathbb{I} = \{A, B, \dots\}$  drive, each car  $C$  of them with a certain speed  $spd(C)$ , acceleration  $acc(C)$  and position on the lane,  $pos(C)$ . There are two different types of occupation a car can have on lanes, either a reservation  $res : \mathbb{I} \rightarrow \mathcal{P}(\mathbb{L})$ , the space it physically occupies (multiple lanes if it is changing lanes at this moment) or a claim  $clm : \mathbb{I} \rightarrow \mathcal{P}(\mathbb{L})$ , the lane a car wishes to change to, which is the equivalent of setting a turn signal. Altogether, this information is represented as a traffic snapshot  $TS = (res, clm, pos, spd, acc)$ .

In a traffic snapshot, there is no information present what the sizes of the cars and their braking distances are, as  $pos$  only stores the rear end of each vehicle. We neglect the concept of a *Sensor Function* here that makes this information available to us and simply assume that size and braking distance of each car is known. Also omitted is the *View* that allows to only consider a finite extension of the infinite extension of a traffic snapshot when evaluating formulae. A graphical representation of three traffic snapshots is depicted in Fig. 1, where we also omitted showing concrete values for the position, speed and acceleration of the cars.

As already hinted at, a traffic snapshot describes the situation on the road at one point in time only. A situation on the road may evolve, which is handled in the model via *transitions*.

**Transitions** We divide the set of transitions usable in a traffic snapshot into transitions regarding the discrete behaviour between lanes and transitions regarding the continuous behaviour along the lanes. The first set consists of car  $C$  claiming a lane  $n$  resp. withdrawing all claims ( $c(C, n)/wd\ c(C)$ ) and car  $C$  reserving a lane resp. withdrawing all reservations except the one on lane  $n$  ( $r(C)/wd\ r(C, n)$ ). The second set is the one we focus more on, as it is considered more intensively in what follows, they handle the change of a car's acceleration to some value  $a$  ( $acc(C, a)$ ) as well as the passing of  $t$  time units ( $t$ ). In

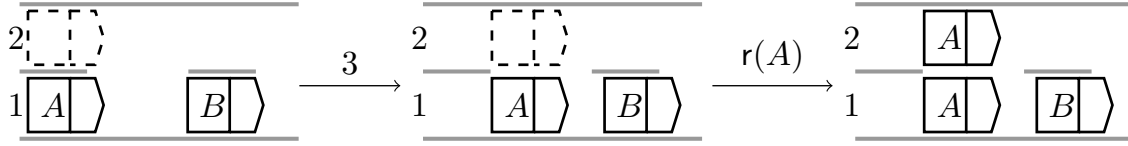


Figure 1: A transition sequence with three traffic snapshots and two actions, taking place with two cars on two lanes. Car A is faster than car B, so it comes closer when  $t = 3$  time units pass. Afterwards, it reserves its formerly claimed space (dashed copy of it on the neighbouring lane), one step further in an attempt to overtake car B.

the following definition,  $\oplus$  is the overriding operator of Z [19]:

$$\begin{aligned}
 \mathcal{TS} \xrightarrow{t} \mathcal{TS}' &\Leftrightarrow \mathcal{TS}' = (res, clm, pos', spd', acc) \\
 &\wedge \forall C \in \mathbb{I} : pos'(C) = pos(C) + spd(C) \cdot t + \frac{1}{2} acc(C) \cdot t^2 \\
 &\wedge \forall C \in \mathbb{I} : spd'(C) = spd(C) + acc(C) \cdot t \\
 \mathcal{TS} \xrightarrow{acc(C, a)} \mathcal{TS}' &\Leftrightarrow \mathcal{TS}' = (res, clm, pos, spd, acc') \\
 &\wedge acc' = acc \oplus \{C \mapsto a\},
 \end{aligned}$$

Over the set of actions, which are the transitions without the one where only time passes, we define *timed words of actions*  $\omega = \langle (\alpha_0, t_0), \dots, (\alpha_n, t_n) \rangle$ , with  $\alpha_i$  an action and  $\langle t_0, \dots, t_n \rangle$  a real-time sequence. For a formal account, we refer to [14].

A graphical representation of a transition sequence including two transitions (one discrete and one continuous) is shown in Fig. 1. It can be interpreted as the timed word  $\omega_1 = \langle (r(A), 3) \rangle$ .

**Remark 1** (Dynamic behaviour of the cars). *The model used for describing the dynamics of the cars is a quiet simple one, ignoring many difficulties that one would encounter in the real world, such as friction or variable acceleration capabilities based on the current speed. Still, we believe that it is a good abstraction of the real world's dynamics. Especially if one considers that the positions (plus size and braking distance) could be over-approximations, this allows some degree of freedom in achieving a behaviour to match the correct positions.*

**Logic** To reason about the traffic situations formalised as traffic snapshots, the logic MLSL [9] was introduced. Here, we consider a variant of MLSL called MLSLS (MLSL *with scopes*) [6], limiting the range of cars over which car variables are evaluated to a finite range. Formulae of MLSLS are constructed according to the grammar

$$\varphi ::= \gamma = \gamma' \mid free \mid re(\gamma) \mid cl(\gamma) \mid l = k \mid \exists c. \varphi \mid \varphi_1 \sim \varphi_2 \mid \begin{array}{l} \varphi_1 \\ \varphi_2 \end{array} \mid cs : \varphi,$$

and standard Boolean combinations of such formulae, where  $\gamma$  and  $\gamma'$  are car variables,  $k \in \mathbb{R}$  and  $cs$  is a (sub-)set of car variables.

The atoms that formulae are constructed from are the comparison of two car variables  $\gamma = \gamma'$ , *free* denotes a segment with free space that is not occupied by others, the reservation  $re(\gamma)$  of a car  $\gamma$ , the

claim  $cl(\gamma)$  of a car  $\gamma$  and the comparison  $l = k$  of the length  $l$  of the considered segment against some value  $k$ . With  $\exists c.\varphi$  one asks for the existence of a car  $c$  that satisfies  $\varphi$ . The horizontal chop operator ( $\varphi_1 \sim \varphi_2$ ) is used to determine if it is possible to divide the current segment into two parts along the lanes s.t. in the first part  $\varphi_1$  and in the part directly ahead of it  $\varphi_2$  holds. The same can be specified with the vertical chop operator  $\left( \begin{array}{c} \varphi_1 \\ \varphi_2 \end{array} \right)$ , but the point to divide is in between two lanes this time.  $cs : \varphi$  limits the scope over which  $\varphi$  is evaluated to the finite domain  $cs$  and effects only formulae that use quantification over the cars in their semantics and thus only  $\exists c.\varphi$  and *free*. In Sect. 3, we focus on length comparisons  $l = k$  and simply abbreviate them as  $\theta$ , as there can be multiple of them regarding the same segment.

A common abbreviation used is the *somewhere* modality  $\langle \varphi \rangle$ , expressing that there is a partition on the road along and in between the lanes s.t.  $\varphi$  holds in some point of the partition.

For a formal definition, especially about the exact semantics of MLSLS formulae, we refer the reader to [6], but would like to point out that MLSLS is, in contrast to pure MLSL, decidable.

## 2.2 Model and Logic of Time

The logic we consider here for the timing aspects is called *State-Clock Logic (SCL)* [15]. Formulae from this logic are constructed over an alphabet of propositions  $\Sigma$  according to the grammar

$$\psi ::= p \mid \psi_1 \vee \psi_2 \mid \neg\psi \mid \psi_1 \mathcal{U} \psi_2 \mid \psi_1 \mathcal{S} \psi_2 \mid \triangleright_{\sim c} \psi \mid \triangleleft_{\sim c} \psi,$$

with  $\sim \in \{<, \leq, =, \geq, >\}$  and  $p \in \Sigma$ .

Apart from well-known Boolean combinations of formulae and the usual until- and since operators ( $\mathcal{U}$  resp.  $\mathcal{S}$ ), SCL allows to measure the time since ( $\triangleleft$ )/ until ( $\triangleright$ ) a formula  $\psi$  held/holds for the last/next time and compare this difference with  $\sim c$ .

The semantics of SCL formulae is evaluated on (usually infinite) *timed sequences of states*  $m = \langle (s_0, I_0), (s_1, I_1), \dots \rangle$  with  $s_i \subseteq \Sigma$  and  $\langle I_0, I_1, \dots \rangle$  a monotonically increasing sequence of adjacent intervals. Intuitively, the formula  $\triangleright_{\sim c} \psi$  holds at time point  $t$  in the  $i$ th state of  $m$ , written  $(m, i, t) \models_{\sim c} \psi$  iff there is some state  $(s_j, I_j)$  at position  $j > i$  where  $\psi$  holds, all states in between  $i$  and  $j$  do not satisfy  $\psi$  and the difference between the left border of the  $j$ th interval  $I_j$  and  $t$  satisfies  $\sim c$ . The analogous applies for the operator  $\triangleleft_{\sim c} \psi$ , the semantics of the remaining operators is as expected. Example for both syntax and semantics are given in the next section, for a more formal account on the topic we refer the reader to [15].

The decidability problem of SCL is known to be decidable. Given a SCL formula  $\psi$ , one can construct a *State-Clock (SC) Automaton*  $A_\psi = (\mathcal{P}, C_\mathcal{P}, L, L_0, E, \mathcal{L}, \Delta, \mathcal{F})$ , with  $\mathcal{P}$  the set of propositions used,  $C_\mathcal{P}$  the set of clocks,  $L$  and  $L_0$  the (initial) locations of  $A_\psi$ , a transition relation  $E \subseteq L \times L$ , a labelling function  $\mathcal{L}$  assigning the propositions valid in it to every location of  $L$ , another labelling function  $\Delta$  assigning constraints over  $C_\mathcal{P}$  to every location of  $L$  and a family of Büchi acceptance sets  $\mathcal{F}$ . From there on a *Region-Automaton*  $\mathcal{R}(A_\psi)$  [15][2] can be constructed. Having operators to compare the time to for both the future and the past, SC automata have a history clock  $x_p$  and a prophecy clock  $y_p$  for each proposition  $p \in \mathcal{P}$  in  $C_\mathcal{P}$ . Both types of clocks need to be respected when constructing the region automaton. A *region*  $[v]$  describes a class of clock valuations  $v$  that cannot be distinguished by any SC automaton and can be represented as a set of (in-)equalities over the set of clock variables  $x_y, y_p$  and the natural numbers  $\mathbb{N}$ . Iff the language of  $\mathcal{R}(A_\psi)$  is not empty, the formula  $\psi$  is satisfiable. In [4], we extended State-Clock automata with *broadcast communication* like the timed automata of UPPAAL [11]. For a formal account on the broadcast communication used we refer to [17] and only point out that sending some data  $d$  over a channel  $c$  is denoted as  $c!\langle d \rangle$  and receiving this data on the same channel is denoted as  $c?\langle d \rangle$ . We

also allowed simple functions dealing with data structures and simple computations on the transitions of communicating SC automata.

### 2.3 TMLSL

To express and reason about both spatial and timing properties in motorway traffic, we introduced *TMLSL* [3]. The idea of this logic is to use *MLSLS*-formulae as the propositions that *SCL* formulae are build from. The intuitive idea for the semantics is that a traffic snapshot  $TS$  with a timed word of actions  $\omega$  is a model for a formula  $\varphi$ ,  $TS_0, \omega \models \varphi$  iff there is a timed sequence of states  $m(TS_0, \omega)$  that is propositionally consistent and complete in the subformulae of  $\varphi$ , describes the evolution of  $TS$  along  $\omega$  and is a model of  $\varphi$  in the *SCL*-semantics,  $m(TS, \omega) \models_{SCL} \varphi$ . We now give an example for a *TMLSL* formula as well as their satisfaction. For simplicity, we use car identifiers instead of car variables in the *MLSLS* propositions in the formula.

**Example 1 (TMLSL).** Consider the *TMLSL* formula

$$\varphi_1 = \langle re(A) \sim free \wedge l = 21 \sim re(B) \rangle \implies \triangleright_{=5} \langle re(A) \sim free \wedge l = 15 \sim re(B) \rangle$$

that specifies that when the distance between the reservations of the two cars  $A$  and  $B$  is equal to 21 distance units somewhere in the traffic snapshot, it needs to be equal to 15 distance units within exactly 5 time units. A satisfying sequence of states is

$$m = \langle (\langle re(A) \sim free \wedge l = 21 \sim re(B) \rangle, [0, 0]), (\langle re(A) \sim free \wedge l \leq 21 \wedge l > 15 \sim re(B) \rangle, (0, 5)), (\langle re(A) \sim free \wedge l = 15 \sim re(B) \rangle, [5, 7]) \rangle,$$

one that is not is

$$m' = \langle (\langle re(A) \sim free \wedge l = 21 \sim re(B) \rangle, [0, 0]), (\langle re(A) \sim free \wedge l \leq 21 \wedge l > 15 \sim re(B) \rangle, (0, 5]), (\langle re(A) \sim free \wedge l = 15 \sim re(B) \rangle, (5, 7)) \rangle,$$

as the distance between the two cars reached the value 15 too late. The only difference between  $m$  and  $m'$  are the shapes of the second and third interval. Please note that we omitted some subformulae of  $\varphi_1$  in the sequences in an attempt to keep them readable.

Ex. 2 gives values for the positions, speeds and accelerations of the cars in the first traffic snapshot  $TS$  of Fig. 1 and a timed word of actions  $\omega$  s.t. applying  $\omega$  to  $TS$  results in the timed sequence of states  $m(TS_0, \omega)$  that satisfies  $\varphi_1$  written  $TS, \omega \models \varphi_1$ , resp.  $m(TS_0, \omega) \models_{SCL} \varphi_1$ .

**Remark 2.** Using *MLSLS* formulae as the propositions of *SCL* rather than *MLSL* actions imposes some difficulties, still some situations on the road can only be described using formulae rather than action. Additionally, actions take zero time in the model, so one could argue that they are not observable from the outside. So for the cars on the motorway that we want to control, we have *MLSL* actions as the input to them, but the system produces evolutions of traffic snapshots as an output, which we observe through *MLSLS* formulae.

**Finite Semantics** In [4], we introduced a finite semantics for *TMLSL*, suited for the finite sequences that are usually available in runtime monitoring/enforcement and the reality on motorways. Intuitively, a finite word of actions  $\omega$  satisfies a formula  $\varphi$  in the finite semantics up to time  $t$ ,  $TS, \omega \models_t \varphi$  iff there exists at least one suffix  $m'$  s.t.  $m(TS, \omega).m' \models_{SCL} \varphi$  in the infinite semantics.



### 3 Decidability Results

Regarding the decidability of TMLSL, we point out that the logic is at least semi-decidable [3]. In answering this question, we considered maximum values on the acceleration of the cars (from  $acc_{min}$  to  $acc_{max}$ ) and the speed (from 0 to  $spd_{max}$ ), as in the real world there are (at least) physical bounds, too. We do the same here. The decidability results of SCL and MLSLS do not directly transfer to TMLSL, as we need to interpret the SCL-propositions, which are MLSLS-formulae, and the cars in the traffic snapshot may not be able to behave as specified in TMLSL. One cause for only semi-decidability are the actions regarding the dynamic behaviour along the lanes, the change of a car's acceleration and the passing of time. While it is easy to see what discrete actions need to be executed and when, given a timed sequence of states  $m$ , it was unknown how many acceleration changes are needed to achieve cars behaving correct with respect to the lengths constraints specified in  $m$ . In this paper, we give an algorithm that decides this question. Before doing so, however, we start with an example (adjusted example of [3]):

**Example 2.** *Consider a traffic snapshot with one lane and two cars A and B, B driving ahead of A, where the distance between the two cars is equal to 21 and both of them have a speed of 4. For simplicity, we furthermore assume that car B cannot change its acceleration, it is fixed at 0, the initial acceleration of A does not matter. For this traffic situation, we have a specification expressing that the distance between the two cars needs to be equal to 15 within 5 time units, formalised as the formula  $\varphi_1$  from Ex. 1. Furthermore, assume that we have  $acc_{min} = -10$ ,  $acc_{max} = 5$  and  $spd_{max} = 13$  as bounds on the dynamic behaviour. In this example, there is no timed word of actions that allows the traffic snapshot to behave as specified, if we only allow acceleration changes at one point in time, as we either obtain a speed too fast or need to accelerate stronger than the specified bounds allow. If we allow acceleration changes at two points, there is a solution:  $\omega = \langle (acc(A, 0.75), 0), (acc(A, -6), 4) \rangle$ . Letting one further time unit pass results in a distance exactly 15.*

As described in Sect. 2.1, the dynamics of each car  $C$  evolves according to the simple mechanical equation  $pos'(C) = pos(C) + spd(C) \cdot t + \frac{1}{2} \cdot acc(C) \cdot t^2$ , with  $t$  being the time that elapses and  $pos'(C)$  the new position of  $C$ . The speed evolves according to  $spd'(C) = spd(C) + t \cdot acc(C)$ .

For a finite timed sequence of states  $m = \langle (s_0, I_0), \dots, (s_m, I_m) \rangle$  and a number  $n$  of points in which we can split the interval  $I = [0, t] = \langle I_0, \dots, I_m \rangle$ , we define  $DYN(m, n, I)$  as the set of equations that describe the solution space of  $m$  on the (timing) interval  $I$  for these  $n$  splitting points. In the equations listed below, we only consider the length measurements/constraints that we need to satisfy in  $s_i$ , as the question which discrete actions one needs to execute between two phases is easy to answer. We summarise these constraints as  $\theta_m(t)$  for the length constraints that occur in the phase  $(s_i, I_i)$  of  $m$  with  $t \in I_i$ . If the difference in the position of any two cars affected by it satisfies  $\theta_m(t)$  at point  $t$ , we denote this as  $\Delta pos(t) \models \theta_m(t)$ . Please note that  $\theta_m(t)$  can consist of an arbitrary number of constraint, e.g. when we require that the distance between two cars is smaller than some value and greater than some other value, for example when we want to exclude (potential) collisions while being quite close to the car in front. Additionally,  $\theta_m(t)$  can constrain the distance between more than two cars.

$$DYN(m, n, I) = \text{pos}_0(C) \text{ and } spd_0(C) \text{ are as in } TS_0, \quad (1)$$

$$\text{pos}_n(C) = \text{pos}_{n-1}(C) + spd_{n-1}(C) \cdot t_{n-1} \cdot \frac{1}{2} \cdot acc_{n-1}(C) \cdot t_{n-1}^2, \quad (2)$$

$$spd_n(C) = spd_{n-1}(C) + acc_{n-1}(C) \cdot t_{n-1}, \quad (3)$$

$$\forall t' \in I: \Delta \text{pos}(t) \models \theta_m(t) \text{ and} \quad (4)$$

$$\forall t' \in I, \forall C \in cs: spd_{t'}(C) \text{ and } acc_{t'}(C) \text{ remain in the specified bounds.} \quad (5)$$

An illustration of the solution space and a solution we are searching for for two cars is depicted in Fig. 2. As one can see, we want to find out how many splitting points there need to be such that the difference in the position of the cars satisfies the spatial constraints  $\theta_i$  of each phase  $(s_i, I_i)$  of  $m$  as well as the constraints on the speeds of the cars. The possible curves for the relative position that the evolution yields are furthermore constrained by the maximal and minimal acceleration forces possible. Initial values for the relative position and the speeds are fixed, as they are determined by the traffic snapshot from which on we ask for a satisfying sequence of actions. In the figure, we have both a maximum and a minimum spatial constraint on the distance between the cars in each of the phases. Please note that also a single constraint (distance is e.g. greater than some value) or even no constraint (the distance between the cars is not important in this phase) is possible. Despite being possible, the later one should usually not occur because requiring collision freedom should always be included in a specification, which immediately imposes length constraints.

If  $DYN(m, n, I)$  is satisfied,  $n$  splitting points are sufficient to obtain a satisfying sequence of actions that satisfies the behaviour specified by  $m$ . We furthermore need a relaxed version  $DYN'(m, n, I)$ , which is equivalent to  $DYN$ , except that we alter equation (4) and remove the length constraint on the last phase of  $m$  not reached. Please note that we can rewrite the equations (4) and (5) in an equivalent form that does not use quantifiers, so we gain an easy to solve equation system not dealing with quantifiers over infinite domains.

Later, we are interested how  $DYN'$  behaves when answering the question whether or not there is a solution to  $DYN$ . For this purpose, we denote with  $max\_extension(DYN')$  the maximal value  $x$  s.t.  $DYN'(m, n, [0, x])$  has a solution. Similarly, we denote with  $max\_outcome\_pos(DYN')$  the largest interval  $[y, y']$  s.t.  $DYN'(m', n, I)$  has a solution, where  $m'$  is identical to  $m$  except that in the last phase, the length constraints are replaced with  $[y, y']$ .  $max\_outcome\_spd(DYN')$  is the largest interval  $[z, z']$  for the speed that a car can have when exceeding  $I$  while  $DYN'(m, n, I)$  still has a solution.

Despite not focusing on that topic, we would like to mention that both  $max\_outcome\_pos(DYN')$  and  $max\_outcome\_spd(DYN')$  are vectors, the first one over the differences in positions that are compared in the phases and the second one over the cars.

Utilising the aforementioned equations, we now can give an answer to the question whether or not there is a satisfying sequence of actions s.t. the cars behave as specified:

**Theorem 1** (Number of Accelerations). *Given a traffic snapshot  $TS$  and timed sequence of states  $m = \langle (s_0, I_0), \dots, (s_n, I_n) \rangle$  with adjacent intervals  $I_i$  and sets  $s_i$  of lengths constraints  $\theta$  between the cars to achieve, one can decide after finitely many steps whether or not there is a sequence of acceleration changes  $\omega = \langle (\alpha_0, t_0), \dots, (\alpha_{n'}, t_{n'}) \rangle$  s.t.  $m(TS_0, \omega) = m$ .*

We can decide this question using Alg. 1.

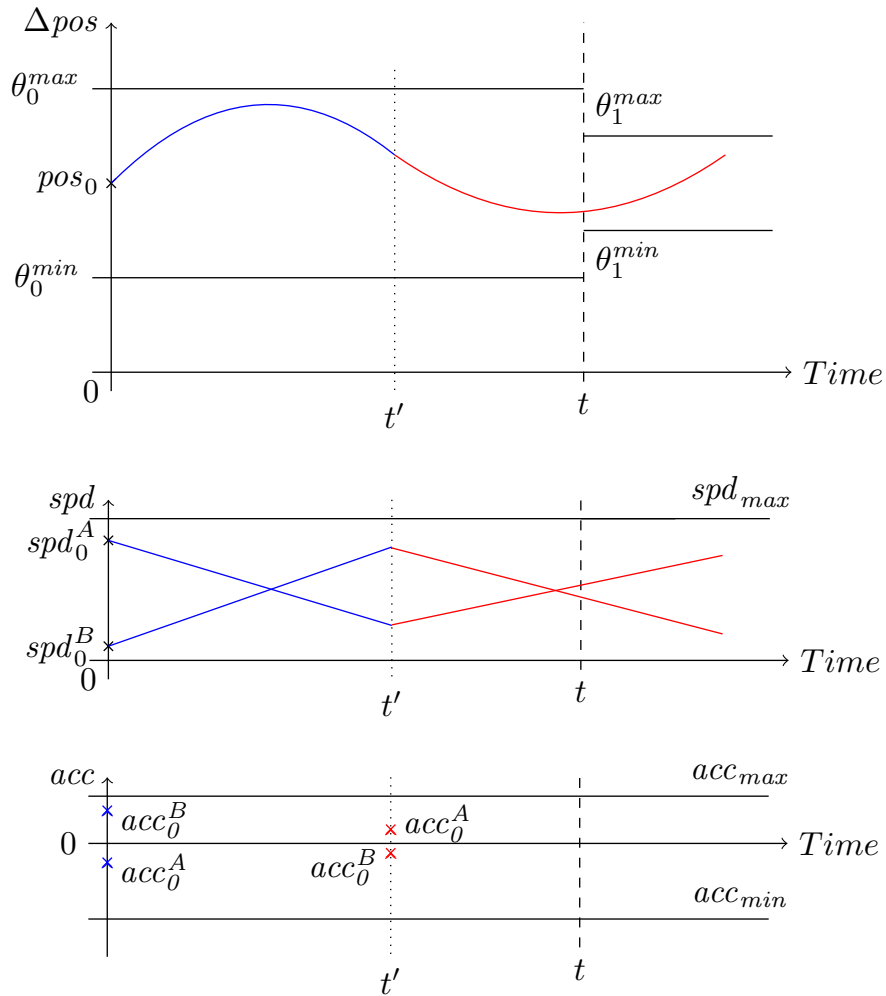


Figure 2: Dynamic Evolution of two cars on a timing interval with one additional acceleration change for each car in between. For a sequence of actions that satisfies the specification, all values need to stay inside the specified intervals. For the relaxed version  $DYN'$ , this does not apply for the last length constraint ( $\theta_1$  here).

**Algorithm 1** Deciding Acceleration

---

**Require:** Input: sequence of states  $m = \langle p_1, \dots, p_m \rangle$ , Interval  $I = [0, t]$

- 1:  $n \leftarrow 0$ ;  $i \leftarrow 1$ ;
- 2: **while**  $i \leq m$  **do**
- 3:      $m' \leftarrow \langle p_1, \dots, p_i \rangle$ ;
- 4:     **while**  $\max\_extension(DYN'(m', I, n)) \neq \max\_extension(DYN'(m', I, n+1))$
- 5:      $\vee \max\_outcome\_pos(DYN'(m', I, n)) \neq \max\_outcome\_pos(DYN'(m', I, n+1))$
- 6:      $\vee \max\_outcome\_spd(DYN'(m', I, n)) \neq \max\_outcome\_spd(DYN'(m', I, n+1))$  **do**
- 7:          $n \leftarrow n + 1$ ;
- 8:     **end while**
- 9:      $i \leftarrow i + 1$ ;
- 10: **end while**
- 11: **if**  $DYN(m, I, n)$  has a solution **then**
- 12:     **return**  $\langle (acc_0(C), t_0), \dots, (acc_n(C), t_n) \rangle$  of  $DYN(m, I, n)$  (for each car  $C$ )
- 13: **else**
- 14:     **return** no Solution existent.
- 15: **end if**

---

**Lemma 1** (Termination and Correctness). *Algorithm 1 terminates iff a solution is found and returns it or there is no solution at all.*

*Proof.* (sketched) Alg. 1 subsequently maximises the outcome that the dynamics (position and speed) may have after each phase of  $m$ , iterating through prefixes  $m'$  of  $m$ , solving  $DYN'$  for this prefix. The first line of the second while-condition (line 4) ensures that we have sufficiently many splitting points such that we reach the right (time) border of the phase. Line 5 ensures that we maximise the difference in the position between two cars at the end of the current phase  $p_i$ , where line 6 maximises the differences in their speeds. If none of these values increases any further within one iteration, no further iteration will. Aborting then is possible because the outcome of the dynamics is strictly monotone in the number of splitting points and thus has converged against a solution that is maximal for the phase.  $\square$

This result is needed in the next section:

## 4 Enforcement

We now present our – in comparison to [4]– enhanced runtime enforcement approach, utilising the results from the previous section.

In this section, we show how cars can find actions to execute in a distributed manner s.t. the overall evolution of the traffic snapshot satisfies the specified properties up to some time bound. In distinction to previous work [4], where the cars non-deterministically guessed actions to execute, they now only propose timed sequences of states that are valid at least for their own properties. The proposed sequences  $\Pi_C$  are afterwards – either by one of the cars or by another central entity – combined into a set of timed sequences of states  $\Pi$ , where each sequence represents a combined behaviour of all cars. It is then checked for the existence of a satisfying sequence of actions  $\omega$ . If existent, the participating cars get informed over the timed actions  $\omega_c$  they themselves must execute to comply to  $\omega$ .

As the sequences  $\Pi$  are in the end checked by a single entity, one could argue that it would be easier to refrain from having the specification distributed over all cars. There are, however, several benefits that one gains when using the more distributed approach. First of all, the whole specification does not need to be known beforehand, neither to the other cars nor to the entity that checks if there is a satisfying sequence of actions. Being able to handle such cases is one big strength of runtime enforcement approaches. Additionally, only the specification up to some time bound needs to be known, not the behaviour beyond that timed horizon, which might not be of interest to the others. Therefore, each car knows its whole specification completely and the central entity/other cars just enough to fulfil its/their task(s). While this argument mostly aimed at privacy concerns, we can also consider it in the light of complexity: The size of the region automaton is exponential in the size of the corresponding SC automaton, which itself is exponential in the size of the specification in SCL. When we considers that only the behaviour up to some time bound is of interest to us, a central SC automaton or even worse, region automaton, would be unnecessary large.

We now focus on the question wherefrom the cars know which timed sequences of states to announce. This includes getting sequences of regions first (Lemma. 2) and computing a satisfying timed sequence of states from them (Lemma. 3).

Before doing so, we would like to mention some results from [4]: Given an region automaton  $\mathcal{R}(A_\psi)$  for a specification  $\psi$  in SCL, one can label some of the states as *bad*, these are the once that, if reached, do not allow the run of the region automaton to be extended in a way that allow  $\psi$  to get satisfied. Vice versa, if a sequence ends in a state that is not bad, one can extend it in a way s.t.  $\psi$  is satisfied.

The second case, however, does not hold if we consider specifications  $\varphi$  in TMLSL rather than SCL: Here, it can be the case that the region automaton claims that there is an extension s.t.  $\varphi$  is satisfied, but the cars are not able to behave in a way that conforms to this extension. Thus, the sequences of regions that the region automaton suggests as satisfiable might actually not be satisfiable (but are candidates):

**Lemma 2** (Sequences of regions as potential solutions). *In every traffic snapshot  $TS$ , one can compute the set  $\Pi$  of sequences of regions  $\pi = \langle [v_i], \dots, [v_j] \rangle$  that start with the region  $[v_i]$  reached in the region automaton  $\mathcal{R}(A_\varphi)$  in the evolution towards  $TS$  and are candidates for satisfying runs of the region automaton. Moreover, there is no sequence  $\pi'$  not in  $\Pi$  but with  $m(\pi') \models_t \varphi$ .*

*Proof.* Using Def. 3 of [4], we can compute the set of locations  $\{l_0, \dots, l_n\}$  that the SC automaton  $A_\varphi$  reaches along the evolution towards  $TS$ . Each location  $l \in \{l_0, \dots, l_n\}$  corresponds to a set of regions  $\{[v_1], \dots, [v_m]\}$ , with  $[v_i] \models \Delta_x(l)$  that is, the region  $[v_i]$  satisfies the clock constraints over the history clocks of the location  $l$  and especially  $p \in \mathcal{L}(l_i)$  iff  $[v_i](p) = 0$  for every proposition  $p$ . We ignore constraints over the prophecy clocks here, because the future is (at least at the end of the sequence) unknown and the history clocks are sufficient for determining the intervals.  $\square$

Given a sequence of regions, we can compute a timed sequence of states that satisfies the sequence of regions:

**Lemma 3** (From regions to timed sequences of states). *For every sequence of regions  $\pi = \langle [v_0], \dots, [v_n] \rangle$  one can construct a timed sequence of state  $m(\pi) = \langle (s_0, I_0), \dots, (s_m, I_m) \rangle$  s.t.  $m \models \pi$ .*

*Proof.* For simplicity, we assume that there is a global clock that is not reset, counting the time from the beginning of the sequence. Starting with  $[v_0]$  and subsequently going trough all  $[v_i]$ , we determine for each point in time  $t$  which propositions  $p$  are valid in it, which is achieved by looking at formulae of the form  $x_p = 0$ . To determine the shape of the intervals ( $[]$ ,  $()$ ,  $[\ ]$  or  $(\ ]$ ), we consider the (in-)equalities in the regions: If some  $p$  is valid in the next point in time,  $y_p = 1$  leads to a closed interval border (“ $]$ ”),  $y_p < 1$  leads to an open one (“ $)$ ”). We do the same for the history clocks  $x_p$ .  $\square$

As all cars announce timed sequences of states, we need to combine them into a single sequence that the central entity can check:

**Lemma 4** (Combining Sequences of states). *Given two finite timed sequences of states  $m_1$  and  $m_2$ , one can construct a timed sequence  $m$  s.t. for every  $\varphi$ :  $m_1 \models \varphi \vee m_2 \models \varphi \implies m \models \varphi$ .*

*Proof.* We start with an “empty” sequence  $m = \langle (\_, [0, 0]), (\_(0, 1)), (\_, [1, 1]), (\_, (1, 2)), \dots (\_, I_n) \rangle$ . Going through each state  $(\_, I_i)$  of  $m$ , we look in both  $m_1$  and  $m_2$  and insert the propositions from the states  $(s_j, I_j)$ , where  $I_j$  contains  $I_i$ . If it happens that for some two neighbouring states  $(s_i, I_i)$  and  $(s_{i+1}, I_{i+1})$ ,  $s_i = s_{i+1}$  holds, we can fuse the two into a single state  $(s_i, I_i + I_{i+1})$ . As a last step, we check if the resulting sequence is consistent. If it happens that there is a contradiction in one of the states, say  $cl(A)$  and  $\neg cl(A)$  need to hold in the same (time) interval, the timed sequence of states cannot be satisfied at all and is thus invalid.  $\square$

Please note that the other direction not necessarily holds, as e.g.  $p_1 \wedge p_2$  could hold in  $[1, 1]$  of  $m$ , but in  $m_1$  only  $p_1$  and in  $m_2$  only  $p_2$  holds in the respective interval.

We utilise the aforementioned results in the controllers of the cars and the central entity that determines whether a solution exists. The controller is depicted in Fig. 3 and the central entity *RSU* (Road-Side Unit) in Fig. 4.

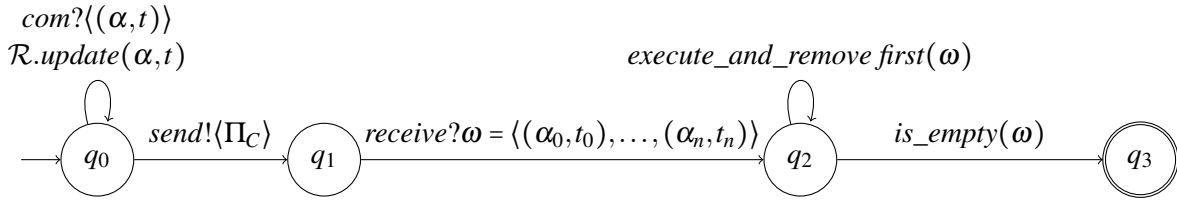


Figure 3: Controller for each car  $C$ , each is equipped with an instance of it. The controller keeps track of the traffic situation in  $q_0$  and updates the sequences it could announce accordingly. When announcing the sequences  $\Pi_C$ , it proceeds to  $q_1$  and waits for a positive response from the central entity and executes the sequence of actions  $\omega$  that it received from there, until there are not further actions to execute. Please note that we omitted clock constraints to actually force the controller to leave a state.

Both of them use several functions on their transitions. In the controller,  $\mathcal{R}.update(\alpha, t)$  is used so that the internal region automaton keeps track about the behaviour on the road and thus is in the correct state(s), before the enforcement mechanism is triggered. In this location, we may already have an evolution that leads to a state s.t. all further extensions are unsatisfiable. If such a behaviour is undesired, constraints should be added s.t. one does not stay in this location.  $execute\_and\_remove\_first(\omega)$  takes the first time stamped action  $(\alpha_1, t_1)$  from the action sequence  $\omega = \langle (\alpha_1, t_1), \dots, (\alpha_n, t_n) \rangle$ , waits until the clock reaches  $t_1$  and executes  $\alpha_1$ . Afterwards, this element is removed from  $\omega$ , so that the next action is ready to be executed.  $is\_empty(\omega)$  is true for the empty sequence  $\langle \rangle$ .

In *RSU*, the function  $D.push(\Pi_C)$  is used to internally store the announced sets of timed sequences of states  $\Pi_C$  for each car  $C$  in some data structure  $D$ . The set of sequences representing all possible satisfying sequences for all cars is constructed using  $combine(C_1, \dots, C_n)$  and stored in  $\Pi$ . Using Alg. 1, it can then decide whether or not one of the sequences in  $\Pi$  is one for which there is a satisfying sequence of actions. If so, the solution  $\omega$  is computed and afterwards split into single solutions  $\omega_c$  for each car  $c$ , so that every car only gets informed of the actions it itself has to execute. After informing a car, it is removed from the data structure  $D$ .

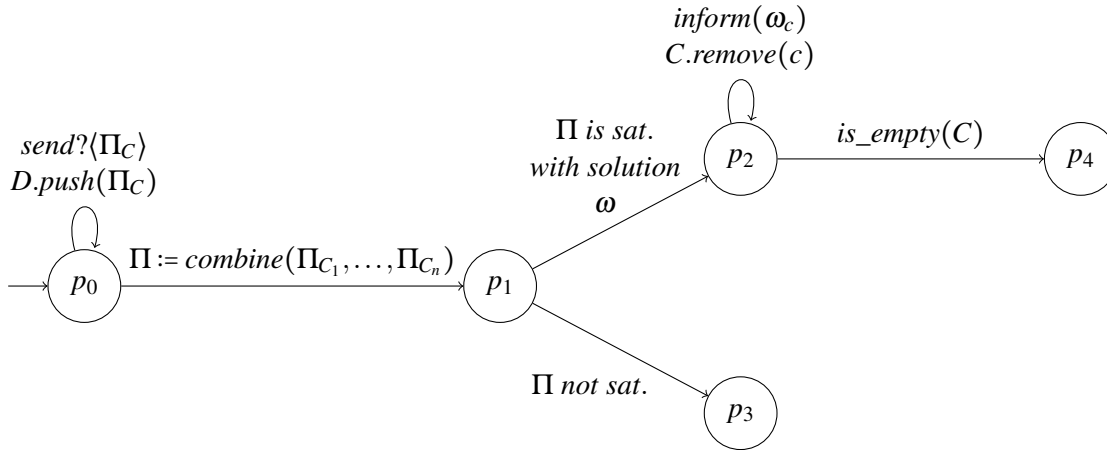


Figure 4: Central Entity *RSU*. *RSU* waits in  $p_0$  for announced sequences  $\Pi_{C_i}$  from the cars and combines them into a single set of sequences  $\Pi$ , which it then can check for a satisfying run. If positively answered, it sends the sequence to execute to each car.

The communication takes place over the channels *com*, over which the cars announce action they execute before the enforcement mechanism is triggered. Channel *send* is used to inform *RSU* about the possible plans  $\Pi_C$  of each car. *receive* is used for the opposite direction, informing the cars which actions to execute.

**Theorem 2** (Correct- and Soundness). *If the controller from Fig. 3 proceeds to location  $q_3$ , the specification of all cars is satisfied up to the given time bound  $t$ . If it cannot proceed to location  $q_3$ , then there is no sequence of actions for the cars to take that respects the specification of all cars.*

*Proof.* Due to Lemma 2, we get all possible satisfying sequences of regions of cars. Due to Lemma 3, we can compute equivalent timed sequences of states from them. Lemma 4 lets us combine them towards some  $\Pi$  on the transition from  $p_0$  to  $p_1$  in *RSU* s.t. all possible combinations of solutions for all cars are considered. Each of them is checked using Alg. 1, so the due to Lemma 1, the solution found is a correct one.  $\square$

If *RSU* reaches location  $p_3$ , the specification is unsatisfiable, so there is no sequence of actions for the cars to execute. In this case collision freedom can still be guaranteed (assuming that there were no collisions yet), as the reservations of the cars occupy a space big enough to come to a standstill within that space.

**Remark 3** (Number of sequences to consider). *If all cars announce all timed sequences of states that satisfy their specification, the central entity/road-side unit needs to check all combinations of these sequences (with exactly one sequence in the combination for each car), resulting in a lot of computation. However, these computations do not depend on each other and can thus be parallelised. If we consider that the cars themselves compute this, rather than some road cite unit, one can think of a more advanced protocol than the one proposed here, where the cars distribute the sequences to check and thus the computational effort between each other.*

**Remark 4** (Discrete Actions). *Through both Sect. 3 and Sect. 4, we only considered how the cars can change their accelerations to ensure that they satisfy the length measurements in the specification, ignoring the discrete actions completely. As said, given a timed sequence of states, it is easy to see what*

*discrete actions are to execute when, as they directly change the formulae valid and thus the phase. Some of them, however, need to be respected when constructing the length comparisons  $\theta_i$  that we check in DYN.*

## 5 Conclusion

**Contribution** In this paper, we proposed a runtime enforcement approach for autonomous car in motorway traffic, employing communication between the cars, where the knowledge about the satisfaction of a property is represented using a region automaton. In answering the question whether or not a specification (now in the form of a timed sequence of states) is satisfiable, we were able to eliminate one of the roots for the semi-decidability of the satisfiability problem of TMLSL.

**Future Work** Future work on the topic includes studying the satisfiability problem of TMLSL again, in an effort to show that the logic is indeed decidable over infinite runs. Further topics also include the extension of the logic and the proposed runtime enforcement approach towards the aforementioned more complex road topologies. Both of them offer some challenges in the semantics and runtime enforcement, as their models are more complicated than the ones for motorway traffic. For urban traffic, the assumption that there is a central entity that all cars can communicate with is not too far from reality, as on almost all intersections traffic lights are present, some of which already communicate with the buses that cross/approach them.

Steps towards an implementation for solving the decidability problem of *TMLSL* were made and could be adjusted to be used in the runtime enforcement setting. With an implementation, we could also examine if the proposed approach is suited for real-time applications like car control on motorways, e.g. the computation happens fast enough.

**Acknowledgements.** We thank the anonymous reviewers for their helpful comments.

## References

- [1] Matthias Althoff, Sebastian Maierhofer & Christian Pék (2021): *Provably-Correct and Comfortable Adaptive Cruise Control*. *IEEE Trans. Intell. Veh.* 6(1), pp. 159–174, doi:10.1109/TIV.2020.2991953.
- [2] Rajeev Alur & David L. Dill (1994): *A Theory of Timed Automata*. *Theor. Comput. Sci.* 126(2), pp. 183–235, doi:10.1016/0304-3975(94)90010-8.
- [3] Christopher Bishopink & Ernst-Rüdiger Olderog (2022): *Spatial and Timing Properties in Highway Traffic*. In Helmut Seidl, Zhiming Liu & Corina S. Pasareanu, editors: *Theoretical Aspects of Computing - ICTAC 2022 - 19th International Colloquium, Tbilisi, Georgia, September 27-29, 2022, Proceedings, Lecture Notes in Computer Science* 13572, Springer, pp. 114–131, doi:10.1007/978-3-031-17715-6\_9.
- [4] Christopher Bishopink & Ernst-Rüdiger Olderog (2023): *Time for Traffic Manoeuvres*, pp. 163–179. Springer Nature Switzerland, Cham, doi:10.1007/978-3-031-40132-9\_11.
- [5] Christopher Bishopink & Maïke Schwammberger (2019): *Verification of Fair Controllers for Urban Traffic Manoeuvres at Intersections*. In Emil Sekerinski, Nelma Moreira, José N. Oliveira, Daniel Ratiu, Riccardo Guidotti, Marie Farrell, Matt Luckcuck, Diego Marmsoler, José Creissac Campos, Troy Astarte, Laure Gonnord, Antonio Cerone, Luis Couto, Brijesh Dongol, Martin Kutrib, Pedro Monteiro & David Delmas, editors: *Formal Methods. FM 2019 International Workshops - Porto, Portugal, October 7-11, 2019, Revised Selected Papers, Part I, Lecture Notes in Computer Science* 12232, Springer, pp. 249–264, doi:10.1007/978-3-030-54994-7\_18.



- [6] Martin Fränzle, Michael R. Hansen & Heinrich Ody (2015): *No Need Knowing Numerous Neighbours - Towards a Realizable Interpretation of MLSL*. In Roland Meyer, André Platzer & Heike Wehrheim, editors: *Correct System Design, Lecture Notes in Computer Science* 9360, Springer, pp. 152–171, doi:10.1007/978-3-319-23506-6\_11.
- [7] Klaus Havelund & Allen Goldberg (2008): *Verify Your Runs*, pp. 374–383. Springer Berlin Heidelberg, Berlin, Heidelberg, doi:10.1007/978-3-540-69149-5\_40.
- [8] Martin Hilscher, Sven Linker & Ernst-Rüdiger Olderog (2013): *Proving Safety of Traffic Manoeuvres on Country Roads*. In Zhiming Liu, Jim Woodcock & Huibiao Zhu, editors: *Theories of Programming and Formal Methods - Essays Dedicated to Jifeng He on the Occasion of His 70th Birthday, Lecture Notes in Computer Science* 8051, Springer, pp. 196–212, doi:10.1007/978-3-642-39698-4\_12.
- [9] Martin Hilscher, Sven Linker, Ernst-Rüdiger Olderog & Anders P. Ravn (2011): *An Abstract Model for Proving Safety of Multi-lane Traffic Manoeuvres*. In Shengchao Qin & Zongyan Qiu, editors: *Formal Methods and Software Engineering - 13th International Conference on Formal Engineering Methods, ICFEM 2011, Durham, UK, October 26-28, 2011. Proceedings, Lecture Notes in Computer Science* 6991, Springer, pp. 404–419, doi:10.1007/978-3-642-24559-6\_28.
- [10] Kim Guldstrand Larsen, Marius Mikucionis & Jakob Haahr Taankvist (2015): *Safe and Optimal Adaptive Cruise Control*. In Roland Meyer, André Platzer & Heike Wehrheim, editors: *Correct System Design, September 8-9, 2015. Proceedings, Lecture Notes in Computer Science* 9360, Springer, pp. 260–277, doi:10.1007/978-3-319-23506-6\_17.
- [11] Kim Guldstrand Larsen, Paul Pettersson & Wang Yi (1997): *UPPAAL in a Nutshell*. *Int. J. Softw. Tools Technol. Transf.* 1(1-2), pp. 134–152, doi:10.1007/s100090050010.
- [12] Sarah M. Loos, André Platzer & Ligia Nistor (2011): *Adaptive Cruise Control: Hybrid, Distributed, and Now Formally Verified*. In Michael J. Butler & Wolfram Schulte, editors: *FM 2011: Formal Methods - 17th International Symposium on Formal Methods, Limerick, Ireland, 4, 2011. Proceedings, Lecture Notes in Computer Science* 6664, Springer, pp. 42–56, doi:10.1007/978-3-642-21437-0\_6.
- [13] Heinrich Ody (2015): *Undecidability Results for Multi-Lane Spatial Logic*. In Martin Leucker, Camilo Rueda & Frank D. Valencia, editors: *Theoretical Aspects of Computing - ICTAC, Lecture Notes in Computer Science* 9399, Springer, pp. 404–421, doi:10.1007/978-3-319-25150-9\_24.
- [14] Heinrich Ody (2020): *Monitoring of traffic manoeuvres with imprecise information*. Ph.D. thesis, University of Oldenburg, Germany. Available at <https://oops.uni-oldenburg.de/4730>.
- [15] Jean-François Raskin & Pierre-Yves Schobbens (1997): *State Clock Logic: A Decidable Real-Time Logic*. In Oded Maler, editor: *Hybrid and Real-Time Systems, International Workshop. HART'97, Grenoble, France, March 26-28, 1997, Proceedings, Lecture Notes in Computer Science* 1201, Springer, pp. 33–47, doi:10.1007/BFb0014711.
- [16] Fred B. Schneider (2000): *Enforceable security policies*. *ACM Trans. Inf. Syst. Secur.* 3(1), pp. 30–50, doi:10.1145/353323.353382.
- [17] Maike Schwammberger (2018): *An abstract model for proving safety of autonomous urban traffic*. *Theor. Comput. Sci.* 744, pp. 143–169, doi:10.1016/j.tcs.2018.05.028.
- [18] Maike Schwammberger (2018): *Introducing Liveness into Multi-lane Spatial Logic lane change controllers using UPPAAL*. In Mario Gleirscher, Stefan Kugele & Sven Linker, editors: *Proceedings 2nd International Workshop on Safe Control of Autonomous Vehicles, SCAV@CPSWeek 2018, Porto, Portugal, 10th April 2018, EPTCS* 269, pp. 17–31, doi:10.4204/EPTCS.269.3.
- [19] J. C. P. Woodcock & Jim Davies (1996): *Using Z - specification, refinement, and proof*. Prentice Hall international series in computer science, Prentice Hall.

# Correct-by-Construction Control for Stochastic and Uncertain Dynamical Models via Formal Abstractions

Thom Badings   Nils Jansen

Radboud University  
Nijmegen, the Netherlands  
thom.badings@ru.nl

Licio Romao   Alessandro Abate

University of Oxford  
Oxford, United Kingdom

Automated synthesis of correct-by-construction controllers for autonomous systems is crucial for their deployment in safety-critical scenarios. Such autonomous systems are naturally modeled as stochastic dynamical models. The general problem is to compute a controller that provably satisfies a given task, represented as a probabilistic temporal logic specification. However, factors such as stochastic uncertainty, imprecisely known parameters, and hybrid features make this problem challenging. We have developed an abstraction framework that can be used to solve this problem under various modeling assumptions. Our approach is based on a robust finite-state abstraction of the stochastic dynamical model in the form of a Markov decision process with intervals of probabilities (iMDP). We use state-of-the-art verification techniques to compute an optimal policy on the iMDP with guarantees for satisfying the given specification. We then show that, by construction, we can refine this policy into a feedback controller for which these guarantees carry over to the dynamical model. In this short paper, we survey our recent research in this area and highlight two challenges (related to scalability and dealing with nonlinear dynamics) that we aim to address with our ongoing research.

## 1 Introduction

Controlled autonomous systems are increasingly deployed in safety-critical settings [30]. When the transitions between states are specified by probabilities, autonomous systems can often be naturally modeled as stochastic dynamical models [26]. For deployment in safety-critical settings, controllers for stochastic models must act safely and reliably with respect to desired specifications. Traditional control design methods use, e.g., Lyapunov functions and optimization to provide guarantees for simple tasks such as stability, convergence, and invariance [11]. However, alternative methods are needed to give formal guarantees about richer temporal specifications relevant to, for example, safety-critical applications [20].

**Formal controller synthesis** Temporal logic is a rich language for specifying the desired behavior of autonomous systems [32]. In particular, probabilistic computation tree logic (PCTL, [25]) is widely used to define temporal requirements on the behavior of probabilistic systems. For example, in a motion control problem for an unmanned aerial vehicle (UAV), a PCTL formula can specify that, with at least 90% probability, the UAV must safely fly to a target location within 2 minutes without crashing into obstacles (commonly known as a *reach-avoid specification* [21]). Leveraging tools from probabilistic verification [9], the problem is to synthesize a controller that ensures the satisfaction of such a PCTL formula for the model under study [24]. Finite abstractions can make continuous models amenable to techniques and tools from formal verification: by discretizing their state and action spaces, abstractions result in, e.g., finite Markov decision processes (MDPs) that soundly capture the continuous dynamics [2]. Verification guarantees on the finite abstraction can thus carry over to the continuous model.

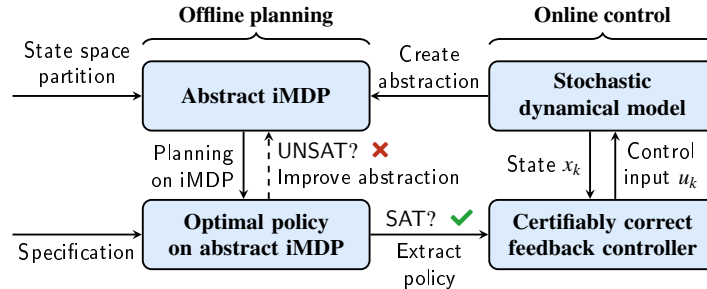


Figure 1: Our overall approach integrated into a safe model-based learning framework.

**Problem statement** In this research, we adopt such an abstraction-based approach to controller synthesis for autonomous systems. Our goal is to compute feedback controllers that *provably satisfy* a given temporal logic specification. In this short paper, we focus on reach-avoid specifications, but most of our approaches can readily be extended to general PCTL properties [35]. We consider the following general problem:

Given (1) a discrete-time stochastic dynamical model and (2) a reach-avoid specification, compute a *feedback controller* together with a *certificate* in the form of a probability threshold, such that the induced closed-loop system satisfies this specification with at least this certified probability.

In this paper, we survey our recent work in which we have considered this general problem under various modeling assumptions. First, we provide a general introduction to our abstraction framework. Thereafter, we summarize our main results from several recent papers [7, 5, 8, 35]. Finally, we highlight two key challenges that remain open, and we describe our current research plans that aim to address these changes.

**Our abstraction framework** Our general abstraction framework is shown in Fig. 1. First, we compute a finite-state abstraction of the stochastic dynamical model [38], which we obtain from a *partition* of its continuous state space into a set of disjoint convex *regions*. Actions in this abstraction correspond to continuous control inputs that yield transitions between these regions. Due to the stochastic noise in the dynamical model, the outcome of an action is stochastic, rendering transitions probabilistic. We capture these probabilities in a Markov decision process (MDP) [34]. A defining characteristic of our approach is that we leverage *backward reachability computations* on the dynamical model to determine which actions are enabled at each discrete region. By contrast, most other abstraction methods rely on *forward reachability computations*, which are associated with errors that grow with the time horizon of the property (see related work for details). Our backward scheme avoids such abstraction errors, at the cost of requiring slightly more restrictive assumptions on the model dynamics (see, e.g., [8] for details).

**Interval MDPs** Computing the transition probabilities of the abstraction is subject to estimation errors. To be *robust* against estimation errors in these probabilities, we use data-driven techniques [13, 36] to compute *upper and lower bounds* on the transition probabilities with a predefined *confidence level*. We formalize our abstractions with the probably approximately correct (PAC) probability intervals using so-called iMDPs, which are an extension of MDPs with intervals of probabilities [22]. Policies for iMDPs have to robustly account for *all possible probabilities* within the intervals [33, 39]. In our implementation, we compute robust policies using robust value iteration within the probabilistic model checker PRISM [27]. We show that any policy on the iMDP can be refined into a *piecewise linear feedback controller* for the dynamical model. Crucially, the probability of satisfying the reach-avoid property on the iMDP is a *lower bound* on the satisfaction probability for the dynamical model, thus solving the problem above.

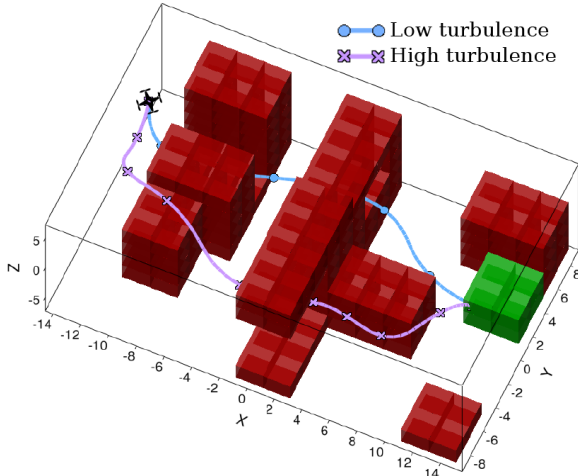


Figure 2: UAV reach-avoid problem (goal in green; obstacles in red), plus simulations with the optimal iMDP-based controller from initial state  $x_0 = [-14, 0, 6, 0, -6, 0]^T$ , under high/low turbulence.

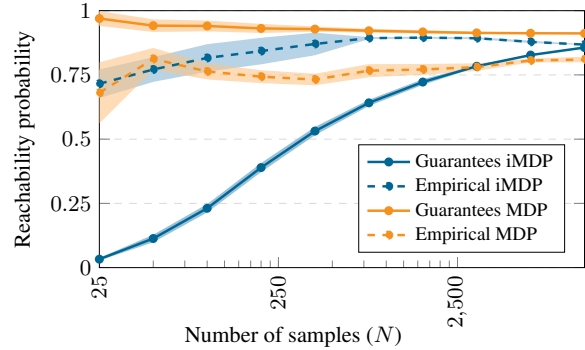


Figure 3: Reach-avoid guarantees on the iMDPs (blue) and MDPs (orange) for their respective policies, versus the resulting empirical (simulated) performance (dashed lines) on the dynamical system. The empirical performance obtained from the MDPs violates the guarantees, whereas that from the iMDPs does not.

**Related work** Abstractions of stochastic models are well-studied [2, 3], with applications to stochastic hybrid [16, 29], switched [28], and partially observable systems [6, 23]. Various tools exist, e.g., StocHy [15] and ProbReach [37]. A distinguishing feature of our abstraction scheme is that we use *backward reachability computations* on the model dynamics to determine the subset of actions enabled in each abstract state. By contrast, standard abstraction methods typically rely on *forward reachability computations* based on discretizing the control input space. In particular, such forward methods propagate *sets of states*  $\mathcal{X} \subset \mathbb{R}^n$  forward through the model dynamics under discretized input  $\hat{\mathbf{u}}_k$  (see the notation from Eq. (1)). Since the noise  $\eta_k$  is stochastic, this yields a *set of distributions* over successor states, which can be difficult to reason over. By contrast, with our backward computations, each abstract action yields a *single distribution*, which is independent of where this action was chosen. However, this requires a higher degree of system controllability, as discussed in more detail in [8, Assumption 2].

## 2 Correct-by-Construction Control via Formal Abstractions

In general, we consider discrete-time, continuous-state dynamical models, where the progression of the state  $x \in \mathbb{R}^n$  depends *linearly* on the current state, on a control input, and on a process noise term. Given a state  $x_k$  at discrete time  $k \in \mathbb{N}$ , the successor state  $x_{k+1}$  at time  $k+1$  is computed as

$$x_{k+1} = Ax_k + Bu_k + q_k + \eta_k, \quad (1)$$

with matrices  $A$  and  $B$ , and a continuous control input  $u_k \in \mathcal{U} \subseteq \mathbb{R}^p$  (i.e., action). The term  $\eta_k \in \Delta \subset \mathbb{R}^n$  is an arbitrary additive process noise term, which is an i.i.d. random variable defined on a probability space  $(\Delta, \mathcal{D}, \mathbb{P})$ , with  $\sigma$ -algebra  $\mathcal{D}$  and probability measure  $\mathbb{P}$  defined over  $\mathcal{D}$ . A controller (i.e., control policy)  $c: \mathbb{R}^n \times \mathbb{N} \rightarrow \mathcal{U}$  chooses a control input based on the current state  $x \in \mathbb{R}^n$  and time  $k \in \mathbb{N}$ .

We now highlight some of the variants of the problem stated in Sect. 1 we have considered thus far.

## 2.1 Stochastic noise of unknown distribution

It is commonly assumed that the distribution of the process noise  $\eta_k$  is known and/or Gaussian [31]. However, in many realistic problems, this assumption yields a poor approximation of the uncertainty [12]. Distributions may even be *unknown*, meaning that one cannot derive a set-bounded or a precise probabilistic representation of the noise. In this case, it is generally hard or even impossible to derive *hard guarantees* on the probability that a given controller ensures the satisfaction of a reach-avoid property.

In papers [5, 8], we thus consider a variant of the controller synthesis problem from Sect. 1 for dynamical systems with additive process noise of an *unknown distribution*. Specifically, the probability measure  $\mathbb{P}$  of the noise  $\eta_k \in \Delta \subset \mathbb{R}^n$  is unknown but time-invariant. To deal with this lack of knowledge, we adapt tools from the scenario approach [14, 13] to compute *PAC* interval estimates for the transition probabilities of the abstract model based on a finite set of samples of the noise. We capture these bounds in the transition probability intervals of a so-called interval Markov decision process (iMDP). This iMDP is, with a user-specified confidence probability, robust against uncertainty in the transition probabilities, and the tightness of the probability intervals can be controlled through the number of samples.

In [8], we use this method to solve a reach-avoid problem for a UAV operating under turbulence (we compare scenarios with different turbulence levels), represented by stochastic noise of unknown distribution. The UAV is modeled by a 6D dynamical model (we refer to [8] for the explicit model). In Fig. 2, we show simulations under the optimal controller for two turbulence levels. Under low noise, the controller prefers the short but narrow path. On the other hand, under high noise, the longer but safer path is preferred. Thus, accounting for process noise is important to obtain controllers that are safe.

We also compared our robust iMDP approach against a naive MDP abstraction. This MDP has the same states and actions as the iMDP, but uses precise (frequentist) probabilities. The maximum reachability probabilities (guarantees) for both methods are shown in Fig. 3. For every value of  $N$ , we apply the resulting controllers to the dynamical system in Monte Carlo simulations with 10,000 iterations to determine the empirical reachability probability. Fig. 3 shows that the non-robust MDPs yield *poor and unsafe performance guarantees*: the actual reachability of the controller is much lower than the reachability guarantees obtained from PRISM. By contrast, our robust iMDP-based approach consistently yields safe lower bound guarantees on the actual performance of controllers.

## 2.2 Set-bounded parameter uncertainty

The approach described in Sect. 2.1 requires *precise knowledge* of the model parameters (namely, the matrices  $A$  and  $B$ ). However, in many realistic cases, there is *epistemic uncertainty* about the precise values of these parameters. For example, consider again the UAV from Sect. 2.1. As shown in Fig. 4, the drone's dynamics depend on uncertain factors, such as the wind and the drone's mass. We assumed that the wind is adequately described by a probabilistic model, reflected in the process noise  $\eta_k$ . Now, let us assume we know that the drone's mass lies between 0.75–1.25 kg, but we do not have information about the likelihood of each value, so employing a probabilistic model is unrealistic. Thus, we treat epistemic uncertainty in such imprecisely known parameters (in this case, the mass) using a *nondeterministic framework* instead.

We have recently extended our abstraction framework in [7] to capture both stochastic noise and set-bounded uncertain parameters. Specifically, we synthesize a controller that (1) is *robust against nondeterminism* due to parameter uncertainty and (2) *reasons over probabilities* derived from stochastic noise. In other words, the controller must satisfy a given specification *under any possible outcome of the nondeterminism* (robustness) and *with at least a certain probability regarding the stochastic noise* (reasoning over probabilities). As before, we wish to synthesize a controller with a *PAC*-style guarantee:

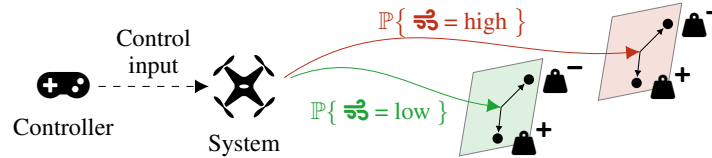


Figure 4: Stochastic uncertainty in the wind ( $w$ ) causes probability distributions over outcomes of controls, while set-bounded uncertainty in the mass ( $m$ ) of the drone causes state transitions to be nondeterministic.

we wish to find a controller that satisfies a reach-avoid specification with at least a desired lower bound threshold probability, and (because our algorithm involves random sampling) that claim should hold with at least a predefined confidence level.

Our experiments in [7] show that we can synthesize controllers that are robust against uncertainty and, in particular, against deviations in the model parameters. Moreover, we show that our method can be used to faithfully capture any uncertainty or error term in the dynamical model that is represented by a bounded set, thus opening the door for the abstraction of nonlinear systems.

### 2.3 Markov jump linear systems

Our approaches described so far are limited to systems with purely continuous dynamics. Thus, these approaches are incompatible with cyber-physical systems, which are characterized by the coupling of digital (discrete) with physical (continuous) components. This results in a *hybrid system* that can jump between discrete modes of operation, each of which is characterized by its own continuous dynamics [29].

To alleviate this restriction, we have extended our abstraction framework in [35] to Markov jump linear systems (MJLSs), which are a well-known class of stochastic, hybrid models suitable for capturing the behavior of cyber-physical systems [19]. An MJLS consists of a finite set of linear dynamics defined by Eq. (1) (also called *operational modes*), where jumps between these modes are governed by a Markov chain (MC). If mode jumping can be controlled, the jumps are governed by an MDP. Due to the jumping between modes, the overall dynamics of an MJLS are nonlinear, making controller synthesis challenging. For brevity, we refer to [35] for further results in this problem setting.

## 3 Current research directions

As discussed above, we have considered the general problem in Sect. 1 under various model assumptions. At the same time, each of those settings suffers from its limitations and necessary assumptions, so the general problem of *optimal control under uncertainty* is far from solved. In this section, we discuss two key limitations of our current framework, which are related to scalability and to linearity of the dynamics. Moreover, we describe how we try to address both of these challenges with our ongoing research.

### 3.1 Neural-guided abstraction of nonlinear systems

Thus far, our research has focused on dynamical models with linear dynamics. Extensions to nonlinear dynamical systems are non-trivial and may require more involved reachability computations [10, 17]. Specifically, the challenge is that the backward reachability computations involved in our approach may become non-convex under nonlinear dynamics.

**Neural network partitioning** A recent paper [1] has proposed to use feedforward neural networks to *learn* state space partitions for nonlinear dynamical models into polyhedral regions. Inspired by this approach, we are developing an abstraction procedure for nonlinear stochastic dynamical models, which (1) learns a polyhedral state space partition using a neural network, and (2) constructs a piecewise linear approximation of the nonlinear dynamics based on this partition. By defining the loss function for the neural network such that it minimizes the linearization error across the partition, we hope to find smarter partitions (into fewer elements and of better geometry) than the rectangular ones we employed thus far.

**Abstraction of linearized dynamics** To account for the error caused by the linearization, we add a set-bounded nondeterministic disturbance to the linearized dynamics. As we have shown in [7], we can robustly capture this set-bounded disturbance in an iMDP abstraction. However, the quality of the abstraction largely depends on the size of the disturbance representing the linearization error. Thus, the main challenge with this approach is to obtain a tight, set-bounded representation of the linearization error.

### 3.2 Abstractions of polyhedral Lyapunov functions

Discrete abstractions are computationally expensive in general due to the discretization of the state space. For example, the number of abstract states scales exponentially with the dimension of the state space, commonly called the *curse of dimensionality*. Moreover, adding robustness to multiple sources of uncertainty (as we have done in [7] further increases the number of transitions modeled in the abstract model). Thus, finding ways to reduce the complexity of abstraction while keeping their expressivity is a challenging direction for further research.

**Abstraction of Lyapunov functions** Inspired by [18] and the large body of literature on Lyapunov and Barrier functions [4], we are developing a method for abstracting stochastic dynamical systems using Lyapunov functions. Specifically, we wish to generate an abstract model whose states represent annuli of the sublevel sets of a Lyapunov function. A similar approach was used by [18]. However, the approach by [18] relies on a *strict decrease condition* on the Lyapunov function and is therefore restricted to *nonstochastic* linear systems only. Instead, we believe that our abstraction procedure based on backward reachability analysis can be used to construct *sound* abstractions of sublevel sets of Lyapunov functions.

**Complexity is independent of state dimension** This envisioned abstraction of Lyapunov sublevel sets avoids the need for an exhaustive partitioning of the state space. Notably, the number of states in the envisioned abstraction is *independent of the dimension of the state space*. Thus, we believe that this approach may significantly reduce the computational complexity of the abstraction.

## 4 Conclusions and Future Work

In this short paper, we have surveyed our recent research on abstraction-based controller synthesis for stochastic and uncertain dynamical models. Based on a robust finite-state abstraction in the form of an iMDP, we are able to synthesize controllers for dynamical models that *provably satisfy* given temporal logic specifications, such as reach-avoid tasks. We have considered this general problem under various modeling assumptions, including unknown noise distributions, imprecisely known model parameters, and hybrid features. Moreover, we have highlighted two key challenges that are related to scalability and extensions to nonlinear systems. With our ongoing research, we aim to address these challenges.

**Acknowledgements** This research has been partially funded by NWO grant NWA.1160.18.238 (PrimaVera), by EPSRC IAA Award EP/X525777/1, and by the ERC Starting Grant 101077178 (DEUCE).

## References

- [1] Alessandro Abate, Alec Edwards & Mirco Giacobbe (2022): *Neural Abstractions*. In: *NeurIPS*.
- [2] Alessandro Abate, Maria Prandini, John Lygeros & Shankar Sastry (2008): *Probabilistic reachability and safety for controlled discrete time stochastic hybrid systems*. *Automatica* 44(11), pp. 2724 – 2734, doi:10.1016/j.automatica.2008.03.027.
- [3] Rajeev Alur, Thomas A. Henzinger, Gerardo Lafferriere & George J. Pappas (2000): *Discrete abstractions of hybrid systems*. *Proc. IEEE* 88(7), pp. 971–984, doi:10.1109/5.871304.
- [4] Aaron D. Ames, Samuel Coogan, Magnus Egerstedt, Gennaro Notomista, Koushil Sreenath & Paulo Tabuada (2019): *Control Barrier Functions: Theory and Applications*. In: *ECC*, IEEE, pp. 3420–3431, doi:10.23919/ECC.2019.8796030.
- [5] Thom S. Badings, Alessandro Abate, Nils Jansen, David Parker, Hasan A. Poonawala & Mariëlle Stoelinga (2022): *Sampling-Based Robust Control of Autonomous Systems with Non-Gaussian Noise*. In: *AAAI*, AAAI Press, pp. 9669–9678, doi:10.1609/aaai.v36i9.21201.
- [6] Thom S. Badings, Nils Jansen, Hasan A. Poonawala & Mariëlle Stoelinga (2023): *Correct-by-construction reach-avoid control of partially observable linear stochastic systems*. *CoRR* abs/2103.02398, doi:10.48550/arXiv.2103.02398.
- [7] Thom S. Badings, Licio Romao, Alessandro Abate & Nils Jansen (2023): *Probabilities Are Not Enough: Formal Controller Synthesis for Stochastic Dynamical Models with Epistemic Uncertainty*. In: *AAAI*, AAAI Press, pp. 14701–14710, doi:10.1609/aaai.v37i12.26718.
- [8] Thom S. Badings, Licio Romao, Alessandro Abate, David Parker, Hasan A. Poonawala, Mariëlle Stoelinga & Nils Jansen (2022): *Robust Control for Dynamical Systems with Non-Gaussian Noise via Formal Abstractions*. *J. Artif. Intell. Res.*, doi:10.1613/jair.1.14253.
- [9] Christel Baier & Joost-Pieter Katoen (2008): *Principles of model checking*. MIT Press.
- [10] Somil Bansal, Mo Chen, Sylvia L. Herbert & Claire J. Tomlin (2017): *Hamilton-Jacobi reachability: A brief overview and recent advances*. In: *CDC*, IEEE, pp. 2242–2253, doi:10.1109/CDC.2017.8263977.
- [11] Calin Belta, Boyan Yordanov & Ebru Aydin Gol (2017): *Formal Methods for Discrete-Time Dynamical Systems*. Springer International Publishing, doi:10.1007/978-3-319-50763-7.
- [12] Lars Blackmore, Masahiro Ono, Askar Bektassov & Brian C. Williams (2010): *A Probabilistic Particle-Control Approximation of Chance-Constrained Stochastic Predictive Control*. *IEEE Trans. Robotics* 26(3), pp. 502–517, doi:10.1109/TRO.2010.2044948.
- [13] Marco C. Campi, Algo Carè & Simone Garatti (2021): *The scenario approach: A tool at the service of data-driven decision making*. *Annu. Rev. Control.* 52, pp. 1–17, doi:10.1016/j.arcontrol.2021.10.004.
- [14] Marco C. Campi & Simone Garatti (2008): *The Exact Feasibility of Randomized Solutions of Uncertain Convex Programs*. *SIAM J. Optim.* 19(3), pp. 1211–1230, doi:10.1137/07069821X.
- [15] Nathalie Cauchi & Alessandro Abate (2019): *StocHy: Automated Verification and Synthesis of Stochastic Processes*. In: *TACAS (2), Lecture Notes in Computer Science* 11428, Springer, pp. 247–264, doi:10.1007/978-3-030-17465-1\_14.
- [16] Nathalie Cauchi, Luca Laurenti, Morteza Lahijanian, Alessandro Abate, Marta Kwiatkowska & Luca Cardelli (2019): *Efficiency through uncertainty: scalable formal synthesis for stochastic hybrid systems*. In: *HSCC*, ACM, pp. 240–251, doi:10.1145/3302504.3311805.
- [17] Xin Chen, Erika Ábrahám & Sriram Sankaranarayanan (2013): *Flow\*: An Analyzer for Non-linear Hybrid Systems*. In: *CAV, Lecture Notes in Computer Science* 8044, Springer, pp. 258–263, doi:10.1007/978-3-642-39799-818.



- [18] Xu Chu Ding, Mircea Lazar & Calin Belta (2012): *Formal Abstraction of Linear Systems via Polyhedral Lyapunov Functions*. In: *ADHS, IFAC Proceedings Volumes 45*, Elsevier, pp. 88–93, doi:10.3182/20120606-3-NL-3011.00096.
- [19] Oswaldo Luiz Valle Do Costa, Ricardo Paulino Marques & Marcelo Dutra Fragoso (2005): *Discrete-Time Markov Jump Linear Systems*. Springer, doi:10.1007/b138575.
- [20] Chuchu Fan, Zengyi Qin, Umang Mathur, Qiang Ning, Sayan Mitra & Mahesh Viswanathan (2022): *Controller Synthesis for Linear System With Reach-Avoid Specifications*. *IEEE Trans. Autom. Control.* 67(4), pp. 1713–1727, doi:10.1109/TAC.2021.3069723.
- [21] Jaime F. Fisac, Mo Chen, Claire J. Tomlin & S. Shankar Sastry (2015): *Reach-avoid problems with time-varying dynamics, targets and constraints*. In: *HSCC*, ACM, pp. 11–20, doi:10.1145/2728606.2728612.
- [22] Robert Givan, Sonia M. Leach & Thomas L. Dean (2000): *Bounded-parameter Markov decision processes*. *Artif. Intell.* 122(1-2), pp. 71–109, doi:10.1016/S0004-3702(00)00047-3.
- [23] Sofie Haesaert, Petter Nilsson, Cristian Ioan Vasile, Rohan Thakker, Ali-akbar Agha-mohammadi, Aaron D. Ames & Richard M. Murray (2018): *Temporal Logic Control of POMDPs via Label-based Stochastic Simulation Relations*. In: *ADHS, IFAC-PapersOnLine 51*, Elsevier, pp. 271–276, doi:10.1016/j.ifacol.2018.08.046.
- [24] Ernst Moritz Hahn, Tingting Han & Lijun Zhang (2011): *Synthesis for PCTL in Parametric Markov Decision Processes*. In: *NASA Formal Methods, Lecture Notes in Computer Science 6617*, Springer, pp. 146–161, doi:10.1007/978-3-642-20398-5\_12.
- [25] Hans Hansson & Bengt Jonsson (1994): *A Logic for Reasoning about Time and Reliability*. *Formal Aspects Comput.* 6(5), pp. 512–535, doi:10.1007/BF01211866.
- [26] Panqanamala Ramana Kumar & Pravin Varaiya (2015): *Stochastic systems: Estimation, identification, and adaptive control*. SIAM, doi:10.1137/1.9781611974263.
- [27] Marta Z. Kwiatkowska, Gethin Norman & David Parker (2011): *PRISM 4.0: Verification of Probabilistic Real-Time Systems*. In: *CAV, Lecture Notes in Computer Science 6806*, Springer, pp. 585–591, doi:10.1007/978-3-642-22110-147.
- [28] Morteza Lahijanian, Sean B. Andersson & Calin Belta (2015): *Formal Verification and Synthesis for Discrete-Time Stochastic Systems*. *IEEE Trans. Autom. Control.* 60(8), pp. 2031–2045, doi:10.1109/TAC.2015.2398883.
- [29] Abolfazl Lavaei, Sadegh Soudjani, Alessandro Abate & Majid Zamani (2022): *Automated verification and synthesis of stochastic hybrid systems: A survey*. *Automatica* 146, p. 110617, doi:10.1016/j.automatica.2022.110617.
- [30] Brian Paden, Michal Cap, Sze Zheng Yong, Dmitry S. Yershov & Emilio Frazzoli (2016): *A Survey of Motion Planning and Control Techniques for Self-Driving Urban Vehicles*. *IEEE Trans. Intell. Veh.* 1(1), pp. 33–55, doi:10.1109/TIV.2016.2578706.
- [31] Sangwoo Park, Erchin Serpedin & Khalid A. Qaraqe (2013): *Gaussian Assumption: The Least Favorable but the Most Useful [Lecture Notes]*. *IEEE Signal Process. Mag.* 30(3), pp. 183–186, doi:10.1109/MSP.2013.2238691.
- [32] Andre Platzer (2012): *Logics of Dynamical Systems*. In: *LICS*, IEEE Computer Society, pp. 13–24, doi:10.1109/LICS.2012.13.
- [33] Alberto Puggelli, Wenchao Li, Alberto L. Sangiovanni-Vincentelli & Sanjit A. Seshia (2013): *Polynomial-Time Verification of PCTL Properties of MDPs with Convex Uncertainties*. In: *CAV, Lecture Notes in Computer Science 8044*, Springer, pp. 527–542, doi:10.1007/978-3-642-39799-835.
- [34] Martin L. Puterman (1994): *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics, Wiley, doi:10.1002/9780470316887.
- [35] Luke Rickard, Thom S. Badings, Licio Romao, Nils Jansen & Alessandro Abate (2022): *Formal Controller Synthesis for Markov Jump Linear Systems with Uncertain Dynamics*. *CoRR* abs/2212.00679, doi:10.48550/arXiv.2212.00679.

- [36] Licio Romao, Antonis Papachristodoulou & Kostas Margellos (2023): *On the Exact Feasibility of Convex Scenario Programs With Discarded Constraints*. *IEEE Trans. Autom. Control*. 68(4), pp. 1986–2001, doi:10.1109/TAC.2022.3165320.
- [37] Fedor Shmarov & Paolo Zuliani (2015): *ProbReach: verified probabilistic delta-reachability for stochastic hybrid systems*. In: *HSCC*, ACM, pp. 134–139, doi:10.1145/2728606.2728625.
- [38] Sadeh Esmail Zadeh Soudjani & Alessandro Abate (2013): *Adaptive and Sequential Gridding Procedures for the Abstraction and Verification of Stochastic Processes*. *SIAM J. Appl. Dyn. Syst.* 12(2), pp. 921–956, doi:10.1137/120871456.
- [39] Eric M. Wolff, Ufuk Topcu & Richard M. Murray (2012): *Robust control of uncertain Markov Decision Processes with temporal logic specifications*. In: *CDC*, IEEE, pp. 3372–3379, doi:10.1109/CDC.2012.6426174.

# Towards Formal Fault Injection for Safety Assessment of Automated Systems \*

Ashfaq Farooqui, Behrooz Sangchoolie

Dependable Transport Systems, RISE Research Institutes of Sweden, Borås, Sweden

{ashfaq.farooqui, behrooz.sangchoolie}@ri.se

Reasoning about safety, security, and other dependability attributes of autonomous systems is a challenge that needs to be addressed before the adoption of such systems in day-to-day life. *Formal methods* is a class of methods that mathematically reason about a system's behavior. Thus, a correctness proof is sufficient to conclude the system's dependability. However, these methods are usually applied to abstract models of the system, which might not fully represent the actual system. *Fault injection*, on the other hand, is a testing method to evaluate the dependability of systems. However, the amount of testing required to evaluate the system is rather large and often a problem. This vision paper introduces *formal fault injection*, a fusion of these two techniques throughout the development lifecycle to enhance the dependability of autonomous systems. We advocate for a more cohesive approach by identifying five areas of mutual support between formal methods and fault injection. By forging stronger ties between the two fields, we pave the way for developing safe and dependable autonomous systems. This paper delves into the integration's potential and outlines future research avenues, addressing open challenges along the way.

## 1 Introduction

Safety- and security-critical systems continue to be integrated into our daily lives. Ensuring the safety and security is a multi-disciplinary challenge, where design, development, and evaluation play a crucial role. Thus a strong emphasis on updating current engineering practices to create an end-to-end verification and validation process that integrates all safety and security concerns into a unified approach is key to adopt these systems [23]. Already, *formal methods* and *fault injection* are used in different parts of the development lifecycle to ensure the system is safe and dependable.

*Formal methods* refers to mathematically rigorous techniques for specifying and verifying software and hardware systems. To many researchers, the necessity of formal methods is now a given [41]. However, this has yet to be the case from an industrial perspective. Several reasons have been suggested for this situation, including a lack of accessible tools, high costs, incompatibility with existing development techniques, and the fact that these methods require a certain level of mathematical sophistication [21, 22, 7].

*Fault injection (FI)*, on the other hand, is an established method used for the measurement, test, and assessment of dependable computer systems in extreme stress or faulty conditions. Functional safety standards such as IEC 61508 [9] and ISO 26262 [38] recommend the use of FI to prove that malfunctions in electrical and/or electronic systems will not lead to violations of safety requirements. In comparison to

---

\*This work was partly supported by the VALU3S project, which has received funding from the ECSEL Joint Undertaking (JU) under grant agreement No 876852. The JU receives support from the European Union's Horizon 2020 research and innovation programme and Austria, Czech Republic, Germany, Ireland, Italy, Portugal, Spain, Sweden, Turkey. This work has also been partly financed by the CyReV project, which is funded by the VINNOVA FFI program – the Swedish Governmental Agency for Innovation Systems (Diary number: 2019-03071).

formal methods, FI could be done after the complete system is built. This way, FI could be used to study the impact of a fault in one system and its propagation and impact in the complete end to end system. FI is also used to evaluate security properties of computer systems by means of *attack injection*. Avizienis et al. [4], consider an attack to be a special type of fault which is human made, deliberate and malicious, affecting hardware or software from external system boundaries and occurring during the operational phase.

Formal methods for safety analysis and fault injection are complementary techniques, and the choice of approach depends on the specific system being evaluated and the goals of the evaluation. Each of these methods has its benefits and shortcomings. While formal methods are commonly employed in the early design phase, fault injection testing is performed towards the later stages of development, where the system or its simulation exists. Unfortunately, the knowledge gained from the formal methods at the design phase is rarely reused in other development lifecycle phases when conducting fault injection experiments. Conversely, feedback from fault injection analysis is rarely used to improve the formal design of the system. The underlying problem is the need for common semantics and knowledge sharing across the different communities. It is clear that to deal with modern autonomous systems, formal methods and fault injection will have to be integrated in a smart way to be able to specify, verify, and validate systems, and be understandable by people without a background in formal methods. This last point is essential for autonomous systems, where they must be certified before they can be used.

By looking at the current state of applied research within both fields this paper introduces *formal fault injection* to help develop dependable autonomous systems.

## 2 Integrating Formal Methods and Fault Injection

The following section highlights the different research directions leveraging the existing state-of-the-art. Additionally, we highlight the new possibilities that will open up as a consequence of this integration and our research plans.

1. **Design level fault analysis:** Some studies focus on analyzing fault impacts during the design phase, using formal methods. Automating Failure Mode and Effects Analysis (FMEA) and Fault Tree Analysis (FTA) for system safety analysis [6, 37, 30], are well-established for within the formal community. Other studies concentrate on developing control strategies to safeguard systems against cyber-security attacks [39, 17, 27]. These approaches primarily target system design and are implemented during early development phases. Despite the industry standards of FMEA and FTA techniques, integrating formal methods into the development cycle has encountered limited adoption, mainly due to the scarcity of practical tools. Furthermore, these methods often face computational challenges, restricting their applicability to only a small portion of the system. Consequently, they typically address a fixed set of faults. It is valuable to explore how insights from the fault injection phase could be incorporated at a more generic abstraction level to streamline analysis process early in the design phase.
2. **Learning the behavior of the faulty system via model learning/learning-based testing:** Model learning [13, 42, 26] seeks to devise techniques for acquiring discrete formal models of systems by observing or interacting with them. These techniques engage in iterative testing of the system, or its simulation, to progressively learn the behavior. Model learning techniques are often applied in combination with other tools and methodologies: model checking [8] for model verification, testing methods [26] to rigorously assess system behavior, and supervisory synthesis [13] to derive supervisory controllers for system control.

we envision integrating fault analysis into the model learning phase. This analysis approach offers insights into a system's fault-handling capabilities. The resultant model encompasses both nominal and faulty behaviors, enabling offline safety analysis. Moreover, these models hold the potential to serve as authoritative proof of system safety, offering a resource for regulatory authorities.

Existing model learning tools, already applied in select industrial contexts [14, 19], provide interfaces external systems. However, an investigation is warranted to devise methods for introducing fault models into these tools and subsequently evaluating their utility. An overarching challenge lies in the scale of the resulting models. Nominal models themselves often attain considerable complexity, leading to challenges associated with state-space explosion. Augmenting these models with fault scenarios is sure to encounter state-space limitations, even for relatively modest systems.

3. **Using formal models for reducing the fault space:** Executing exhaustive fault injection campaigns is not practical. In most cases, such an approach would result in executions that do not contribute significantly to safety analysis. The challenge lies in identifying the optimal set of test instances that effectively analyze a system's dependability. Numerous testing methods have been proposed to address this challenge, including probabilistic approaches [20], coverage-based techniques [10], and heuristic as well as machine learning-based methods [29, 28, 35, 36]. However, most of these methods rely on some level of prior knowledge about the target system. The availability of such knowledge poses limitations in practical applications [28, 20].

Yet, when formal specifications exist for a specific system, they inherently contain valuable information that can be leveraged. These specifications offer insights into critical faults and their configurations, which are most likely to lead to system failures. This knowledge about fault configurations proves invaluable when designing fault campaigns. Unfortunately, such utilization is often overlooked.

To address this gap, we propose an investigation into the development of common semantics that can harmonize formal specifications and fault injection techniques. Additionally, we recommend the creation of tools to facilitate the seamless integration of results from both methods. This integration holds the potential to enhance the effectiveness of fault analysis while leveraging the rich information contained within formal specifications.

4. **Falsification for fault analysis:** Falsification methods typically come into play once the system has been implemented, typically in the later stages of the development lifecycle. Since falsification and fault injection share similar approaches and setups, integrating both these methods represents one of the most straightforward way towards a formal fault injection analysis. Both these methods operate with limited knowledge of the system under test. Falsification primarily focuses on testing the input space, while fault injection adopts a broader perspective by also assuming the presence of fault(s) within the system.

Several generalized methods and tools have proven effective in the formal community for these purposes. For instance, tools like Scenic and VerifyAI [15, 12] take a probabilistic approach to generate scenarios intelligently and test the system. Tools like Breach [11] and HyConf [1] interface with MATLAB/Simulink models for falsification. In the realm of fault injection, AV-fuzzer [24] aligns closely with falsification approaches.

While most falsification approaches aim to find input values that lead to violations of the system's specifications, they typically do not distinguish between valid and faulty inputs. To assess a system's safety, falsification engines can be employed to identify boundaries within the state-space.

Subsequently, fault injection analysis focuses on these boundary values and faulty inputs to analyze the impact of faults. In this context, exploring techniques to augment the falsification engine with both nominal and faulty behavior represents a promising avenue for further investigation.

5. **A formal specification language for fault injection:** To bridge the envisioned integration of formal methods and fault injection into practical application, a critical missing element is a well-defined *formal specification language* that can act as an interface between the two domains. Therefore, it becomes imperative to delve into the realm of formal specification languages from both theoretical and practical standpoints. A precedent is set by Bessayah et al. [5], who successfully demonstrated the use of Hoare Logic [18] as a specification language for implementing fault injection in communication systems. However, there have been limited efforts to evaluate the suitability of such a language for autonomous systems. Hence, we propose a comprehensive exploration of available formal specification languages to assess their compatibility with fault injection methodologies and to pinpoint areas of research inquiry apart from studies to develop such a language.

### 3 Formal Fault Injection

In the past few years, a paradigm called “shift-to-left” has inspired researchers to go towards simulation-based and model-based verification and validation of automated systems. The rise of simulation-based development is a key reason we believe this to be the right time to start investigating the integration of formal methods and fault injection-based testing. Notably, simulation-based methodologies have firmly established themselves in both the formal methods and fault injection communities. This shared foundation provides a common ground for the implementation of the proposed integration methods.

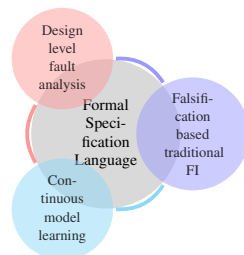


Figure 1: Conceptual overview of the proposed methodology.

Figure 1 provides a birds-eye view of the proposed integration. Several methodologies exist in literature and practice that define the development lifecycle, such as V-method [16], Waterfall [33], and Agile [32], to name a few. Most of these methodologies include the three phases: design, implementation, and testing phases, in an iterative manner and are depicted in Figure 2. This cycle of development is valid at various abstraction levels of the product lifecycle from the initial conceptual design, feature development, simulation and the final product development. During each phase, a particular set of methods as discussed in Section 2 can be mapped to the phases and are depicted using the similarly colored bubbles in Figure 1. These different phases are never isolated and are continuously updated with feedback from one another. The formal specification language makes it possible for this feedback to be easily translatable between the different phases.

It is crucial to recognize that there is no universal solution applicable to all scenarios. Therefore, the proposal does not revolve around creating a singular specification language and its corresponding

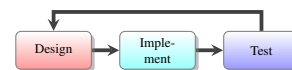


Figure 2: Generalized overview of the development lifecycle.

toolkit. Instead, it presents a high-level methodology for evaluating and constructing dependable systems. It acknowledges that specific system characteristics may demand different formalisms. Hence, the idea is to establish a collection of formal specification languages, each supported by its toolset, to facilitate this methodology. Furthermore, these languages, tools, and application approaches can vary across industries, necessitating a multifaceted strategy to address existing limitations and explore new possibilities. Above all, the aim is to establish a consistent and reproducible framework for assessing system dependability.

The assertion of a system's dependability must always be substantiated by the possibility of reproducing the results of all conducted tests. By formally specifying the entire injection methodology, it becomes possible to perform analysis, and potential replication or extension of the results by interested parties. The overarching vision of this work is to enable a standardized interface for all stakeholders involved in ensuring system safety. For instance, developers and companies can employ formal proofs to demonstrate a product's dependability, governmental certification agencies can utilize available data to enhance certification processes, and third-party auditors can scrutinize systems from security and safety perspectives using existing information. Aligning different phases of the development lifecycle with a common language paves the way for formalizing safety evaluations.

Furthermore, ongoing national and international efforts within the autonomous driving domain aim to define and develop a safety assurance framework [40, 31, 34] for verification and validation of autonomous systems. These methods aim to develop a database of testing scenarios to validate a system. Therefore, in addition to developing tools and techniques, we propose investigating the feasibility of recommending formal fault injection as a best practice through responsible standardization organizations.

## 4 Insights from early experiments

In this section we share our initial experiences and insights from applying formal techniques in simulation based fault injection. Although these experiences are common within the formal community, we believe they offer valuable insights to those interested in bridging the formal and fault injection domains.

### 4.1 The case study

Maleki and Sangchoolie [25] investigated the effects of faults on Advanced Driver Assistance Systems using the Simulation of Urban Mobility (SUMO) [2]. We use this work as a basis to study the integration of formal techniques and fault injection. The scenario used in that work [25] revolves around a three-lane road, spanning 750 meters. Two vehicles, a leader and a follower, navigate this road. The overarching requirements mandate that these vehicles not collide, successfully traverse the road, and maintain a speed not exceeding 36 m/s—the maximum permissible speed. Furthermore, these requirements should endure even when additional vehicles are introduced to the scenario, thus preserving safety and functionality amidst traffic.

To enhance the above with formal methods, we explore the following approaches.

- *Utilizing SAT Techniques for Vehicle Controller Modeling*: This approach involves modeling the vehicle models from SUMO into a SAT solver. By doing so, we enable the solver to identify potential (faulty) parameters that could lead to the violation of requirements. This approach effectively narrows down the fault space that needs to be tested. Subsequently, these identified inputs can be verified within SUMO to assess their impact on the system's behavior.

- *Applying Model Learning for Faulty Model Generation:* Here, we connect SUMO to a model learning tool controlled over TCP/IP and allow the system to learn an abstracted model that closely describes the behavior of the simulation.

## 4.2 Insights

Below we provide some of our insights from early experiments on the previously mentioned case study.

1. **Finding suitable abstractions:** One of the most significant takeaways from our study was the challenge of achieving a suitable abstraction level that aligns with the use case. As the system was implemented within the SUMO framework, we encountered limitations in deriving valuable insights from higher levels of abstraction. Operating at elevated levels of abstraction led to limited applicability of formal analysis to the fault injection process. Specifically, our attempts at model learning within the SUMO context resulted in excessively intricate models, often never terminating.

For instance, to learn a faulty model, employing a model learning approach akin to Angluin's  $L^*$  algorithm [3] necessitated the establishment of a system *alphabet*. In this context, the alphabet signifies the set of symbols providing context to the system states and transitions. Ensuring that this alphabet adeptly captures both faulty and non-faulty state transitions demanded creative thinking and held a pivotal role in shaping the effectiveness of the acquired model.

2. **The state-space explosion problem:** Both the SAT-based and model learning approaches encountered state space explosion. This challenge emerged due to the logical abstraction and simulation granularity, both significantly influencing the efficiency and effectiveness of the formal analysis. A key disparity between fault injection analysis and formal analysis became evident. While fault injection focuses solely on inputs and their corresponding outputs, formal models encompass all inputs and potential outputs, unveiling the system's comprehensive behavior. Consequently, what initially appears as the challenge of examining a finite set of input parameters in fault injection transforms into the state space issue within formal methods. Striking a balance between these two extremes emerges as the sought-after equilibrium.
3. **Non-deterministic behavior:** SUMO is a deterministic simulator. While constructing the formal model, however, this determinism is lost when augmented with faulty parameters. Additionally, the fault injection community employs randomness within faulty inputs to assess system dependability. This introduces complexity in constructing models that can effectively accommodate such randomness, necessitating techniques like abstraction or alternative formalisms. The non-deterministic aspect might not be immediately evident, potentially yielding unfavorable outcomes if appropriate analysis methods are not selected.

## 5 Conclusions

In summary, this paper introduces the concept of formal fault injection, an approach that synergizes formal methods and fault injection techniques to enhance the dependability and safety of autonomous systems. By harmonizing development and evaluation phases, this approach facilitates cross-phase knowledge sharing, fostering a unified approach to ensuring dependability attributes. This alignment not only benefits certification bodies and developers but also strengthens the foundation for proving system reliability.



## References

- [1] Arend Aerts, Mohammad Reza Mousavi & Michel Reniers (2015): *A Tool Prototype for Model-Based Testing of Cyber-Physical Systems*. In Martin Leucker, Camilo Rueda & Frank D. Valencia, editors: *Theoretical Aspects of Computing - ICTAC 2015*, Lecture Notes in Computer Science, Springer International Publishing, Cham, p. 563–572, doi:10.1007/978-3-319-25150-9\_32.
- [2] Pablo Alvarez Lopez, Michael Behrisch, Laura Bieker-Walz, Jakob Erdmann, Yun-Pang Flötteröd, Robert Hilbrich, Leonhard Lüken, Johannes Rummel, Peter Wagner & Evamarie Wiessner: *Microscopic Traffic Simulation using SUMO*. doi:10.1109/ITSC.2018.8569938.
- [3] Dana Angluin (1987): *Learning regular sets from queries and counterexamples*. *Information and Computation* 75(2), pp. 87 – 106, doi:10.1016/0890-5401(87)90052-6.
- [4] A. Avizienis, J.-C. Laprie, B. Randell & C. Landwehr (2004): *Basic concepts and taxonomy of dependable and secure computing*. *IEEE Transactions on Dependable and Secure Computing* 1(1), pp. 11–33, doi:10.1109/TDSC.2004.2.
- [5] Fayçal Bessayah, Ana Cavalli & Eliane Martins (2009): *A formal approach for specification and verification of fault injection process*. In: *Proceedings of the 2nd International Conference on Interaction Sciences: Information Technology, Culture and Human*, ICIS '09, Association for Computing Machinery, New York, NY, USA, p. 883–890. Available at <https://doi.org/10.1145/1655925.1656086>.
- [6] Marco Bozzano, Alessandro Cimatti, Cristian Mattarei & Stefano Tonetta (2014): *Formal Safety Assessment via Contract-Based Design*. In Franck Cassez & Jean-François Raskin, editors: *Automated Technology for Verification and Analysis*, Lecture Notes in Computer Science, Springer International Publishing, Cham, p. 81–97, doi:10.1007/978-3-319-11936-6\_7.
- [7] Holloway C. Michael (1997): *Why Engineers Should Consider Formal Methods*. Technical Report, NASA Langley Technical Report Server, doi:10.1109/DASC.1997.635021.
- [8] Sofia Cassel, Falk Howar, Bengt Jonsson & Bernhard Steffen (2016): *Active learning for extended finite state machines*. *Formal Aspects of Computing* 28(2), pp. 233–263, doi:10.1007/s00165-016-0355-5.
- [9] International Electrotechnical Commission (2010): *IEC 61508: Functional safety of electrical/electronic/programmable electronic safety-related systems*. Available at <https://webstore.iec.ch/publication/5515>.
- [10] M. Cukier, D. Powell & J. Ariat (1999): *Coverage estimation methods for stratified fault-injection*. *IEEE Transactions on Computers* 48(7), p. 707–723, doi:10.1109/12.780878.
- [11] Alexandre Donzé (2010): *Breach, A Toolbox for Verification and Parameter Synthesis of Hybrid Systems*. In Tayssir Touili, Byron Cook & Paul Jackson, editors: *Computer Aided Verification*, Lecture Notes in Computer Science, Springer, Berlin, Heidelberg, p. 167–170, doi:10.1007/978-3-642-14295-6\_17.
- [12] Tommaso Dreossi, Daniel J. Fremont, Shromona Ghosh, Edward Kim, Hadi Ravanbakhsh, Marcell Vazquez-Chanlatte & Sanjit A. Seshia (2019): *VerifAI: A Toolkit for the Formal Design and Analysis of Artificial Intelligence-Based Systems*, p. 432–442. *Lecture Notes in Computer Science* 11561, Springer International Publishing, Cham, doi:10.1007/978-3-030-25540-4\_25. Available at [http://link.springer.com/10.1007/978-3-030-25540-4\\_25](http://link.springer.com/10.1007/978-3-030-25540-4_25).
- [13] Ashfaq Farooqui (2021): *On supervisor synthesis via active automata learning*. Chalmers Tekniska Hogskola (Sweden). Available at <https://research.chalmers.se/en/publication/523934>.
- [14] Ashfaq Farooqui, Fredrik Hagebring & Martin Fabian (2021): *MIDES: A Tool for Supervisor Synthesis via Active Learning*. In: *2021 IEEE 17th International Conference on Automation Science and Engineering (CASE)*, pp. 792–797, doi:10.1109/CASE49439.2021.9551435.
- [15] Daniel J. Fremont, Edward Kim, Tommaso Dreossi, Shromona Ghosh, Xiangyu Yue, Alberto L. Sangiovanni-Vincentelli & Sanjit A. Seshia (2020): *Scenic: A Language for Scenario Specification and Data Generation* (arXiv:2010.06580). doi:10.48550/arXiv.2010.06580. Available at <http://arxiv.org/abs/2010.06580>. ArXiv:2010.06580 [cs].

- [16] Iris Graessler & Julian Hentze (2020): *The new V-Model of VDI 2206 and its validation*. at - *Automatisierungstechnik* 68(5), p. 312–324, doi:10.1515/auto-2020-0015. Available at <https://www.degruyter.com/document/doi/10.1515/auto-2020-0015/html>.
- [17] Christoforos N. Hadjicostis, Stéphane Lafortune, Feng Lin & Rong Su (2022): *Cybersecurity and Supervisory Control: A Tutorial on Robust State Estimation, Attack Synthesis, and Resilient Control*. In: *2022 IEEE 61st Conference on Decision and Control (CDC)*, p. 3020–3040, doi:10.1109/CDC51059.2022.9992966.
- [18] C A R Hoare (1969): *An axiomatic basis for computer programming* 12(10). doi:10.1145/363235.363259.
- [19] Malte Isberner, Falk Howar & Bernhard Steffen (2015): *The open-source LearnLib: A Framework for Active Automata Learning*. In: *International Conference on Computer Aided Verification*, Springer International Publishing, Cham, pp. 487–495, doi:10.1007/978-3-319-21690-4\_32.
- [20] Saurabh Jha, Subho Banerjee, Timothy Tsai, Siva K. S. Hari, Michael B. Sullivan, Zbigniew T. Kalbarczyk, Stephen W. Keckler & Ravishankar K. Iyer (2019): *ML-Based Fault Injection for Autonomous Vehicles: A Case for Bayesian Fault Injection*. In: *2019 49th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, p. 112–124, doi:10.1109/DSN.2019.00025.
- [21] John C. Knight (1998): *Challenges in the Utilization of Formal Methods*. In Anders P. Ravn & Hans Rischel, editors: *Formal Techniques in Real-Time and Fault-Tolerant Systems*, Lecture Notes in Computer Science, Springer, Berlin, Heidelberg, pp. 1–17, doi:10.1007/BFb0055331.
- [22] John C. Knight, Colleen L. DeJong, Matthew S. Gible & Luís G. Nakano (1997): *Why Are Formal Methods Not Used More Widely?* In: *Fourth NASA Formal Methods Workshop*, pp. 1–12.
- [23] Philip Koopman & Michael Wagner (2017): *Autonomous vehicle safety: An interdisciplinary challenge*. *IEEE Intelligent Transportation Systems Magazine* 9(1), pp. 90–96, doi:10.1109/MITS.2016.2583491.
- [24] Guanpeng Li, Yiran Li, Saurabh Jha, Timothy Tsai, Michael Sullivan, Siva Kumar Sastry Hari, Zbigniew Kalbarczyk & Ravishankar Iyer (2020): *AV-FUZZER: Finding Safety Violations in Autonomous Driving Systems*. In: *2020 IEEE 31st International Symposium on Software Reliability Engineering (ISSRE)*, IEEE, Coimbra, Portugal, p. 25–36, doi:10.1109/ISSRE5003.2020.00012. Available at <https://ieeexplore.ieee.org/document/9251068/>.
- [25] Mehdi Maleki & Behrooz Sangchoolie (2021): *SUFI: A Simulation-based Fault Injection Tool for Safety Evaluation of Advanced Driver Assistance Systems Modelled in SUMO*. In: *2021 17th European Dependable Computing Conference (EDCC)*, pp. 45–52, doi:10.1109/EDCC53658.2021.00014.
- [26] Karl Meinke & Muddassar A Sindhu (2013): *LBTest: a learning-based testing tool for reactive systems*. In: *2013 IEEE Sixth International Conference on Software Testing, Verification and Validation*, IEEE, pp. 447–454, doi:10.1109/ICST.2013.62.
- [27] Rômulo Meira-Goes, Eunsuk Kang, Raymond H. Kwong & Stéphane Lafortune (2020): *Synthesis of sensor deception attacks at the supervisory layer of Cyber-Physical Systems*. *Automatica* 121, p. 109172, doi:10.1016/j.automatica.2020.109172. Available at <https://linkinghub.elsevier.com/retrieve/pii/S0005109820303708>.
- [28] Mehrdad Moradi, Bentley Oakes & Joachim Denil: *Machine Learning-assisted Fault Injection*. Available at <https://hal.laas.fr/hal-02931709v1>.
- [29] Mehrdad Moradi, Bentley James Oakes, Mustafa Saraoglu, Andrey Morozov, Klaus Janschek & Joachim Denil (2020): *Exploring Fault Parameter Space Using Reinforcement Learning-based Fault Injection*. In: *2020 50th Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W)*, p. 102–109, doi:10.1109/DSN-W50199.2020.00028.
- [30] Frank Ortmeier & Gerhard Schellhorn (2007): *Formal Fault Tree Analysis - Practical Experiences*. *Electronic Notes in Theoretical Computer Science* 185, pp. 139–151, doi:10.1016/j.entcs.2007.05.034. Available at <https://www.sciencedirect.com/science/article/pii/S1571066107004549>. Proceedings of the 6th International Workshop on Automated Verification of Critical Systems (AVoCS 2006).
- [31] PEGASUS: *PEGASUS project*. Available at <https://www.pegasusprojekt.de/en/home>.

- [32] D.J. Reifer (2002): *How good are agile methods?* *IEEE Software* 19(4), p. 16–18, doi:10.1109/MS.2002.1020280.
- [33] Dr Winston W Rovce: *MANAGING THE DEVELOPMENT OF LARGE SOFTWARE SYSTEMS*. Available at <https://dl.acm.org/doi/10.5555/41765.41801>.
- [34] SAKURA: *SAKURA project*. Available at [https://www.sakura-prj.go.jp/project\\_info/](https://www.sakura-prj.go.jp/project_info/).
- [35] Behrooz Sangchoolie, Karthik Pattabiraman & Johan Karlsson (2022): *An Empirical Study of the Impact of Single and Multiple Bit-Flip Errors in Programs*. *IEEE Transactions on Dependable and Secure Computing* 19(3), pp. 1988–2006, doi:10.1109/TDSC.2020.3043023.
- [36] Ali Sedaghatbaf, Mehrdad Moradi, Jaafar Almasizadeh, Behrooz Sangchoolie, Bert Van Acker & Joachim Denil (2022): *DELFASE: A Deep Learning Method for Fault Space Exploration*. In: *2022 18th European Dependable Computing Conference (EDCC)*, pp. 57–64, doi:10.1109/EDCC57035.2022.00020.
- [37] Yuvaraj Selvaraj, Zhennan Fei & Martin Fabian (2020): *Supervisory Control Theory in System Safety Analysis*. In António Casimiro, Frank Ortmeier, Erwin Schoitsch, Friedemann Bitsch & Pedro Ferreira, editors: *Computer Safety, Reliability, and Security. SAFECOMP 2020 Workshops*, Lecture Notes in Computer Science, Springer International Publishing, Cham, p. 9–22, doi:10.1007/978-3-030-55583-2\_1.
- [38] International Organization for Standardization (2018): *ISO 26262: Road vehicles — Functional safety*. Available at <https://www.iso.org/standard/68383.html>.
- [39] Rong Su (2018): *Supervisor synthesis to thwart cyber attack with bounded sensor reading alterations*. *Automatica* 94, p. 35–44, doi:10.1016/j.automatica.2018.04.006. Available at <https://www.sciencedirect.com/science/article/pii/S0005109818301912>.
- [40] SUNRISE: *SUNRISE project*. Available at <https://sunrise-europe.eu/>.
- [41] Jim Woodcock, Peter Gorm Larsen, Juan Bicarregui & John Fitzgerald (2009): *Formal methods: Practice and experience*. *ACM Computing Surveys* 41(4), p. 1–36, doi:10.1145/1592434.1592436. Available at <https://dl.acm.org/doi/10.1145/1592434.1592436>.
- [42] H. Zhang, L. Feng & Z. Li (2018): *A Learning-Based Synthesis Approach to the Supremal Nonblocking Supervisor of Discrete-Event Systems*. *IEEE Trans. on Automatic Control* 63(10), pp. 3345–3360, doi:10.1109/TAC.2018.2793662.

# Formal Verification of Long Short-Term Memory based Audio Classifiers: A Star based Approach

Neelanjana Pal

Institute for Software Integrated Systems  
Vanderbilt University  
Nashville, USA  
neelanjana.pal@vanderbilt.edu

Taylor T Johnson

Institute for Software Integrated Systems  
Vanderbilt University  
Nashville, USA  
taylor.johnson@vanderbilt.edu

Formally verifying audio classification systems is essential to ensure accurate signal classification across real-world applications like surveillance, automotive voice commands, and multimedia content management, preventing potential errors with serious consequences. Drawing from recent research, this study advances the utilization of star-set-based formal verification, extended through reachability analysis, tailored explicitly for Long Short-Term Memory architectures and their Convolutional variations within the audio classification domain. By conceptualizing the classification process as a sequence of set operations, the star set-based reachability approach streamlines the exploration of potential operational states attainable by the system. The paper serves as an encompassing case study, validating and verifying sequence audio classification analytics within real-world contexts. It accentuates the necessity for robustness verification to ensure precise and dependable predictions, particularly in light of the impact of noise on the accuracy of output classifications.

## 1 Introduction

Deep Neural Networks (DNNs) have demonstrated remarkable capabilities in addressing intricate tasks like image classification, object detection, speech recognition, natural language processing, and document analysis, at times even surpassing human performance [21, 23, 24]. This success has ignited a surge in exploring the viability of DNNs across diverse real-world domains, including biometric authentication, mobile facial recognition for security, and malware detection. However, given the sensitive nature of the data in these critical applications, incorporating safety, security, and robust verification into their design has become paramount.

However, studies have revealed that even slight modifications in input data can effectively mislead cutting-edge, well-trained networks, causing inaccuracies in their predictions [12, 32, 40]. The arena of network verification has primarily concentrated on image inputs, particularly emphasizing the assurance of safety and robustness in various classification neural networks [2, 7, 19, 31, 43, 44]. Previous investigations have scrutinized a range of network architectures, encompassing feed-forward neural networks (FFNNs [42]), convolutional neural networks (CNNs [44]), semantic segmentation networks (SSNs [43]), and a few using Recurrent Neural Networks (RNNs [41]) employing diverse set-based reachability tools such as Neural Network Verification (NNV [26, 45]) and JuliaReach [6], among others.

Models utilizing NNs for audio classification have found application in diverse tasks, ranging from Music Genre Classification [8, 10, 11] and Environmental Sound Classification [4, 9, 13] to Audio Generation [33, 36]. Therefore, formal verification of audio classification systems holds paramount importance in ensuring their reliability and safety, particularly in safety-critical applications such as autonomous vehicles [35, 46], medical diagnosis [15, 30], and industrial monitoring [47].

This study introduces an extension, building upon the foundations laid by two recent studies [34, 41] in the domain of formal verification. The objective is to leverage set-based reachability techniques to

verify audio classification models based on the Long Short Term Memory (LSTM) and CNN-LSTM architectures. Drawing inspiration from [41], which highlights the star-based verification of basic vanilla RNNs, and from [34], which demonstrates the formal verification of convolutional neural networks operating on time series data, work shown in this paper amalgamates both concepts. Specifically, it employs two LSTM models and one CNN-LSTM model for these classifications, following the ones depicted in [27–29].

### Contributions.

1. This paper presents a thorough case study on the formal verification of audio classification models using the LSTM and CNN-LSTM architectures with two different datasets. Our focus is to rigorously assess the robustness verification of these models within a formal verification framework, analyzing their behavior and performance against input noises. We develop our work as an extension of the NNV tool<sup>1</sup> to formally analyze and explore CNN-LSTM architecture verification for audio data using sound and deterministic reachability methods.
2. Building on insights from existing research [34,41], this paper extends formal verification to more complex RNN architectures. This involves addressing the challenges of the complex structure of the LSTM layers, comprehensively evaluating their behavior, and ensuring robustness compliance through formal verification. This study pushes formal verification’s boundaries, embracing design complexities for heightened assurance and reliability.
3. In this assessment, we conduct a thorough and comprehensive evaluation of three distinct network architectures across diverse audio classification scenarios.
4. Finally, we develop insights on evaluating the reachability analysis on those networks and possible future direction.

**Outline.** The paper is organized as follows: Section 2 mentions the works already done in the literature and the inspiration works for this paper; Section 3 provides the necessary context for the background, and defines the verification properties for this work; Section 4 explains the reachability calculations for the LSTM layer; and Section 5 describes the methodology, including dataset, network models, and input perturbations. Section 6 presents the experimental results, evaluation metrics, and their implications. Finally, Section 7 summarizes the main findings and suggests future research directions.

## 2 Related Work

In recent times, an upsurge of methodologies and tools have arisen to confront the verification complexities inherent in intricate systems like Deep Neural Networks (DNNs), as evident from the literature [14, 17, 25, 44]. Correspondingly, tools have emerged to tackle the robustness challenges of Convolutional Neural Networks (CNNs) [2, 19, 20, 37–39]. Earlier undertakings in the verification of Recurrent Neural Networks (RNNs) are showcased through projects like RnnVerify [18] and RNSVerify [1]. RNSVerify employs an unrolling technique to translate RNNs into extensive Feedforward Neural Networks (FFNNs), simplifying verification through Mixed-Integer Linear Program (MILP) approaches [1]. However, this unrolling method faces scalability constraints, particularly with bounded n-step RNNs, as

---

<sup>1</sup>The code for this paper is available at <https://github.com/verivital/nnv/tree/master/code/nnv/examples/Submission/FMAS2023>

the verification complexity scales dramatically. Conversely, RnnVerify [18] employs invariant inference for RNN verification, bypassing unrolling. Their strategy involves crafting an FFNN with matching dimensions to over-approximate the RNN, followed by verifying RNN properties over this approximation using SMT-based methodologies. Our work gets inspiration from [41], where authors introduce a pioneering approach founded on star reachability for RNN verification, aiming to amplify the dependability and safety of RNNs and show the results based on some vanilla RNN models.

**Distinction from the previous works [34,41].** While both papers share the common goal of validating the robustness of RNNs, the preceding study can be perceived as an initial step in that research trajectory. In contrast, this paper represents a more comprehensive evolution of the concepts initially introduced.

1. The work in [41] focused on the Vanilla RNN, while this paper delves into models of significantly greater intricacy, such as the LSTM and CNN-LSTM architectures. Vanilla RNNs and LSTMs are both types of recurrent neural networks. However, Vanilla RNNs are simpler in terms of architecture and have fewer parameters, whereas LSTMs are more complex due to their gated units and larger parameters.
  - (a) Vanilla RNNs handle input sequences in a sequential manner, updating a hidden state at each step. In contrast, LSTMs also maintain a hidden state, but their structure is more complex, featuring multiple gates (input, forget, and output gates) that regulate the information flow.
  - (b) Vanilla RNNs face challenges capturing long-term dependencies within sequences due to the vanishing gradient problem. This problem limits their ability to learn connections between distant time steps. LSTMs, on the other hand, were specifically designed to tackle the vanishing gradient problem and excel at capturing long-term dependencies, rendering them better suited for tasks involving intricate temporal relationships.
  - (c) Vanilla RNNs possess a constrained memory capacity, often rapidly discarding information from earlier time steps. This limitation can hinder their performance in tasks with extended sequences. In contrast, LSTMs feature an improved memory mechanism that allows them to retain or discard information from prior time steps selectively. This capability equips them to handle longer sequences and capture complex patterns effectively.
2. The study conducted in [34] focused on examining time-series regression models in the Prognostics and Health Management domain. Drawing inspiration from this work, our study extends the investigation to encompass the time domain's influence, specifically concerning sequential audio noise and Japanese vowel audio samples. This basic experiment provides a foundation for understanding the robustness and reliability of audio classification systems. They offer insights that can be directly applied to real-world scenarios, making them valuable for a broad audience in the field of audio classification.
  - (a) Utilizing real-world datasets, this experiment can yield practical insights into audio classification system performance, benefiting fields such as speech recognition, audio surveillance, and multimedia content management by offering real-world applicability.
  - (b) This experiment can provide valuable insights into the robustness of audio classifiers when exposed to different noise levels and perturbations, offering crucial implications for applications where audio data is frequently affected by a noise like voice commands in automobiles or audio analysis in noisy settings.
  - (c) While this work concentrates on two particular datasets, the verification methodologies showcased can be extended to diverse audio classification endeavors, allowing readers engaged in

various audio classification challenges to customize and apply the methodologies to their specific contexts.

- (d) The metrics used in this paper can be potential for real-world applications to evaluate and enhance the reliability and efficiency of audio classification systems.

### 3 Preliminaries

This section introduces some basic definitions and descriptions necessary to understand the progression of this paper and the necessary evaluations on audio classification models.

#### 3.1 Neural Network Verification Tool and Star Sets

The Neural Network Verification (NNV) tool constitutes a framework designed to verify the safety and robustness of neural networks [26,45]. This tool meticulously scrutinizes neural network behavior across diverse input conditions, warranting secure and accurate functionality across all scenarios. NNV employs reachability algorithms, including the exact and over-approximate star set methodologies [42, 44], to compute reachable sets for each neural network layer. These sets encapsulate all feasible network states for a given input, thereby facilitating the verification of specific safety properties.

NNV holds particular significance in safety-critical domains like autonomous vehicles and medical devices, where the trustworthiness and reliability of neural networks are paramount. NNV bolsters public confidence in these applications by ensuring consistent performance across all conditions. In this paper, we extend the capabilities of the NNV tool to implement our work, utilizing the star-based reachability analysis to ascertain the reachable sets of neural networks at their outputs.

$$\Theta = c + \alpha v = \begin{matrix} \begin{matrix} 0 & 4 & 1 & 2 \\ 2 & 3 & 2 & 3 \\ 1 & 3 & 1 & 2 \\ 2 & 1 & 3 & 2 \end{matrix} \\ c \in \mathbb{R}^{4 \times 4} \end{matrix} + \alpha \begin{matrix} \begin{matrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \\ v \in \mathbb{R}^{4 \times 4} \end{matrix}, P \equiv \begin{pmatrix} 1 \\ -1 \end{pmatrix} \alpha \leq \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

Figure 1: Star for a sequence input data with four Feature Values (rows) with four time-steps (columns)

**Definition 3.1** A generalized star set (or simply star)  $\Theta$  is a tuple  $\langle c, V, P \rangle$  where  $c \in \mathbb{R}^n$  is the center,  $V = \{v_1, v_2, \dots, v_m\}$  is a set of  $m$  vectors in  $\mathbb{R}^n$  called basis vectors, and  $P: \mathbb{R}^m \rightarrow \{\top, \perp\}$  is a predicate. The basis vectors are arranged to form the star's  $n \times m$  basis matrix. The set of states represented by the star is given as:

$$[[\Theta]] = \{x \mid x = c + \sum_{i=1}^m (\alpha_i v_i) \text{ and } P(\alpha_1, \dots, \alpha_m) = \top\}. \tag{1}$$

In this work, we restrict the predicates to be a conjunction of linear constraints,  $P(\alpha) \triangleq C\alpha \leq d$  where, for  $p$  linear constraints,  $C \in \mathbb{R}^{p \times m}$ ,  $\alpha$  is the vector of  $m$ -variables, i.e.,  $\alpha = [\alpha_1, \dots, \alpha_m]^T$ , and  $d \in \mathbb{R}^{p \times 1}$ .

## 3.2 Network Architecture Specifics

### 3.2.1 Long Short Term Memory (LSTM) Layer

An LSTM layer, a subtype of the Recurrent Neural Network (RNN) layer, excels at capturing long-term dependencies in time series and sequential data [16]. It comprises two critical elements: the hidden state ( $h_t$ , also called the output state) and the cell state ( $c_t$ ). At each time step ‘t,’ the hidden state captures the layer’s output for that instance, while the cell state accumulates insights from preceding time steps.

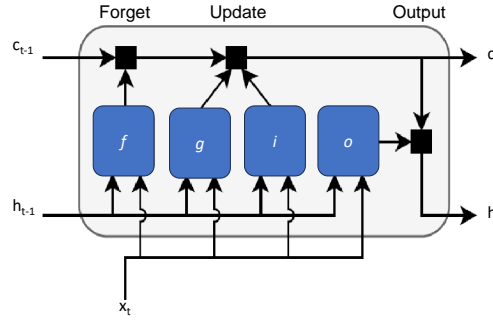


Figure 2: The flow of data at time step t in an LSTM layer

$$\begin{aligned} c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ h_t &= o_t \odot \sigma_c(c_t) \end{aligned} \quad (2)$$

During each time step, the layer refines the cell state by incorporating or omitting information. This process is steered by distinct gates that control these adjustments, as shown in Fig. 2.

$$\begin{aligned} i_t &= \sigma_g(W_i x_t + R_i h_{t-1} + b_i) \\ f_t &= \sigma_g(W_f x_t + R_f h_{t-1} + b_f) \\ g_t &= \sigma_c(W_g x_t + R_g h_{t-1} + b_g) \\ o_t &= \sigma_g(W_o x_t + R_o h_{t-1} + b_o) \end{aligned} \quad (3)$$

In these equations,  $\odot$  represents the Hadamard product (element-wise multiplication),  $\sigma_c$  denotes the activation function applied element-wise to the cell state  $c_t$  and to the cell state gate  $g_t$ ;  $\sigma_g$  denotes the activation function applied element-wise to the hidden state gates. Here,  $W$ ,  $R$ , and  $b$  are, respectively, hidden state weights, recurrent weights, and biases for each of the gates.

### 3.2.2 Convolutional Neural Network + Long Short Term Memory (CNN+LSTM) Architecture

When processing sequences, a CNN uses sliding convolutional filters over the input, extracting information from spatial and temporal dimensions. Conversely, an LSTM network progresses through time steps, capturing lasting connections between them. The synergy of CNN and LSTM layers, as seen in CNN+LSTM architectures [49], harnesses the strengths of both convolutional and LSTM units for insightful data analysis.

The convolutional component forms the foundation for acquiring local feature modules that grasp both local and hierarchical correlations. This fusion enables the identification of intricate data relationships. Additionally, the inclusion of an LSTM layer enhances the network’s capacity to capture prolonged dependencies by leveraging information from these localized features.



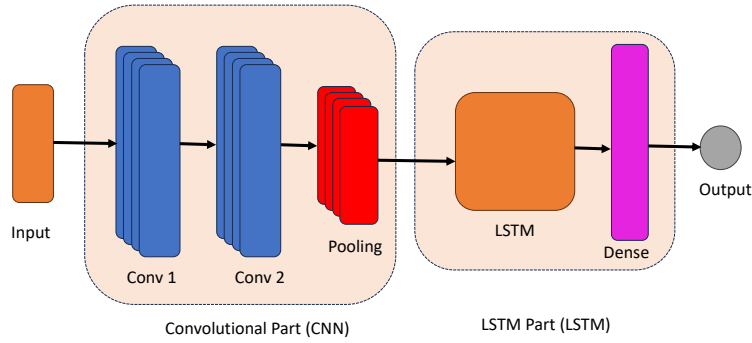


Figure 3: Layers of a demo CNN+LSTM Architecture model

### 3.3 Reachability Analysis Computation

This section describes how the reachability of an NN layer and the NN as a whole is computed for this study.

In this context, we adopt an alternative technique for defining a Star set. This method involves utilizing the input's upper and lower bounds with noise, subsequently aligning them around the original input. We establish a comprehensive array of constraints by incorporating these bounds for each input parameter alongside predicates. These constraints are then presented to the optimizer for a solution, ultimately yielding the initial set of states.

**Definition 3.2** A *layer*  $L$  of a NN is a function  $h : u \in \mathbb{R}^j \rightarrow v \in \mathbb{R}^p$ , defined as follows

$$v = h(u) \quad (4)$$

where the function  $h$  is determined by parameters  $\theta$ , typically defined as a tuple  $\theta = \langle \sigma, W, b \rangle$  for fully-connected and convolutional layers, where  $W$  is the weight matrix,  $b$  is the bias vector, and activation function is  $\sigma$ . For CNN layers,  $\theta$  may include parameters like the filter size, padding, or dilation factor.

**Definition 3.3** Let  $h : u \in \mathbb{R}^j \rightarrow v \in \mathbb{R}^p$ , be an NN layer as described in Eq. 4. The *reachable set*  $\mathcal{R}_h$ , with input,  $\mathcal{I} \in \mathbb{R}^n$  is defined as

$$\mathcal{R}_h \triangleq \{v \mid v = h(u), u \in \mathcal{I}\} \quad (5)$$

**Reachability analysis (or, shortly, reach) of an NN  $f$  on Star input set  $\mathcal{I}$**  is similar to the reachable set calculations for CNN [44] or FFNN [42].

$$\text{Reach}(f, \mathcal{I}) : \mathcal{I} \rightarrow \mathcal{R}_{ts} \quad (6)$$

We call  $\mathcal{R}_{ts}(I)$  the *output reachable set* of the NN corresponding to the input set  $\mathcal{I}$ .

For an NN, the output reachable set can be calculated as a step-by-step process of constructing the reachable sets for each network layer.

$$\begin{aligned} \mathcal{R}_{L_1} &\triangleq \{v_1 \mid v_1 = h_1(x), x \in \mathcal{I}\}, \\ \mathcal{R}_{L_2} &\triangleq \{v_2 \mid v_2 = h_2(v_1), v_1 \in \mathcal{R}_{L_1}\}, \\ &\vdots \\ \mathcal{R}_{ts} = \mathcal{R}_{L_k} &\triangleq \{v_k \mid v_k = h_k(v_{k-1}), v_{k-1} \in \mathcal{R}_{L_{k-1}}\}, \end{aligned} \quad (7)$$

where  $h_k$  is the function represented by the  $k^{th}$  layer  $L_k$ . The reachable set  $\mathcal{R}_{L_k}$  contains all outputs of the neural network corresponding to all input vectors  $x$  in the input set  $\mathcal{I}$ .

### 3.4 Adversarial Perturbation

An audio classification system may face real-world scenarios involving elements like background noise, interference, or distortions. While potentially perceptible, these factors remain within the scope of challenges that practical systems are designed to address. However, this paper exclusively used l-infinity perturbations, focusing on assessing how audio classification models respond to variations within specific constraints.

Considering an input sequence characterized by  $t_s$  time instances and  $n_f$  features, various perturbation types ( $l_\infty$  norm) [34] arise based on their distribution across the sequence. These adversarial perturbation categories can be delineated as follows:

1. **Single Feature Single-instance (SFSI):** This entails perturbing the value of a specific feature solely at a particular instance ( $t$ ), deviating by a certain percentage from the actual value:

$$s^{perturb} = g_{\mathcal{E},s^{perturb}}(s) = s + \mathcal{E}_t \cdot s_t^{perturb} \quad (8)$$

2. **Single Feature All-instances (SFAI):** In this scenario, a particular feature across all time instances undergoes perturbation by a certain percentage relative to its original values:

$$s^{perturb} = g_{\mathcal{E},s^{perturb}}(s) = s + \sum_{i=1}^n \mathcal{E}_i \cdot s_i^{perturb} \quad (9)$$

3. **Multifeature Single-instance (MFSI):** All feature values experience perturbation, but solely at a specific instance ( $t$ ), following the principle outlined in Eq. 8 for each feature.
4. **Multifeature All-instance (MFAI):** Perturbation affects all feature values across all instances, aligning with the approach delineated in Eq. 9 for every feature.

### 3.5 Robustness Verification Properties

**Robustness.** Robustness pertains to the capacity of a system or model to sustain its performance and functionality amid diverse challenging conditions, uncertainties, or perturbations. This highly desirable trait ensures the system's dependability, resilience, and adaptability in the presence of altering or unfavorable circumstances. To formally articulate the concept of robustness for quantifying the desired classification task, the following formulation can be employed:

$$\|x' - x\|_\infty < \delta \implies f(x') == f(x) \quad (10)$$

Here,  $x$  signifies the original input from the input space  $R^{n_f \times t_s}$ ,  $x'$  represents the perturbed input,  $f(x')$  and  $f(x)$  correspond to the classifiers' outputs for  $x'$  and  $x$ , respectively.  $\delta$  stands for the maximum magnitude of the introduced perturbation ( $\delta \in \mathbf{R} > 0$ ). By disregarding the softmax and classification layers within the models and focusing on the output of the layer immediately preceding the softmax, the formulation for robustness simplifies as follows:

$$\|x' - x\|_\infty < \delta \implies \maxID(g(x')) == \maxID(g(x)) \quad (11)$$

In this context, the function  $g$  symbolizes the operation performed by the neural network classifier model until the softmax layer, and  $\maxID$  denotes the function responsible for identifying the class with the highest value in the output.

**Verification Properties.** Verification properties can be broadly classified into two distinct categories: local and global. A local property must be valid for specific predefined inputs, while a global property [48] is established across the entire input space  $R^{n_f \times t_s}$  of the network model, holding true for all inputs without exceptions.

**Local Robustness.** Given a sequence classifier  $f$  and an input sequence  $S$ , the network is called **locally robust** to any perturbation  $\mathcal{A}$  if and only if: reachable bounds of the desired class will be max compared to the bounds of the other classes, even in the presence of any perturbation.

**Robustness Value (RV)** of a sequence  $S$  is a binary variable, which indicates the local robustness of the system. RV is 1 when the reachable output range of the desired class is greater than the reachable bounds of other classes, making it locally robust; otherwise, RV is 0.

$$RV = 1 \iff LB_{desired} \geq UB_{other} \text{ else, } RV = 0$$

where  $LB_{desired}$  and  $UB_{other}$  are the lower reachable bound of the desired class and  $UB_{other}$  are the upper bounds of all other classes.

**Percentage Robustness (PR).** We apply the concept of Percentage Robustness (PR), previously utilized in image-based classification or segmentation neural networks [43], to the context of sequence audio inputs. The PR for a sequence classifier, corresponding to any adversarial perturbation, is defined as:

$$PR = \frac{N_{robust}}{N_{total}} \times 100 \quad (12)$$

where  $N_{robust}$  represents the total number of robust sequences, and  $N_{total}$  is the overall count of sequences in the test dataset. Percentage robustness can be used as an indicator of **global robustness** [48] with respect to various types of perturbations.

## 4 Reachability of a Long Short Term Memory Layer

To compute the reachability of an LSTM layer in relation to a star input set  $S_t$ , a series of stepwise reachability computations are necessary to ultimately determine the reachable set of the LSTM layer's output, as depicted in Eq. 2-3. Ensuring accurate results relies on verifying the validity of specific conditions, which are crucial for this process to be sound and accurate:

1. **Affine Mapping Validity.** The transformation of a star set through an affine mapping using a given weight and bias must result in another valid star set [42].
2. **Star Set Summation.** Combining two star sets through Minkowski summation should lead to the formation of yet another valid star set [5].
3. **Activation Function Application.** Upon applying the activation function to a star set, the output should also result in a star set(s). The outcome could manifest as a single star set or a composition of multiple star sets, contingent on factors such as the activation functions employed and the specific reachability technique utilized [42–45].
4. **Hadamard Product Validity.** The Hadamard Product of two star sets should yield another valid star set.

While the validity of the first three conditions for star sets has been established in prior research, this current study aims to extend that validation to include the fourth condition as well.

**Definition 4.1 (Hadamard product of two star sets)** Given two star-sets  $\Theta_1 = \langle c_1, V_1, P_1 \rangle$  and  $\Theta_2 = \langle c_2, V_2, P_2 \rangle$ , the Hadamard product of them  $\bar{\Theta} = \Theta_1 \odot \Theta_2 = \{y \mid y = x_1 \odot x_2, x_1 \in \Theta_1, x_2 \in \Theta_2\}$  is another star with the following characteristics.

$$\bar{\Theta} = \langle \bar{c}, \bar{V}, \bar{P} \rangle, \bar{c} = c_1 \odot c_2, \bar{V} = \begin{bmatrix} V_1 & 0 \\ 0 & V_2 \end{bmatrix}, \bar{P} \equiv \begin{bmatrix} P_1 & 0 \\ 0 & P_2 \end{bmatrix}$$

Therefore we can conclude that for a given input set  $S_i$  and an LSTM layer, the output is also a star set.

## 5 Experimental Setup

### 5.1 Hardware Used:

The actual experimental results shown in this paper are conducted in a Windows-10 computer with the 64-bit operating system, Intel(R) Core(TM) i7-8850H processor, and 16 GB RAM.

### 5.2 Dataset Description

For evaluation, we consider two different audio datasets for noise classification and Japanese vowel classification.

**Audio Noise Data:** To curate this dataset, we generated a collection of 1000 white noise signals, 1000 brown noise signals, and 1000 pink noise signals using MATLAB. Each signal corresponds to a 0.5-second duration and adheres to a 44.1 kHz sample rate. From this pool of 1000 signals, a training set is fashioned, comprising 800 white noise signals, 800 brown noise signals, and 800 pink noise signals. Given the multidimensionality inherent in audio data, often containing redundant information, a dimensionality reduction strategy is employed. We begin by extracting features and subsequently training the model using only two extracted features. These features are generated from the centroid and slope of the mel spectrum over time.

**Japanese Vowel [3, 22]:** This dataset is collected from [3] from the University of Irvine Machine Learning Repository. Two Japanese vowels were sequentially pronounced by nine male speakers. A 12-degree linear prediction analysis was subjected to each instance of utterances. Each speaker’s utterance constitutes a time series ranging from 7 to 29 points in length, with each point featuring 12 coefficients. For 9 classes (i.e., vowels), the dataset has a total of 640 time series. Among these, 270 time series were designated for training purposes, while the remaining 370 were allocated for testing.

### 5.3 Network Description

**Audio Noise Data:** The network architecture used for training the audio noise dataset, partially adopted from [27], is an LSTM network. The network has two input features which correspond to one noise type at the output. Following 11, the network for this dataset can be represented as:

$$f: x \in \mathbb{R}^{2 \times l_s} \rightarrow y \in \mathbb{R}^3$$

$$\text{noise}\hat{C}lass = \max ID(g(x)) \tag{13}$$

**Japanese Vowel:** Here we have trained two different classifiers for the Japanese Vowel dataset. The LSTM architecture is partially adopted from [28] and the CNN+LSTM is partially adopted from [29]. Both the networks have twelve input features which correspond to one vowel at the output. Therefore, the networks for this dataset can be represented as:

$$f : x \in \mathbb{R}^{12 \times l_s} \rightarrow y \in \mathbb{R}^9$$

$$\text{vowelClass} = \text{maxID}(g(x)) \quad (14)$$

Here  $l_s$  is the audio sequence length and the function  $\text{maxID}$  provides the class with the maximum value.

Table 1: Performances of different networks used in this paper

Networks	Accuracy(%)
<i>audio_noise_lstm</i>	100
<i>japanese_vowel_lstm</i>	93.51
<i>japanese_vowel_cnnlstm</i>	96.49

## 6 Evaluation

### 6.1 Robustness Verification of Audio Noise Classifier

To conduct robustness verification on the audio noise dataset, we encompass all four categories of perturbations, following [34]. First, we curate 100 sequences each of white, brown, and pink noise as test datasets. Then, we generate adversarial sequences centered around the original ones by applying  $l_\infty$  norms. This involves utilizing 5 distinct percentage values for perturbation ( $\epsilon$ ), specifically 50%, 60%, 70%, 80%, and 90% of the mean ( $\mu$ ) value. These newly created adversarial inputs are subjected to assessment through the exact-star reachability analysis [Sec. 4] to determine their robustness. Notably, in the case of Single Feature Single-instance Noise (SFSI) and Single Feature All-instances Noise (SFAI), we opt for random selection of feature 1 for input perturbation.

Table 2: Global Robustness: percentage robustness (PR) and total verification runtime (sumRT in seconds) for 100 test audio noise sequences

<i>noise</i>	$PR_{SFSI}$	$PR_{SFAI}$	$PR_{MFSI}$	$PR_{MFAI}$	$sumRT_{SFSI}$	$sumRT_{SFAI}$	$sumRT_{MFSI}$	$sumRT_{MFAI}$
50	98	80.33	98	80.33	0.3071	0.2626	0.3018	0.2625
60	96	71.67	96	71.67	0.3034	0.2571	0.3018	0.2578
70	94	25.67	94	25.67	0.3039	0.2637	0.3070	0.2663
80	85.33	12.33	85.33	12.33	0.3073	0.2559	0.3111	0.2556
90	63.33	8.33	63.33	8.33	0.3060	0.2504	0.3093	0.2537

**Observations and Analysis.** Table 2 and Fig. 4 present the network’s overall performance, i.e., the percentage robustness measures, PR [Sec. 3.5], and total verification runtime (sumRT) in seconds, with respect to each adversarial perturbation. The observations derived from both the table and the figure provide the following insights:

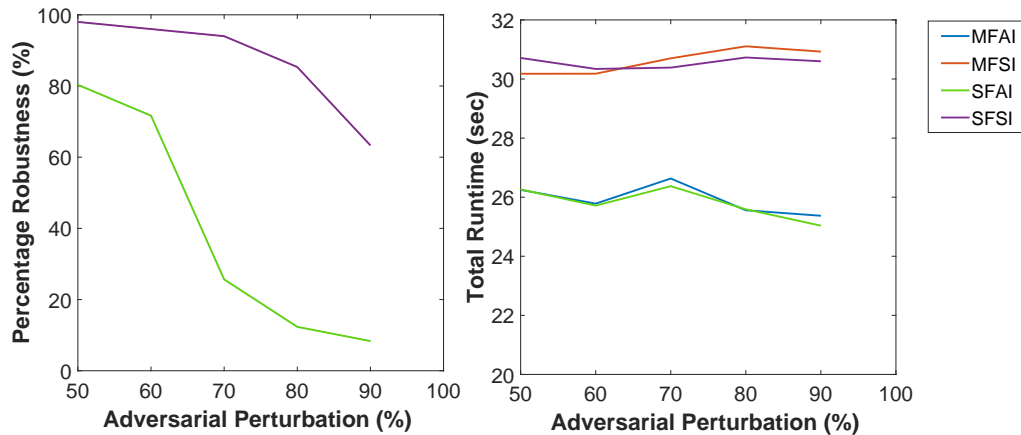


Figure 4: Percentage Robustness and Runtime plots w.r.t increasing perturbations

1. **Trend of Percentage Robustness (PR).** As the adversary level increases from 50 to 90, we observe a consistent decrease in PR values for all perturbation scenarios (SFSI, SFAI, MFSI, MFAI), which aligns with the concept of the robustness verification property. This decrease in PR signifies a reduction in the system’s ability to maintain its classification accuracy in the presence of higher adversary levels.
2. **Comparative Analysis of Perturbation Scenarios.** Within each noise level, comparing PR values across different perturbation scenarios (SFSI, SFAI, MFSI, MFAI), it’s evident that PR values for SFSI and MFSI are generally higher than those for SFAI and MFAI. This finding indicates that perturbing features at a single instance or all features at a single instance generally leads to better robustness against varying noise levels.
3. **Similar PR Values for Different Perturbation Scenarios.** Another notable observation is the similarity in robustness matrices between SFSI and MFSI scenarios, accompanied by closely comparable computation times for their respective verification processes. This parallelism is also evident for SFAI and MFAI perturbations as well. This pattern could be ascribed to the dataset’s limited feature set of only two dimensions, where the foremost feature likely holds paramount importance in influencing the class determination in the presence of noise. Consequently, when single-instance perturbations target the first feature, perturbing both features results in an effect akin to perturbing the first feature alone. This interpretation is applicable to both MFAI and SFAI scenarios as well.

## 6.2 Robustness Verification of Japanese Vowel Classifiers

To verify the robustness of both the LSTM and the CNN+LSTM models in the context of the Japanese vowel classifier, we extend the evaluation to encompass all four perturbation categories, mirroring the approach undertaken for the audio noise classifier. During this procedure, we focus on the complete set of correctly classified test sequences. Subsequently, we create adversarial inputs centered around the original sequences by applying  $l_\infty$  norms to evaluate their robustness. This perturbation process involves applying 5 distinct percentage values ( $\epsilon$ ) for perturbation: specifically, 50%, 60%, 70%, 80%,

and 90% of the mean ( $\mu$ ) value. The resulting set of adversarial inputs is then assessed using the exact-star reachability analysis for both the classifiers to ascertain their robustness. Like the earlier scenario, for SFSI and SFAI, feature 1 is chosen for perturbation.

Table 3: Global Robustness: percentage robustness (PR) and total verification runtime (sumRT in seconds) for all test Japanese Vowel audio sequences

<i>noise</i>	$PR_{SFSI}$	$PR_{SFAI}$	$PR_{MFSI}$	$PR_{MFAI}$	$sumRT_{SFSI}$	$sumRT_{SFAI}$	$sumRT_{MFSI}$	$sumRT_{MFAI}$
50	100	68.21	100	74.86	1.1502	1.0398	1.0392	1.0442
60	100	60.40	100	50.29	0.9989	0.9992	0.9970	0.9966
70	100	50.29	100	27.17	0.9981	0.9968	0.9965	0.9952
80	100	43.93	100	13.01	0.9920	0.9930	0.9978	0.9882
90	100	39.02	100	8.38	1.0044	1.0083	1.004	0.9985

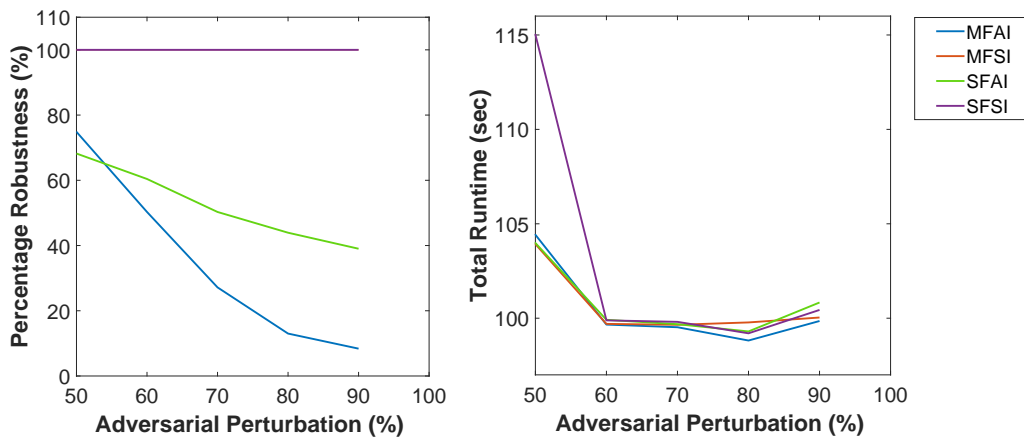


Figure 5: Percentage Robustness and Runtime plots w.r.t increasing perturbations, for LSTM architecture

**Observations and Analysis: LSTM Model** Table 3 and Fig. 5 present the LSTM network’s overall performance, i.e., the percentage robustness measures, PR [Sec. 3.5], and total verification runtime (sumRT), with respect to each adversarial perturbation. The notable findings are outlined as follows

1. **Trend of Percentage Robustness (PR).** Similar to the audio noise classifier, the trends in PR values here also suggest that as noise levels increase, the percentage robustness tends to decrease across all scenarios. This aligns with the intuitive expectation that higher adversary levels lead to increased challenges in maintaining robustness.

The  $PR_{SFSI}$  and  $PR_{MFSI}$  values remain consistently at 100% across all noise levels, indicating that perturbing either a single feature or all features at a specific instance does not significantly affect the robustness of the audio sequences. On the other hand,  $PR_{SFAI}$  and  $PR_{MFAI}$  show distinct trends. As adversary levels increase,  $PR_{SFAI}$  gradually decreases, suggesting that perturbing all instances but only a single feature starts impacting the robustness. Similarly,  $PR_{MFAI}$  also experiences a decline with increasing noise levels, reflecting that perturbing all instances and features has an impact on the sequences’ robustness.

2. **Comparative Analysis of Perturbation Scenarios.** The comparison between single-instance perturbation scenarios (*SFSI* and *SFAI*) and multifeature perturbation scenarios (*MFSI* and *MFAI*) reveals a pattern. The former scenarios (single-instance) generally maintain higher robustness compared to the latter (multifeature) scenarios. This suggests that perturbing all features has a larger impact on robustness than perturbing just a single feature.

The interrelation between  $PR_{SFSI}$  and  $PR_{MFSI}$  is also notable. Both scenarios exhibit identical trends, regardless of the noise level. Similarly,  $PR_{SFAI}$  and  $PR_{MFAI}$  also demonstrate similar behaviors, with both scenarios showing a decline in robustness as noise increases.

**Observations and Analysis: CNN+LSTM Model** Table 4 and Fig. 6 present the CNN+LSTM network's overall performance.

Table 4: Global Robustness: percentage robustness (PR) and total verification runtime (sumRT in seconds) for all correctly-classified test Japanese Vowel audio sequences

<i>noise</i>	$PR_{SFSI}$	$PR_{SFAI}$	$PR_{MFSI}$	$PR_{MFAI}$	$sumRT_{SFSI}$	$sumRT_{SFAI}$	$sumRT_{MFSI}$	$sumRT_{MFAI}$
50	96.82	49.13	97.39	65.89	5.5019	4.2007	4.2197	4.1148
60	96.82	40.75	97.39	43.64	4.5148	3.9238	3.9123	3.9200
70	96.82	34.68	97.10	18.78	4.6223	4.0966	4.0778	4.0626
80	96.82	30.63	97.10	3.17	4.5658	4.0998	4.0550	4.0714
90	96.82	26.58	97.10	0	4.5773	4.0785	4.0715	4.0605

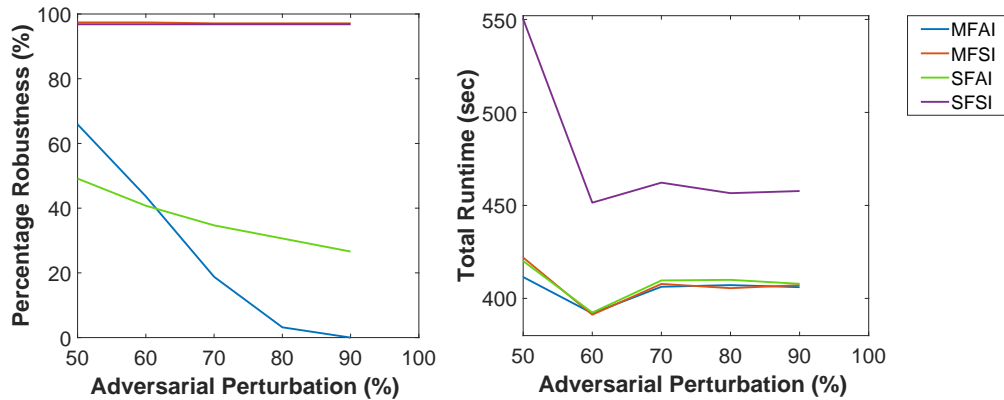


Figure 6: Percentage Robustness and Runtime plots w.r.t increasing perturbations, for LSTM architecture

Key insights gleaned from both the table and the plot include:

1. **Trend of Percentage Robustness (PR).** Across all perturbation levels, the  $PR_{SFSI}$  remain consistently at around 96% and the  $PR_{MFSI}$  at around 97%, indicating that the perturbations applied in these scenarios do not significantly affect the robustness of the audio sequences. For SFAI and MFAI perturbations, PR also decreases with rising noise levels, although the decline is more pronounced. PR values for SFSI and MFSI perturbations are significantly higher compared to SFAI



and MFAI perturbations at all noise levels, indicating that sequences with perturbations at a single instance are more robust to noise.

2. **Trend of Verification Runtimes ( $sumRT$ ).** Verification runtimes tend to rise with elevated noise levels across all perturbation scenarios. However, in the case of the Japanese Vowel dataset, an initial decrease is observed in the runtime trend, followed by an increase at perturbation level 70% and then again decreases at 80%, followed by another increase at 90%. It's also worth noting that contrary to the expected trend,  $sumRT_{SFSI}$  exhibits a higher runtime value in comparison to  $sumRT_{SFAI}$  and  $sumRT_{MFAI}$ .

Overall, the above tables demonstrate how different perturbation scenarios and adversary levels impact the percentage robustness of the audio noise and Japanese Vowel audio classifiers. The trends and inter-relations provide insights into the varying effects of perturbations on different scenarios and noise levels, helping to understand the robustness behavior of the neural network models under different conditions.

## 7 Conclusion and Future Directions

This study delves into formal method-based reachability analysis for various LSTM-based neural networks (NNs) using exact and approximate Star methods, specifically in the context of audio sequence classification – a critical aspect for safety-critical applications. The investigation encompasses four distinct adversarial perturbation types, as introduced in the existing literature. The evaluation occurs across two audio sequence datasets: audio noise sequences and Japanese vowel audio sequences. The unified reachability analysis accommodates shifting features within time sequences while scrutinizing the output against the desired audio class. Robustness properties are verified for both datasets. Although real-world datasets are employed, further research is essential to strengthen the connection between practical issues and performance metrics. The evaluation can also be conducted with multiple repetitions to ensure that the reported results are not dependent on specific instances or random fluctuations, thus enhancing the overall validity and reliability of the findings. Exploring real-world scenarios encompassing a wider array of perturbation types and magnitudes will also be fascinating, potentially yielding diverse effects on system behavior. The study paves the way for exploring the impact of perturbations on the output and expanding reachability analysis to three-dimensional sequence data like videos. An intriguing direction for exploration can involve analyzing the peculiar runtime patterns observed in the plots for the Japanese Vowel audio dataset. Potential future applications can also encompass medical video analysis. Notably, this work concentrates on offline data analysis, omitting considerations for real-time stream processing and memory limitations, which offers intriguing prospects for future investigation.

**Acknowledgements.** The material presented in this paper is based upon work supported by the National Science Foundation (NSF) through grant numbers 1910017, 2028001, 2220426, and 2220401, and the Defense Advanced Research Projects Agency (DARPA) under contract number FA8750-18-C-0089 and FA8750-23-C-0518, and the Air Force Office of Scientific Research (AFOSR) under contract number FA9550-22-1-0019 and FA9550-23-1-0135. Any opinions, findings, conclusions, or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of AFOSR, DARPA, or NSF.

## References

- [1] Michael E Akintunde, Andreea Kevorchian, Alessio Lomuscio & Edoardo Pirovano (2019): *Verification of RNN-Based Neural Agent-Environment Systems*. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, 33, pp. 6006–6013, doi:[10.1609/aaai.v33i01.33016006](https://doi.org/10.1609/aaai.v33i01.33016006).
- [2] Greg Anderson, Shankara Pailoor, Isil Dillig & Swarat Chaudhuri (2019): *Optimization and abstraction: a synergistic approach for analyzing neural network robustness*. In: *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pp. 731–744, doi:[10.1145/3314221.3314614](https://doi.org/10.1145/3314221.3314614).
- [3] Arthur Asuncion & David Newman (2007): *UCI machine learning repository*.
- [4] Yusuf Aytar, Carl Vondrick & Antonio Torralba (2016): *Soundnet: Learning sound representations from unlabeled video*. *Advances in Neural Information Processing Systems 29: December 5-10, 2016, Barcelona, Spain*.
- [5] Stanley Bak & Parasara Sridhar Duggirala (2017): *Simulation-equivalent reachability of large linear systems with inputs*. In: *International Conference on Computer Aided Verification*, Springer, pp. 401–420, doi:[10.1007/978-3-319-63387-9\\_20](https://doi.org/10.1007/978-3-319-63387-9_20).
- [6] Sergiy Bogomolov, Marcelo Forets, Goran Frehse, Kostiantyn Potomkin & Christian Schilling (2019): *JuliaReach: a toolbox for set-based reachability*. In: *Proceedings of the 22nd ACM International Conference on Hybrid Systems: Computation and Control*, pp. 39–44, doi:[10.1145/3302504.3311804](https://doi.org/10.1145/3302504.3311804).
- [7] Elena Botoeva, Panagiotis Kouvaros, Jan Kronqvist, Alessio Lomuscio & Ruth Misener (2020): *Efficient verification of relu-based neural networks via dependency analysis*. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, 34, pp. 3291–3299, doi:[10.1609/aaai.v34i04.5729](https://doi.org/10.1609/aaai.v34i04.5729).
- [8] Keunwoo Choi, György Fazekas, Mark Sandler & Kyunghyun Cho (2017): *Convolutional recurrent neural networks for music classification*. In: *2017 IEEE International conference on acoustics, speech and signal processing (ICASSP)*, IEEE, pp. 2392–2396, doi:[10.1109/ICASSP.2017.7952585](https://doi.org/10.1109/ICASSP.2017.7952585).
- [9] Fatih Demir, Daban Abdulsalam Abdullah & Abdulkadir Sengur (2020): *A new deep CNN model for environmental sound classification*. *IEEE Access* 8, pp. 66529–66537, doi:[10.1109/ACCESS.2020.2984903](https://doi.org/10.1109/ACCESS.2020.2984903).
- [10] Mingwen Dong (2018): *Convolutional neural network achieves human-level accuracy in music genre classification*. *arXiv preprint arXiv:1802.09697*.
- [11] Peace Busola Falola, Emmanuel Oluwadunsin Alabi, Folashade Titilope Ogunajo & Oluwakemi Dunsin Fasae (2022): *Music genre classification using machine and deep learning techniques: a review*. *ResearchJet J Anal Invent* 3(03), pp. 35–50.
- [12] Ian J Goodfellow, Jonathon Shlens & Christian Szegedy (2015): *Explaining and harnessing adversarial examples*. *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- [13] Andrey Guzhov, Federico Raue, Jörn Hees & Andreas Dengel (2021): *Esresnet: Environmental sound classification based on visual domain models*. In: *2020 25th International Conference on Pattern Recognition (ICPR)*, IEEE, pp. 4933–4940, doi:[10.1109/ICPR48806.2021.9413035](https://doi.org/10.1109/ICPR48806.2021.9413035).
- [14] Navid Hashemi, Bardh Hoxha, Tomoya Yamaguchi, Danil Prokhorov, Georgios Fainekos & Jyotirmoy Deshmukh (2023): *A Neurosymbolic Approach to the Verification of Temporal Logic Properties of Learning-enabled Control Systems*. In: *Proceedings of the ACM/IEEE 14th International Conference on Cyber-Physical Systems (with CPS-IoT Week 2023)*, pp. 98–109, doi:[10.1145/3576841.3585928](https://doi.org/10.1145/3576841.3585928).
- [15] Ezz El-Din Hemdan, Walid El-Shafai & Amged Sayed (2023): *CR19: A framework for preliminary detection of COVID-19 in cough audio signals using machine learning algorithms for automated medical diagnosis applications*. *Journal of Ambient Intelligence and Humanized Computing* 14(9), pp. 11715–11727, doi:[10.1007/s12652-022-03732-0](https://doi.org/10.1007/s12652-022-03732-0).
- [16] Sepp Hochreiter & Jürgen Schmidhuber (1997): *Long short-term memory*. *Neural computation* 9(8), pp. 1735–1780, doi:[10.1162/neco.1997.9.8.1735](https://doi.org/10.1162/neco.1997.9.8.1735).

- [17] Xiaowei Huang, Daniel Kroening, Wenjie Ruan, James Sharp, Youcheng Sun, Emese Thamo, Min Wu & Xinping Yi (2020): *A survey of safety and trustworthiness of deep neural networks: Verification, testing, adversarial attack and defence, and interpretability*. *Computer Science Review* 37, p. 100270, doi:[10.1016/j.cosrev.2020.100270](https://doi.org/10.1016/j.cosrev.2020.100270).
- [18] Yuval Jacoby, Clark Barrett & Guy Katz (2020): *Verifying recurrent neural networks using invariant inference*. In: *Automated Technology for Verification and Analysis: 18th International Symposium, ATVA 2020, Hanoi, Vietnam, October 19–23, 2020, Proceedings 18*, Springer, pp. 57–74, doi:[10.1007/978-3-030-59152-6\\_3](https://doi.org/10.1007/978-3-030-59152-6_3).
- [19] Guy Katz, Derek A Huang, Duligur Ibeling, Kyle Julian, Christopher Lazarus, Rachel Lim, Parth Shah, Shantanu Thakoor, Haoze Wu, Aleksandar Zeljić et al. (2019): *The marabou framework for verification and analysis of deep neural networks*. In: *International Conference on Computer Aided Verification*, Springer, pp. 443–452, doi:[10.1007/978-3-030-25540-4\\_26](https://doi.org/10.1007/978-3-030-25540-4_26).
- [20] Panagiotis Kouvaros & Alessio Lomuscio (2018): *Formal verification of cnn-based perception systems*. *arXiv preprint arXiv:1811.11373*.
- [21] Alex Krizhevsky, Ilya Sutskever & Geoffrey E Hinton (2012): *Imagenet classification with deep convolutional neural networks*. *Advances in Neural Information Processing Systems 25: December 3-6, 2012, Lake Tahoe, Nevada, United States*, pp. 1106–1114.
- [22] Mineichi Kudo, Jun Toyama & Masaru Shimbo (1999): *Multidimensional curve classification using passing-through regions*. *Pattern Recognition Letters* 20(11-13), pp. 1103–1111, doi:[10.1016/S0167-8655\(99\)00077-X](https://doi.org/10.1016/S0167-8655(99)00077-X).
- [23] Steve Lawrence, C Lee Giles, Ah Chung Tsoi & Andrew D Back (1997): *Face recognition: A convolutional neural-network approach*. *IEEE transactions on neural networks* 8(1), pp. 98–113, doi:[10.1109/72.554195](https://doi.org/10.1109/72.554195).
- [24] Yann LeCun, Léon Bottou, Yoshua Bengio & Patrick Haffner (1998): *Gradient-based learning applied to document recognition*. *Proceedings of the IEEE* 86(11), pp. 2278–2324, doi:[10.1109/5.726791](https://doi.org/10.1109/5.726791).
- [25] Changliu Liu, Tomer Arnon, Christopher Lazarus, Christopher Strong, Clark Barrett, Mykel J Kochenderfer et al. (2021): *Algorithms for verifying deep neural networks*. *Foundations and Trends® in Optimization* 4(3-4), pp. 244–404, doi:[10.1561/24000000035](https://doi.org/10.1561/24000000035).
- [26] Diego Manzananas Lopez, Sung Woo Choi, Hoang-Dung Tran & Taylor T Johnson (2023): *NNV 2.0: the neural network verification tool*. In: *International Conference on Computer Aided Verification*, Springer, pp. 397–412, doi:[10.1007/978-3-031-37703-7\\_19](https://doi.org/10.1007/978-3-031-37703-7_19).
- [27] *Classify Sound Using Deep Learning - MATLAB & Simulink — mathworks.com*. <https://www.mathworks.com/help/audio/gs/classify-sound-using-deep-learning.html>.
- [28] *Sequence Classification Using Deep Learning - MATLAB & Simulink — mathworks.com*. <https://www.mathworks.com/help/deeplearning/ug/classify-sequence-data-using-lstm-networks.html>.
- [29] *Sequence Classification Using 1-D Convolutions - MATLAB & Simulink — mathworks.com*. <https://www.mathworks.com/help/deeplearning/ug/sequence-classification-using-1-d-convolutions.html>.
- [30] Toshio Modegi & Shun-ichi Iisaku (1997): *Application of MIDI technique for medical audio signal coding*. In: *Proceedings of the 19th Annual International Conference of the IEEE Engineering in Medicine and Biology Society: 'Magnificent Milestones and Emerging Opportunities in Medical Engineering'* (Cat. No. 97CH36136), 4, IEEE, pp. 1417–1420, doi:[10.1109/IEMBS.1997.756970](https://doi.org/10.1109/IEMBS.1997.756970).
- [31] Jeet Mohapatra, Tsui-Wei Weng, Pin-Yu Chen, Sijia Liu & Luca Daniel (2020): *Towards verifying robustness of neural networks against a family of semantic perturbations*. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 244–252, doi:[10.1109/CVPR42600.2020.00032](https://doi.org/10.1109/CVPR42600.2020.00032).
- [32] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi & Pascal Frossard (2016): *Deepfool: a simple and accurate method to fool deep neural networks*. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2574–2582, doi:[10.1109/CVPR.2016.282](https://doi.org/10.1109/CVPR.2016.282).

- [33] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior & Koray Kavukcuoglu (2016): *Wavenet: A generative model for raw audio*. *The 9th ISCA Speech Synthesis Workshop, Sunnyvale, CA, USA, 13-15 September 2016*, p. 125.
- [34] Neelanjana Pal, Diego Manzananas Lopez & Taylor T Johnson (2023): *Robustness verification of deep neural networks using star-based reachability analysis with variable-length time series input*. In: *International Conference on Formal Methods for Industrial Critical Systems*, Springer, pp. 170–188, doi:[10.1007/978-3-031-43681-9\\_10](https://doi.org/10.1007/978-3-031-43681-9_10).
- [35] Gunasekaran Raja, Senbagapriya Senthilkumar, Sivaseelan Ganesan, Rithika Edhayachandran, Geetha Vijayaraghavan & Ali Kashif Bashir (2021): *AV-CPS: audio visual cognitive processing system for critical intervention in autonomous vehicles*. In: *2021 IEEE International Conference on Communications Workshops (ICC Workshops)*, IEEE, pp. 1–6, doi:[10.1109/ICCWorkshops50388.2021.9473647](https://doi.org/10.1109/ICCWorkshops50388.2021.9473647).
- [36] Adam Roberts, Jesse Engel, Colin Raffel, Curtis Hawthorne & Douglas Eck (2018): *A hierarchical latent vector model for learning long-term structure in music*. In: *International conference on machine learning*, PMLR, pp. 4364–4373.
- [37] Wenjie Ruan, Min Wu, Youcheng Sun, Xiaowei Huang, Daniel Kroening & Marta Kwiatkowska (2019): *Global Robustness Evaluation of Deep Neural Networks with Provable Guarantees for the L<sub>0</sub> Norm*. *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, pp. 5944–5952, doi:[10.24963/ijcai.2019/824](https://doi.org/10.24963/ijcai.2019/824).
- [38] Gagandeep Singh, Timon Gehr, Matthew Mirman, Markus Püschel & Martin Vechev (2018): *Fast and effective robustness certification*. In: *Advances in Neural Information Processing Systems*, pp. 10825–10836.
- [39] Gagandeep Singh, Timon Gehr, Markus Püschel & Martin Vechev (2019): *An abstract domain for certifying neural networks*. *Proceedings of the ACM on Programming Languages* 3(POPL), p. 41, doi:[10.1145/3291645](https://doi.org/10.1145/3291645).
- [40] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow & Rob Fergus (2014): *Intriguing properties of neural networks*. In: *2nd International Conference on Learning Representations, ICLR 2014*.
- [41] Hoang Dung Tran, Sung Woo Choi, Xiaodong Yang, Tomoya Yamaguchi, Bardh Hoxha & Danil Prokhorov (2023): *Verification of Recurrent Neural Networks with Star Reachability*. In: *Proceedings of the 26th ACM International Conference on Hybrid Systems: Computation and Control*, pp. 1–13, doi:[10.1145/3575870.3587128](https://doi.org/10.1145/3575870.3587128).
- [42] Hoang-Dung Tran, Diago Manzananas Lopez, Patrick Musau, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang & Taylor T Johnson (2019): *Star-based reachability analysis of deep neural networks*. In: *International Symposium on Formal Methods*, Springer, pp. 670–686, doi:[10.1007/978-3-030-30942-8\\_39](https://doi.org/10.1007/978-3-030-30942-8_39).
- [43] Hoang-Dung Tran, Neelanjana Pal, Patrick Musau, Diego Manzananas Lopez, Nathaniel Hamilton, Xiaodong Yang, Stanley Bak & Taylor T Johnson (2021): *Robustness verification of semantic segmentation neural networks using relaxed reachability*. In: *International Conference on Computer Aided Verification*, Springer, pp. 263–286, doi:[10.1007/978-3-030-81685-8\\_12](https://doi.org/10.1007/978-3-030-81685-8_12).
- [44] Hoang-Dung Tran, Weiming Xiang & Taylor T Johnson (2020): *Verification approaches for learning-enabled autonomous cyber-physical systems*. *IEEE Design & Test* 39(1), pp. 24–34, doi:[10.1109/MDAT.2020.3015712](https://doi.org/10.1109/MDAT.2020.3015712).
- [45] Hoang-Dung Tran, Xiaodong Yang, Diego Manzananas Lopez, Patrick Musau, Luan Viet Nguyen, Weiming Xiang, Stanley Bak & Taylor T Johnson (2020): *NNV: The neural network verification tool for deep neural networks and learning-enabled cyber-physical systems*. In: *International Conference on Computer Aided Verification*, Springer, pp. 3–17, doi:[10.1007/978-3-030-53288-8\\_1](https://doi.org/10.1007/978-3-030-53288-8_1).
- [46] Finley Walden, Sagar Dasgupta, Mizanur Rahman & Mhafuzul Islam (2022): *Improving the Environmental Perception of Autonomous Vehicles using Deep Learning-based Audio Classification*. *CoRR* abs/2209.04075, doi:[10.48550/arXiv.2209.04075](https://doi.org/10.48550/arXiv.2209.04075).

- [47] Avery Wang et al. (2003): *An industrial strength audio search algorithm*. In: *ISMIR 2003, 4th International Conference on Music Information Retrieval, Baltimore, Maryland, USA, October 27-30, 2003, Proceedings*, Washington, DC, pp. 7–13.
- [48] Zhilu Wang, Yixuan Wang, Feisi Fu, Ruochen Jiao, Chao Huang, Wenchao Li & Qi Zhu (2022): *A Tool for Neural Network Global Robustness Certification and Training*. CoRR abs/2208.07289, doi:[10.48550/arXiv.2208.07289](https://doi.org/10.48550/arXiv.2208.07289).
- [49] Jianfeng Zhao, Xia Mao & Lijiang Chen (2019): *Speech emotion recognition using deep 1D & 2D CNN LSTM networks*. *Biomedical signal processing and control* 47, pp. 312–323, doi:[10.1016/j.bspc.2018.08.035](https://doi.org/10.1016/j.bspc.2018.08.035).

# 3vLTL: A Tool to Generate Automata for Three-valued LTL

Francesco Belardinelli

Imperial College, London, United Kingdom  
francesco.belardinelli@imperial.ac.uk

Angelo Ferrando

University of Genoa, Genoa, Italy  
angelo.ferrando@unige.it

Vadim Malvone

Telecom Paris, Paris, France  
vadim.malvone@telecom-paris.fr

Multi-valued logics have a long tradition in the literature on system verification, including run-time verification. However, comparatively fewer model-checking tools have been developed for multi-valued specification languages. We present 3vLTL, a tool to generate Büchi automata from formulas in Linear-time Temporal Logic (LTL) interpreted on a three-valued semantics. Given an LTL formula, a set of atomic propositions as the alphabet for the automaton, and a truth value, our procedure generates a Büchi automaton that accepts all the words that assign the chosen truth value to the LTL formula. Given the particular type of the output of the tool, it can also be seamlessly processed by third-party libraries in a natural way. That is, the Büchi automaton can then be used in the context of formal verification to check whether an LTL formula is true, false, or undefined on a given model.

## 1 Introduction

Multi-valued logics have a long tradition in the literature on system verification, as demonstrated by various references [8, 15, 3, 23, 17, 18], and they play a crucial role in run-time verification as well [4, 5]. Of particular interest are three-valued logics, including temporal extensions of Kleene's logic [19], where the third value, in addition to true and false, is interpreted as "unknown" or "unspecified". Such semantics prove especially convenient when constructing smaller abstractions of complex reactive and distributed systems. These abstractions are typically approximations of the original model, containing strictly less information. Consequently, the challenge lies in finding the right trade-off between reducing complexity and minimizing information loss during the abstraction process. In system verification, one of the most widely used temporal logics for specifying requirements is Linear-time Temporal Logic (*LTL*) [22]. The model checking problem for *LTL* is typically addressed through automata-theoretic techniques [2]. Given a model  $M$  of a transition system and an *LTL* formula  $\varphi$ , the approach involves generating Büchi automata for both  $M$  and the negation of  $\varphi$ . This allows us to determine whether  $\varphi$  is satisfied in  $M$  by examining whether the language accepted by the product of these two automata is empty.

Several tools are available now to generate Büchi automata from *LTL* formulas. Notable examples include [11, 13]. However, to the best of our knowledge, no tool has yet been proposed to directly generate Büchi automata for multi-valued temporal logics.

**Contribution.** In this paper we present 3vLTL, a tool to generate (generalized) Büchi automata from *LTL* formulas interpreted on a three-valued semantics. Specifically, given an *LTL* formula, a set of atomic propositions (representing the automaton alphabet), and a truth value (true, false or undefined), our procedure generates a Büchi automaton that accepts all the words that assign the chosen truth value to the input *LTL* formula. Furthermore, 3vLTL has the functionality to process our output (i.e., the

automaton) by third-party libraries in a natural way. The present work is motivated by the use of three-valued logics in system verification. Indeed, our tool can be used in several works, such as [20, 16, 24, 25], to provide results for the verification of three-valued *LTL* formulas. Furthermore, our tool is already used in [6]. In this work, the authors present an abstraction-refinement method to partially solve the model checking of multi-agent systems under imperfect information and perfect recall strategies. Note that, a three-valued semantics becomes particularly significant in situations involving imperfect information, as the absence of information can lead to the emergence of a third value. This is particularly evident in autonomous and distributed systems, where a component may not have access to the complete system's information [12].

**Related Work.** Concerning the three-valued automata technique for *LTL* employed in this work, the most closely related approaches can be found in [21, 10, 9, 26]. Notably, in [9], there is an exploration of a reduction from multi-valued to two-valued *LTL*, but it does not encompass automata-theoretic techniques. Conversely, in [10], an automata-theoretic approach for general multi-valued *LTL* is presented, following the tableau-based construction as outlined in [14]; however, this work is more suitable for on-the-fly verification w.r.t. to our approach. In a different vein, [21] delves into general multi-valued automata, defining lattices, deterministic and non-deterministic automata, as well as their extensions with Büchi acceptance conditions. As part of their theoretical findings, they introduce an automata construction for multi-valued *LTL*, though it lacks a clear explanation of states and transitions. With respect to our work, in [21], the model checking is only briefly discussed, and their approach is tailored more toward multi-valued logics in a broader sense.

To summarize, unlike [21, 10, 9], our proposed approach makes minimal modifications to the automata-theoretic construction for two-valued *LTL* [2] and extends it to a three-valued interpretation.

## 2 Preliminaries

In this part we present a three-valued semantics for Linear-time Temporal Logic *LTL* and recall the definition of generalized non-deterministic Büchi automata. To fix the notation, we assume that  $AP = \{q_1, q_2, \dots\}$  is the set of *atomic propositions*, or simply atoms. We denote the length of a tuple  $t$  as  $|t|$ , and its  $i$ -th element as  $t_i$ . For  $i \leq |t|$ , let  $t_{\geq i}$  be the suffix  $t_i, \dots, t_{|t|}$  of  $t$  starting at  $t_i$  and  $t_{\leq i}$  its prefix  $t_1, \dots, t_i$ . Notice that we start enumerations with index 1.

**Models.** We begin by giving a formal definition of Transition Model [2].

**Definition 1 (Transition Model)** *Given a set  $AP$  of atoms, a Transition Model is a tuple  $M = \langle S, s_0, \longrightarrow, V \rangle$  such that (i)  $S$  is a finite, non-empty set of states, with initial state  $s_0 \in S$ ; (ii)  $\longrightarrow \subseteq S \times S$  is a serial transition relation; (iii)  $V : S \times AP \rightarrow \{\text{tt}, \text{ff}, \text{uu}\}$  is the three-valued labelling function.*

A path  $p \in S^\omega$  is an infinite sequence  $s_1 s_2 s_3 \dots$  of states where  $s_i \longrightarrow s_{i+1}$ , for each  $i \geq 1$ .

**Syntax.** Here, we recall the syntax of *LTL*.

**Definition 2 (LTL)** *Formulas in *LTL* are defined as follows, where  $q \in AP$ :*

$$\varphi ::= q \mid \neg\varphi \mid \varphi \wedge \varphi \mid X\varphi \mid (\varphi U \varphi)$$

The meaning of operators *next*  $X$  and *until*  $U$  is standard [2]. Operators *release*  $R$ , *finally*  $F$ , and *globally*  $G$  can be introduced as usual:  $\varphi R \psi \equiv \neg(\neg\varphi U \neg\psi)$ ,  $F\varphi \equiv \text{tt} U \varphi$ ,  $G\varphi \equiv \text{ff} R \varphi$ .

**Semantics.** Formally we define the three-valued semantics for *LTL* as follows.

**Definition 3 (Satisfaction)** *The three-valued satisfaction relation  $\models^3$  for a Transition Model  $M$ , path  $p \in S^\omega$ , atom  $q \in AP$ ,  $v \in \{\text{tt}, \text{ff}\}$ , and formula  $\phi$  is defined as follows:*

$$\begin{aligned}
((M, s) \models^3 A\psi) = \text{ff} & \quad \text{iff} \quad \text{for some path } p \text{ in } M, ((M, p) \models^3 \psi) = \text{ff} \\
((M, p) \models^3 q) = v & \quad \text{iff} \quad V(p_1, q) = v \\
((M, p) \models^3 \neg\psi) = v & \quad \text{iff} \quad ((M, p) \models^3 \psi) = \neg v \\
((M, p) \models^3 \psi \wedge \psi') = \text{tt} & \quad \text{iff} \quad ((M, p) \models^3 \psi) = \text{tt} \text{ and } ((M, p) \models^3 \psi') = \text{tt} \\
((M, p) \models^3 \psi \wedge \psi') = \text{ff} & \quad \text{iff} \quad ((M, p) \models^3 \psi) = \text{ff} \text{ or } ((M, p) \models^3 \psi') = \text{ff} \\
((M, p) \models^3 X\psi) = v & \quad \text{iff} \quad ((M, p_{\geq 2}) \models^3 \psi) = v \\
((M, p) \models^3 \psi U \psi') = \text{tt} & \quad \text{iff} \quad \text{for some } k \geq 1, ((M, p_{\geq k}) \models^3 \psi') = \text{tt}, \text{ and} \\
& \quad \text{for all } j, 1 \leq j < k \Rightarrow ((M, p_{\geq j}) \models^3 \psi) = \text{tt} \\
((M, p) \models^3 \psi U \psi') = \text{ff} & \quad \text{iff} \quad \text{for all } k \geq 1, ((M, p_{\geq k}) \models^3 \psi') = \text{ff}, \text{ or} \\
& \quad \text{for some } j \geq 1, ((M, p_{\geq j}) \models^3 \psi) = \text{ff}, \text{ and} \\
& \quad \text{for all } j', 1 \leq j' \leq j \Rightarrow ((M, p_{\geq j'}) \models^3 \psi') = \text{ff}
\end{aligned}$$

In all other cases the value of  $\phi$  is *uu*.

**Generalized non-deterministic Büchi automaton.** Now, we recall the definition of the class of automata that we will use in our construction and in the tool.

**Definition 4 (GNBA)** *A generalized non-deterministic Büchi automaton is a tuple  $A = \langle Q, Q_0, \Sigma, \pi, \mathcal{F} \rangle$  where (i)  $Q$  is a finite set of states with  $Q_0 \subseteq Q$  as the set of initial states; (ii)  $\Sigma$  is an alphabet; (iii)  $\pi : Q \times \Sigma \rightarrow 2^Q$  is the (non-deterministic) transition function; (iv)  $\mathcal{F}$  is a (possibly empty) set of subsets of  $Q$ , whose elements are called acceptance sets.*

Given an infinite run  $\rho = q_0q_1q_2 \dots \in Q^\omega$ , let  $\text{Inf}(\rho)$  be the set of states  $q$  for which there are infinitely many indices  $i$  with  $q = q_i$ , that is,  $q$  appears infinitely often in  $\rho$ . Then, run  $\rho$  is *accepting* if for each acceptance set  $F \in \mathcal{F}$ ,  $\text{Inf}(\rho) \cap F \neq \emptyset$ , that is, there are infinitely many indices  $i$  in  $\rho$  with  $q_i \in F$ . The *accepted language*  $L(A)$  of automaton  $A$  consists of all infinite words  $w \in \Sigma^\omega$  for which there exists at least one accepting run  $\rho = q_0q_1q_2 \dots \in Q^\omega$  such that for all  $i \geq 0$ ,  $q_{i+1} \in \pi(q_i, w_i)$ .

### 3 Automata Construction

In this section we provide a slightly variant of the automata-theoretic approach to the verification of the three-valued linear-time logic *LTL* as proposed in [7]. In particular, in what follows we generalize the construction in [7] for the truth values *tt*, *ff*, and *uu*.

**Definition 5 (Closure and Elementarity)** *The closure  $cl(\psi)$  of an LTL formula  $\psi$  is the set consisting of all subformulas  $\phi$  of  $\psi$  and their negation  $\neg\phi$ . A set  $B \subseteq cl(\psi)$  is consistent w.r.t. propositional logic iff for all  $\psi_1 \wedge \psi_2, \neg\phi \in B$ : (i)  $\psi_1 \wedge \psi_2 \in B$  iff  $\psi_1 \in B$  and  $\psi_2 \in B$ ; (ii)  $\neg(\psi_1 \wedge \psi_2) \in B$  iff  $\neg\psi_1 \in B$  or  $\neg\psi_2 \in B$ ; (iii) if  $\phi \in B$  then  $\neg\phi \notin B$ ; (iv)  $\neg\neg\phi \in B$  iff  $\phi \in B$ . Further,  $B$  is locally consistent w.r.t. the until operator iff for all  $\psi_1 U \psi_2 \in B$ : (i) if  $\psi_2 \in B$  then  $\psi_1 U \psi_2 \in B$ ; (ii) if  $\neg(\psi_1 U \psi_2) \in B$  then  $\neg\psi_2 \in B$ ; (iii) if  $\psi_1 U \psi_2 \in B$  and  $\psi_2 \notin B$  then  $\psi_1 \in B$ ; (iv) if  $\neg\psi_1, \neg\psi_2 \in B$ , then  $\neg(\psi_1 U \psi_2) \in B$ .*

Finally,  $B$  is elementary iff it is both consistent and locally consistent.

Note that, unlike the standard construction for two-valued *LTL* [2], we do not require elementary sets to be maximal (i.e., either  $\phi \in B$  or  $\neg\phi \in B$ ), but we do require extra conditions (ii) and (iv) on consistency, and (ii) and (iv) on local consistency. These extra conditions can be derived in the classic, two-valued semantics, but need to be assumed as primitive here.

Hereafter  $\text{Lit} = AP \cup \{\neg q \mid q \in AP\}$  is the set of *literals*.



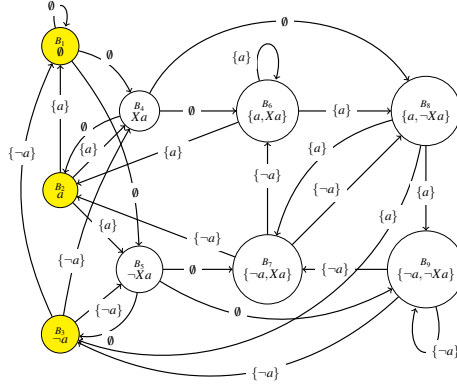


Figure 1: The automaton  $A_{\psi, \text{uu}}$  for formula  $\psi = Xa$ . Initial states are marked in yellow.

**Definition 6 (Automaton  $A_{\psi, v}$ )** Let  $\psi$  be a formula in LTL. We define the automaton  $A_{\psi, v} = \langle Q, Q_0, 2^{\text{Lit}}, \pi, \mathcal{F} \rangle$ , where  $v \in \{\text{tt}, \text{ff}, \text{uu}\}$ , as follows:  $Q$  is the set of all elementary sets  $B \subseteq \text{cl}(\psi)$ . if  $v = \text{tt}$  then  $Q_0 = \{B \in Q \mid \psi \in B\}$ ; else if  $v = \text{ff}$  then  $Q_0 = \{B \in Q \mid \neg\psi \in B\}$ ; otherwise  $Q_0 = \{B \in Q \mid \psi \notin B \text{ and } \neg\psi \notin B\}$ . The transition relation  $\pi$  is given by: let  $A \subseteq \text{Lit}$ . If  $A \neq B \cap \text{Lit}$ , then  $\pi(B, A) = \emptyset$ ; otherwise  $\pi(B, A)$  is the set of all elementary sets  $B'$  of formulas such that for every  $X\phi, \psi_1 U \psi_2 \in \text{cl}(\psi)$ : (i)  $X\phi \in B$  iff  $\phi \in B'$ ; (ii)  $\neg X\phi \in B$  iff  $\neg\phi \in B'$ ; (iii)  $\psi_1 U \psi_2 \in B$  iff  $\psi_2 \in B$  or,  $\psi_1 \in B$  and  $\psi_1 U \psi_2 \in B'$ ; (iv)  $\neg(\psi_1 U \psi_2) \in B$  iff  $\neg\psi_2 \in B$  and,  $\neg\psi_1 \in B$  or  $\neg(\psi_1 U \psi_2) \in B'$ .  $\mathcal{F} = \{F_{\psi_1 U \psi_2} \mid \psi_1 U \psi_2 \in \text{cl}(\psi)\} \cup \{Q\}$ , where  $F_{\psi_1 U \psi_2} = \{B \in Q \mid \psi_1 U \psi_2 \in B \text{ implies } \psi_2 \in B \text{ and } \neg\psi_2 \in B \text{ implies } \neg(\psi_1 U \psi_2) \in B\}$ .

According to Def. 6, the transition relation operates as follows: when the automaton reads a set  $A$  of literals that do not exist in the current state, the transition remains undefined. However, if these literals are present in the state, the automaton proceeds to verify the enabled transitions based on the semantics of the LTL operators. It is worth noting that in Def. 6, we must also specify conditions for negated formulas. This is necessary because elementary sets may not necessarily be maximal in this context.

We present an example of automaton for the next operator and truth value undefined.

**Example 1** Consider  $\psi = Xa$ . The GNBA  $A_{\psi, \text{uu}}$  in Fig. 1 is obtained as indicated in Def. 6. Namely, the state space  $Q$  consists of all elementary sets of formulas contained in  $\text{cl}(\psi) = \{a, \neg a, Xa, \neg Xa\}$ :  $B_1 = \emptyset$ ,  $B_2 = \{a\}$ ,  $B_3 = \{\neg a\}$ ,  $B_4 = \{Xa\}$ ,  $B_5 = \{\neg Xa\}$ ,  $B_6 = \{a, Xa\}$ ,  $B_7 = \{a, \neg Xa\}$ ,  $B_8 = \{\neg a, Xa\}$ ,  $B_9 = \{\neg a, \neg Xa\}$ . The initial states of  $A_{\psi, \text{uu}}$  are the elementary sets  $B \in Q$  with  $\psi, \neg\psi \notin B$ . That is,  $Q_0 = \{B_1, B_2, B_3\}$ . The transitions are depicted in Fig. 1. The set  $\mathcal{F}$  is  $\{Q\}$  as  $\psi$  does not contain until operators. Hence, every infinite run in  $A_{\psi, \text{uu}}$  is accepting.

Now, we provide a generalization of the main theoretical result proved in [7].

**Theorem 1** For every LTL formula  $\psi$  and truth value  $v \in \{\text{tt}, \text{ff}, \text{uu}\}$ , there exists a GNBA  $A_{\psi, v}$  (given as in Def. 6) s.t.  $L(A_{\psi, v}) = \text{Paths}(\psi, v)$ , where  $\text{Paths}(\psi, v)$  is the set of paths  $p \in (2^{\text{Lit}})^\omega$  such that  $(p \models^3 \psi) = v$ . Moreover, the size of  $A_{\psi, v}$  is exponential in the size of  $\psi$ .

## 4 Implementation

**Tool Architecture.** The 3vLTL tool<sup>1</sup> developed for this paper aims at generating highly reusable generalized non-deterministic Büchi automata (GNBA) [7]. Hence, instead of generating only a graphical

<sup>1</sup><https://github.com/AngeloFerrando/3vLTL>

result, 3vLTL produces a machine-readable file which can be easily parsed by third-party tools and libraries as well. From an engineering perspective, a pure graphical representation would help the final user to visualise the generated automaton, but it would not make it accessible for further evaluations.

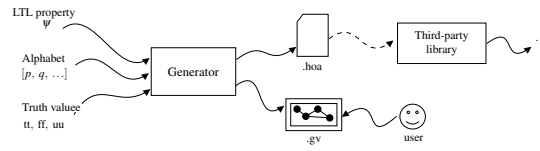


Figure 2: Overview of the tool.

Figure 2 provides an overview of 3vLTL. The tool begins by parsing the user’s input, which consists of three essential arguments. The first argument is the LTL property  $\psi$  of interest, serving as a guide for generating the GNBA. The second argument represents the alphabet of  $\psi$  and informs 3vLTL about the atomic propositions to consider when constructing the automaton. Since the automaton explicitly specifies the atomic propositions associated with its transitions, it is crucial to identify the relevant events of interest. The third argument specifies the truth value against which the LTL formulas are verified. Following the approach proposed in [7], 3vLTL offers support for generating three different GNBA versions. To elaborate further, if *tt* (representing satisfaction), *ff* (representing violation), or *uu* (representing neither satisfaction nor violation) is provided as the third argument, 3vLTL generates the respective GNBA, denoted as  $A_{\psi,tt}$ ,  $A_{\psi,ff}$ , or  $A_{\psi,uu}$ , recognizing traces that satisfy, violate, or neither satisfy nor violate  $\psi$ . 3vLTL produces two distinct output files. The first file, primarily graphical, contains the GNBA description in the DOT graph description language. The choice of DOT format stems from its widespread usage (supported by many programming languages) and its native compatibility with Graphviz<sup>2</sup>. The second file generated by 3vLTL adheres to the HOA (Hanoi Omega-Automata) format<sup>3</sup>, a machine-readable format. This format enjoys support from well-known automata-based libraries, including Spot [11] and LTL3BA [1]<sup>4</sup>. This choice enhances compatibility with third-party tools, promoting the broader utility of the GNBA generated by 3vLTL. It is important to note that while 3vLTL operates independently, it seamlessly integrates with existing automata-based tools, ensuring a smooth transition for users and enabling further advancements and applications of the GNBA it produces.

**Technical details.** We go further into the detail of the implementation. First of all, 3vLTL has been implemented in Java (version 17). The resulting runnable jar can be directly used off-the-shelf.

3vLTL is divided into three components: input handler, automaton generator, and output handler.

*Input handler.* 3vLTL handles three input data: the LTL property  $\psi$ , the alphabet, and the truth value. While the handling of the second and third arguments is straightforward, the first argument requires a bit more of work. Specifically, a parser has been implemented to parse LTL formulas using Antlr<sup>5</sup>, which is directly supported in Java. The resulting visitor for the LTL grammar is not only used to parse the LTL property  $\psi$  given in input, but it is also used to extract the corresponding closure  $cl(\psi)$ .

*Automaton generator.* After the LTL property  $\psi$  given in input has been successfully parsed, and

<sup>2</sup><https://graphviz.org/>

<sup>3</sup><http://adl.github.io/hoaf/>

<sup>4</sup>LTL3BA operates on two-valued automata, but its output is defined using three truth values, allowing it to effectively handle run-time verification scenarios. For a comprehensive examination of the distinctions between "undefined" and "unknown" truth values in the context of run-time verification, you can find more detailed information in [12].

<sup>5</sup><https://www.antlr.org>

the resulting closure  $cl(\psi)$  has been generated, the tool proceeds with the generation of the GNBA. The corresponding Java object, instantiation of the custom *Automaton* class, is generated and stored in memory. Inside such object all information about states and transitions, along with details on the initial and accepting states, are stored. In particular, the set of initial states is determined by the last input given to 3vLTL. If the user desires to produce a GNBA to recognise the traces which satisfy  $\psi$ , then the initial states in the automaton are all the states containing  $\psi$ . Note that, this is possible because the elementary subsets of  $cl(\psi)$  which determine the automaton's states are not necessarily maximal, differently from the standard automaton construction. Interestingly, the accepting states in all three cases are the same.

*Output handler.* Once the GNBA has been generated and the corresponding Java object is stored, 3vLTL moves forward to produce the resulting DOT and HOA output files. Both files are generated using two different custom methods of the *Automaton* class. Such methods pass through all states/transitions, and port these data in the wanted format.

**Experiments.** To show 3vLTL's scalability, we carried out some experiments w.r.t. the size of the LTL formula given in input. Figure 3 reports the results so obtained. As it is easy to note, the results show 3vLTL is exponential w.r.t. the size of LTL formula; where the size denotes the number of operators in the formula (e.g., the LTL formula  $XFp$  has size 2, while  $Gp \wedge Fq$  has size 3). Note that, this was expected because the transformation from LTL to GNBA is known to be exponential w.r.t. the size of the LTL formula. So, 3vLTL extends the standard algorithm, but maintains the same complexity.

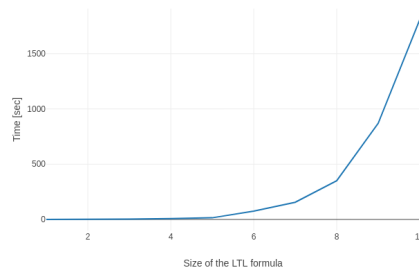


Figure 3: Experimental results.

## 5 Conclusions

In this paper, we have introduced a tool designed for generating automata from LTL formulas, interpreted within a three-valued semantics framework. To implement this tool, we have closely followed the automata construction methodology outlined in [7]. Looking ahead, our future work entails extending the capabilities of our tool to accommodate more than three truth values. This extension would enable us to create a generator capable of handling multi-valued LTL formulas. Additionally, we envision applying the automata construction and its associated implementation in various domains related to multi-valued logics. One such domain is Runtime Verification, where three-valued LTL also finds relevance. However, it is worth noting that the third value in Runtime Verification serves to maintain the impartiality of the monitor, while in our context, the third value signifies imperfect information about the system. As a result, our approach has the potential to address scenarios involving imperfect information, similar to the approach presented in [12]. Unfortunately, due to space constraints and the paper's primary focus, a comparative analysis with other tools has not been included.

## References

- [1] Tomáš Babiak, Mojmír Křetínský, Vojtěch Řehák & Jan Strejček (2012): *LTL to Büchi automata translation: Fast and more deterministic*. In: *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, Springer, pp. 95–109. Available at <https://doi.org/10.48550/arXiv.1201.0682>.
- [2] C. Baier & J. P. Katoen (2008): *Principles of Model Checking*. MIT Press.
- [3] T. Ball & O. Kupferman (2006): *An abstraction-refinement framework for multi-agent systems*. In: *LICS06*, IEEE, pp. 379–388, doi:10.1109/LICS.2006.10.
- [4] Andreas Bauer, Martin Leucker & Christian Schallhart (2006): *Monitoring of Real-Time Properties*. In S. Arun-Kumar & Naveen Garg, editors: *FSTTCS 2006: Foundations of Software Technology and Theoretical Computer Science, 26th International Conference, Kolkata, India, December 13-15, 2006, Proceedings, Lecture Notes in Computer Science 4337*, Springer, pp. 260–272. Available at [https://doi.org/10.1007/11944836\\_25](https://doi.org/10.1007/11944836_25).
- [5] Andreas Bauer, Martin Leucker & Christian Schallhart (2007): *The Good, the Bad, and the Ugly, But How Ugly Is Ugly?* In Oleg Sokolsky & Serdar Taşıran, editors: *Runtime Verification*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 126–138, doi:10.1007/978-3-540-77395-5\_11.
- [6] Francesco Belardinelli, Angelo Ferrando & Vadim Malvone (2023): *An abstraction-refinement framework for verifying strategic properties in multi-agent systems with imperfect information*. *Artif. Intell.* 316. Available at <https://doi.org/10.1016/j.artint.2022.103847>.
- [7] Francesco Belardinelli & Vadim Malvone (2020): *A Three-valued Approach to Strategic Abilities under Imperfect Information*. In Diego Calvanese, Esra Erdem & Michael Thielscher, editors: *Proceedings of the 17th International Conference on Principles of Knowledge Representation and Reasoning, KR 2020, Rhodes, Greece, September 12-18, 2020*, pp. 89–98. Available at <https://doi.org/10.24963/kr.2020/10>.
- [8] G. Bruns & P. Godefroid (1999): *Model Checking Partial State Spaces*. In: *Proceedings of the 11th International Conference on Computer Aided Verification (CAV99)*, LNCS 1633, Springer-Verlag, pp. 274–287, doi:10.1007/3-540-48683-6\_25.
- [9] G. Bruns & P. Godefroid (2003): *Model Checking with Multi-Valued Logics*. Technical Report ITD-03-44535H, Bell Labs.
- [10] Marsha Chechik, Benet Devereux & Arie Gurfinkel (2001): *Model-checking in finite state-space systems with fine-grained abstractions using SPIN*. In: *International SPIN Workshop on Model Checking of Software*, Springer, pp. 16–36, doi:10.1007/3-540-45139-0\_3.
- [11] Alexandre Duret-Lutz, Alexandre Lewkowicz, Amaury Fauchille, Thibaud Michaud, Etienne Renault & Laurent Xu (2016): *Spot 2.0 - A Framework for LTL and  $\omega$ -Automata Manipulation*. In Cyrille Artho, Axel Legay & Doron Peled, editors: *Automated Technology for Verification and Analysis - 14th International Symposium, ATVA 2016, Chiba, Japan, October 17-20, 2016, Proceedings, Lecture Notes in Computer Science 9938*, pp. 122–129. Available at [https://doi.org/10.1007/978-3-319-46520-3\\_8](https://doi.org/10.1007/978-3-319-46520-3_8).
- [12] Angelo Ferrando & Vadim Malvone (2022): *Runtime Verification with Imperfect Information Through Indistinguishability Relations*. In Bernd-Holger Schlingloff & Ming Chai, editors: *Software Engineering and Formal Methods - 20th International Conference, SEFM 2022, Berlin, Germany, September 26-30, 2022, Proceedings, Lecture Notes in Computer Science 13550*, Springer, pp. 335–351. Available at [https://doi.org/10.1007/978-3-031-17108-6\\_21](https://doi.org/10.1007/978-3-031-17108-6_21).
- [13] Paul Gastin & Denis Oddoux (2001): *Fast LTL to Büchi Automata Translation*. In Gérard Berry, Hubert Comon & Alain Finkel, editors: *Computer Aided Verification, 13th International Conference, CAV 2001, Paris, France, July 18-22, 2001, Proceedings, Lecture Notes in Computer Science 2102*, Springer, pp. 53–65. Available at [https://doi.org/10.1007/3-540-44585-4\\_6](https://doi.org/10.1007/3-540-44585-4_6).

- [14] R. Gerth, D. Peled, M. Vardi & P. Wolper (1995): *Simple On-the-fly Automatic Verification of Linear Temporal Logic*. In: *Proceedings of IFIP/WG6.1 Symposium Protocol Specification, Testing and Verification (PSTV95)*, Chapman & Hall, pp. 3–18, doi:10.1007/978-0-387-34892-6\_1.
- [15] P. Godefroid & R. Jagadeesan (2003): *On the Expressiveness of 3-Valued Models*. In: *Proceedings of the 4th International Conference on Verification, Model Checkig, and Abstract Interpretation (VMCAI03)*, LNCS 2575, Springer-Verlag, pp. 206–222, doi:10.1007/3-540-36384-X\_18.
- [16] Patrice Godefroid & Nir Piterman (2009): *LTL generalized model checking revisited*. In: *International Workshop on Verification, Model Checking, and Abstract Interpretation*, Springer, pp. 89–104. Available at <https://doi.org/10.1007/s10009-010-0169-3>.
- [17] Michael Huth, Radha Jagadeesan & David A. Schmidt (2004): *A domain equation for refinement of partial systems*. *Mathematical Structures in Computer Science* 14(4), pp. 469–505. Available at <https://doi.org/10.1017/S0960129504004268>.
- [18] Michael Huth & Shekhar Pradhan (2004): *Consistent Partial Model Checking*. *Electronic Notes in Theoretical Computer Science* 73, pp. 45–85. Available at <https://doi.org/10.1016/j.entcs.2004.08.003>.
- [19] S. C. Kleene (1952): *Introduction to Metamathematics*. North-Holland.
- [20] Beata Konikowska (1998): *A three-valued linear temporal logic for reasoning about concurrency*. ICS PAC, Warsaw, Poland, Tech. Rep.
- [21] Orna Kupferman & Yoad Lustig (2007): *Lattice Automata*. In: *Verification, Model Checking, and Abstract Interpretation, 8th International Conference, VMCAI 2007, Nice, France, January 14-16, 2007, Proceedings*, pp. 199–213. Available at [https://doi.org/10.1007/978-3-540-69738-1\\_14](https://doi.org/10.1007/978-3-540-69738-1_14).
- [22] A. Pnueli (1977): *The Temporal Logic of Programs*. In: *FOCS'77*, IEEE Computer Society, pp. 46–57, doi:10.1109/SFCS.1977.32.
- [23] S. Shoham & O. Grumberg (2004): *Monotonic Abstraction-Refinement for CTL*. In: *TACAS04*, pp. 546–560, doi:10.1007/978-3-540-24730-2\_40.
- [24] Nils Timm & Stefan Gruner (2016): *Parameterised three-valued model checking*. *Science of Computer Programming* 126, pp. 94–110, doi:10.1016/j.scico.2016.01.006.
- [25] Rachel Tzoref & Orna Grumberg (2006): *Automatic refinement and vacuity detection for symbolic trajectory evaluation*. In: *International Conference on Computer Aided Verification*, Springer, pp. 190–204, doi:10.1007/11817963\_20.
- [26] Stefan J. J. Vijzelaar & Wan J. Fokkink (2017): *Creating Büchi Automata for Multi-valued Model Checking*. In Ahmed Bouajjani & Alexandra Silva, editors: *Formal Techniques for Distributed Objects, Components, and Systems - 37th IFIP WG 6.1 International Conference, FORTE 2017, Held as Part of the 12th International Federated Conference on Distributed Computing Techniques, DisCoTec 2017, Neuchâtel, Switzerland, June 19-22, 2017, Proceedings, Lecture Notes in Computer Science 10321*, Springer, pp. 210–224. Available at [https://doi.org/10.1007/978-3-319-60225-7\\_15](https://doi.org/10.1007/978-3-319-60225-7_15).

# Towards Proved Formal Specification and Verification of STL Operators as Synchronous Observers

Céline Bellanger

ENAC, Université de Toulouse

Pierre-Loïc Garoche

ENAC, Université de Toulouse

Matthieu Martel

Université de Perpignan Via Domitia

Célia Picard

ENAC, Université de Toulouse

Signal Temporal Logic (STL) is a convenient formalism to express bounded horizon properties of autonomous critical systems. STL extends LTL to real-valued signals and associates a non-singleton bound interval to each temporal operators. In this work we provide a rigorous encoding of non-nested discrete-time STL formulas into Lustre synchronous observers.

Our encoding provides a three-valued online semantics for the observers and therefore enables both the verification of the property and the search of counter-examples. A key contribution of this work is an instrumented proof of the validity of the implementation. Each node is proved correct with respect to the original STL semantics. All the experiments are automated with the Kind2 model-checker and the Z3 SMT solver.

## 1 Introduction

In the context of autonomous critical systems, an undesirable behaviour can lead to significant material or human damage. Thus, the specification of properties and their formal verification play a paramount role in ensuring the safety, reliability and compliance of such systems.

Dynamical systems continuously respond to environmental changes. Signal Temporal Logic (STL) has emerged as a powerful formalism for expressing temporal properties within these systems [16]. The main particularity of STL language is the association of each temporal operator with a finite, non-singleton time interval, during which the operator is studied. Consider the temporal  $\diamond$  (*Eventually*) operator, which evaluates whether a property  $\varphi$  is satisfied or not at least once. A correct formalism for  $\diamond$  in STL is  $\diamond_{[a,b]}\varphi$ , where  $a$  and  $b$  are times such that  $a < b$ . Most of the time, STL properties are assess offline: we execute the system from start to finish, and we observe after the end of the execution if the system behaviour and its outputs are compliant to specified requirements.

However, the complexity of certain autonomous dynamical systems may require runtime verification. This involves continuous assessment of the system's compliance to its specification throughout execution. Synchronous observers can be employed for this purpose. These specialized observers react when a property is satisfied or violated, providing instantaneous information about the system's state. This approach offers advantages such as consistent real-time information transmission and the ability to halt executions immediately upon property satisfaction or violation, without waiting for completion. Notably, this enables quicker reactions to external events, crucial for critical systems. For instance, let us admit that we wish to satisfy a property  $\varphi$  at least once during a time interval  $[a, b]$ . If the property is satisfied in the interval, then there is no need to wait for time  $b$  to affirm that the property is indeed verified.

This paper introduces preliminary works on the specification and verification of STL operators, using synchronous observers. The rest of the document focuses on discrete times, and non nested temporal

$$\begin{aligned}
(\mathcal{X}, t) \models \mu &\Leftrightarrow \mu(t) & (1) \\
(\mathcal{X}, t) \models \neg \varphi &\Leftrightarrow \neg((\mathcal{X}, t) \models \varphi) & (2) \\
(\mathcal{X}, t) \models \varphi_1 \wedge \varphi_2 &\Leftrightarrow (\mathcal{X}, t) \models \varphi_1 \wedge (\mathcal{X}, t) \models \varphi_2 & (3) \\
(\mathcal{X}, t) \models \varphi_1 \vee \varphi_2 &\Leftrightarrow (\mathcal{X}, t) \models \varphi_1 \vee (\mathcal{X}, t) \models \varphi_2 & (4) \\
(\mathcal{X}, t) \models \varphi_1 \mathcal{U}_{[a,b]} \varphi_2 &\Leftrightarrow \exists t' \in t + [a, b] : (\mathcal{X}, t') \models \varphi_2 \wedge \forall t'' \in [t, t'] : (\mathcal{X}, t'') \models \varphi_1 & (5) \\
(\mathcal{X}, t) \models \diamond_{[a,b]} \varphi &\Leftrightarrow \exists t' \in t + [a, b] : (\mathcal{X}, t') \models \varphi & (6) \\
(\mathcal{X}, t) \models \square_{[a,b]} \varphi &\Leftrightarrow \forall t' \in t + [a, b] : (\mathcal{X}, t') \models \varphi & (7)
\end{aligned}$$

**Figure 1:** STL offline semantics

operators. For example, STL properties like  $\square_{[a,b]}(\diamond_{[c,d]}\varphi)$  with  $\varphi$  an atomic proposition and  $a, b, c$  and  $d$  distinct times such that  $a < b$  and  $c < d$ , are excluded due to the nested  $\diamond$  operator.

Our main contribution concerns the formal verification of the correctness of STL operators. To this end, we provide a three-valued online STL semantics as well as the implementation of each STL operator in the synchronous language Lustre. The soundness of the implementation is expressed as a set of lemmas and automatically proved with the Kind2 model-checker.

Section 2 covers the preliminary concepts, including the Signal Temporal Logic, the synchronous language Lustre, the model checker Kind2, and an introduction to three-valued logic. We formalise an online semantics for STL operators in Section 3, and detail its Lustre implementation in Section 4. Finally, Section 5 describes the formal correction of the operators implementation.

## 2 Preliminaries

### 2.1 Signal Temporal Logic

Let  $\mathbb{T}$  denote a set of discrete times such that  $\mathbb{T} = \mathbb{N}$  and let  $\mathcal{X}$  be a finite sets of signals. Let  $a, b \in \mathbb{T}$  with  $a < b$ . Without loss of generality, we assume that all signals are defined as functions in  $\mathbb{T} \rightarrow \mathbb{R}$  from time to real values. To simplify notations, we denote the set of time  $[t + a, t + b]$  as  $t + [a, b]$ .

**Definition 1** (STL formal grammar). *Let  $\mu$  be an atomic predicate whose value is determined by the sign of a function of an underlying signal  $x \in \mathcal{X}$ , i.e.,  $\mu(t) \equiv \mu(x(t)) > 0$ . Let  $\varphi, \psi$  be STL formulas. STL formula  $\varphi$  is defined inductively as:*

$$\varphi ::= \mu \mid \neg \varphi \mid \varphi \wedge \psi \mid \varphi \vee \psi \mid \square_{[a,b]} \psi \mid \diamond_{[a,b]} \psi \mid \varphi \mathcal{U}_{[a,b]} \psi$$

**Definition 2** (STL semantics). *The semantics of a formula  $\varphi$  is defined at a time  $t \in \mathbb{T}$  and for a set of signals  $\mathcal{X}$  as  $(\mathcal{X}, t) \models \varphi$  as described in the Figure 1.*

$\mu$  is evaluated locally, at time  $t$  over the current values of the signals, Eq. (1). Equation (2) (Negation) is the logical negation of  $\varphi$ . Equation (3) (And) is the logical conjunction between  $\varphi_1$  and  $\varphi_2$ . Equation (4) (Or) is the logical disjunction between  $\varphi_1$  and  $\varphi_2$ .

It is worth mentioning that, in STL, all temporal operators have to be associated to a bounded, non-singleton time interval. Equation (5) (Until) describes a temporal operator that is satisfied if  $\varphi_1$  holds from time  $t$  until  $\varphi_2$  becomes True within the time horizon  $t + [a, b]$ . Equation (6) (Eventually) describes a temporal operator that is satisfied if  $\varphi$  is verified at least once within the time horizon  $t + [a, b]$ . Finally,

Equation (7) (Always or Globally) describes a temporal operator that is satisfied if  $\varphi$  is always verified within the time horizon  $t + [a, b]$ . Note that the usual definitions of  $\diamond_{[a,b]}$  and  $\square_{[a,b]}$  based on  $\mathcal{U}_{[a,b]}$  still apply:

$$\diamond_{[a,b]}\varphi = \text{True } \mathcal{U}_{[a,b]} \varphi, \text{ and} \quad (8)$$

$$\square_{[a,b]}\varphi = \neg(\diamond_{[a,b]}\neg\varphi). \quad (9)$$

**Remark 1.** While evaluation of predicates is performed at time  $t$  in  $(\mathcal{X}, t) \models p \Leftrightarrow \mu(t)$ , all occurrences of time intervals  $[a, b]$  in the definitions of  $\mathcal{U}_{[a,b]}$ ,  $\square_{[a,b]}$  or  $\diamond_{[a,b]}$  are used to delay the current time  $t$ :  $t + [a, b] = [t + a, t + b]$ . These times  $a$  and  $b$  are then relative times while  $t$  acts more as an absolute time.

## 2.2 Lustre

**Lustre**[5] is a synchronous language for modeling systems of synchronous reactive components. A Lustre program  $L$  is a finite collection of nodes  $[N_0, N_1, \dots, N_m]$ . The nodes satisfy the grammar described in Table 1 in which  $td$  denotes type constructors, including enumerated types, and  $v$  either constants of enumerated types  $C$  or primitive constants such as integers  $i$ . Each node is declared by the grammar construct  $d$  of Table 1. A Lustre node  $N$

$td$	$::=$	<b>type</b> $bt$   <b>type</b> $t = \text{enum } \{ C_i, \dots \}$
$bt$	$::=$	<i>real</i>   <i>bool</i>   <i>int</i>   <i>enum\_ident</i>
$d$	$::=$	<b>node</b> $f(p)$ <b>returns</b> $(p)$ ; <b>vars</b> $p$ <b>let</b> $D$ <b>tel</b>
$p$	$::=$	$x : bt; \dots; x : bt$
$D$	$::=$	$pat = e; D$   $pat = e$ ;
$pat$	$::=$	$x$   $(pat, \dots, pat)$
$e$	$::=$	$v$   $x$   $(e, \dots, e)$   $e \rightarrow e$   $op(e, \dots, e)$   <b>if</b> $e$ <b>then</b> $e$ <b>else</b> $e$   <b>pre</b> $e$
$v$	$::=$	$C$   $i$

**Table 1:** A subset of Lustre syntax

transforms infinite streams of *input* flows to streams of *output* flows, with possible local variables denoting *internal* flows. A notion of a symbolic “abstract” universal clock is used to model system progress. At each time step  $k$ , a node reads the value of each input stream and instantaneously computes and returns the value of each output stream. Note that all the equations of a node are computed at each time step. Therefore an if-then-else statement is purely functional and both of its branches are evaluated while only one of the computed value is returned.

**Stateful constructs.** Two important Lustre operators are the unary right-shift `pre` (for `previous`) operator and the binary initialization  $\rightarrow$  (for `followed-by`) operator. Their semantics is as follows. For the operator `Pre`: at first step  $k = 0$ , `pre`  $p$  is undefined, while for each step  $k > 0$  it returns the value of  $p$  at  $k - 1$ . For the operator  $\rightarrow$ : At step  $k = 0$ ,  $p \rightarrow q$  returns the value of  $p$  at  $k = 0$ , while for  $k > 0$  it returns the value of  $q$  at  $k$  step.

For example, the Lustre equation  $y = x_0 \rightarrow \text{pre}(u)$ ; will be defined for each time step  $k$  by:

$$y(k) = \begin{cases} x_0(0) & \text{if } k = 0 \\ u(k-1) & \text{if } k > 0 \end{cases}$$

## 2.3 Specifying and verifying assume-guarantee contracts with Kind2

The annotation language CoCoSpec [6] was proposed for Lustre models to lift the notion of Hoare triple [12] and Assume/Guarantee statements as dataflow contracts. A contract is associated to a node and has only access to the input/output streams of that node. The body of a contract may contain a



```

1      node timeab (const a,b: int) returns (time: bool);
2      (*@contract
3         var clk : int = 0 -> 1 + pre clk;
4         assume a >=0;
5         guarantee time = (clk >= a and clk <= b);
6      *)

```

**Figure 2:** Example of a Lustre contract implementation

set of assume ( $A$ ) and guarantee ( $G$ ) statements and mode declarations. Modes are named and consist of require ( $R$ ) and ensure ( $E$ ) statements. Assumes, guarantees, requires, and ensures are all Boolean expressions over streams. In particular, assumptions and requires are expressions over input streams, while guarantees and ensures are expressions over input/output streams. A synchronous observer corresponds to such a contract with only a guarantee statement. A node *satisfies* a contract  $C = (A, G')$  if it satisfies  $\text{Historically}(A) \Rightarrow G'$ , where  $G' = G \cup \{R_i \Rightarrow E_i\}$  and  $\text{Historically}(A)$  when  $A$  is true at all time.

Contracts can also define local flows, acting as *ghost variables*. These potentially stateful flows can then be used in guarantees and ensure statements.

The following is an example of function `timeab` in Lustre using a local contract Figure 2. `timeab` is a Lustre node indicating whether the current time is inside a given time interval  $[a, b]$ . It takes as inputs the integers  $a$  and  $b$ , and returns a boolean value `time` that is `True` if the current time is inside  $[a, b]$ . First line of the contract (line 3) defines a local variable `clk` as an integer, which initially takes the value 0 and is then incrementing at each time. The *assume* at line 4 indicates to model checker that it has to prove the Lustre node only in the cases where the condition  $a \geq 0$  is satisfied. If another Lustre node is using `timeab`, Kind-2 also checks that this node could not provide an input  $a$  which runs counter to this assumption. Finally, Kind-2 must guarantee the equality, line 5, for all the inputs respecting the previous assumption, whatever the current time is. This equality compares the `time` output to value of the specification clock in a valid interval  $[a, b]$ . `timeab` shall calculate the same output with an internal clock bounded at time  $b$ . So here, we verify that bounding the clock has no effect on the provided output.

The Kind-2 model-checker [7] implements various SMT-based model-checking algorithms such as k-induction [20] or IC3/PDR [4] and allows to verify contracts with respect to nodes.

## 2.4 Three-valued logic

We present here the interest of three-valued logic, and introduce Kleene's three-valued logic, which we use in the next section to formalise an online version of STL operators.

For most tools, when performing monitoring of STL predicates, for a given value of simulation data, the trajectory is typically finite. It is produced by a simulation engine and stored in a data file. It is then loaded by the monitoring tool and analyzed with respect to the STL specification. In this offline setting, the final outcome indicates whether or not the input signal satisfies the specification. It is a boolean output.

Temporal operators are used to evaluate properties that change over time. Most of the time in these situations, we need to wait to decide whether a temporal property is satisfied or violated. Based on this observation, how to evaluate a property before being able to conclude, i.e. before the beginning of the time interval of a STL operator? Should we suppose that the operator is `True` or `False` before being able to decide?

A and B		B		
		F	U	T
A	F	F	F	F
	U	F	U	U
	T	F	U	T

(a) AND Operation

A or B		B		
		F	U	T
A	F	F	U	T
	U	U	U	T
	T	T	T	T

(b) OR Operation

A	$\neg A$
F	T
U	U
T	F

(c) Negation Operation

$A \Rightarrow B$		B		
		F	U	T
A	F	T	T	T
	U	U	U	T
	T	F	U	T

(d) Logical implication

**Table 2:** Truth Tables showing Kleene’s 3-valued strong logic operations

Let us consider a property  $\square_{[0,10]}P$ , a set of signals  $\mathcal{X}$  and an initial time  $t_0$ , e.g.,  $t_0 = 0$ . We are interested in checking  $(\mathcal{X}, t_0) \models \square_{[0,10]}P$ . Let us assume that we are given with a trace for  $\mathcal{X}$  of length  $l < 10$ , e.g.,  $l = 8$ , where the predicate  $P$  is valid along the whole trace. What is the validity of such a predicate? On the one hand, it is always valid, but on the other hand, it has no real definition within the time  $[l, 10]$ . Indeed,  $P$  could be false at time  $t = 9$  or  $t = 10$  and the property would be violated.

Since STL semantics requires all temporal connectors to be associated with bounded intervals, any STL predicate has a bounded horizon limit, after which it is always possible to determine the validity of a formula. We can then use this limit to evaluate temporal operators from it, considering that the values returned before may be irrelevant. But, in some case, the validity of the predicate can already be determined. In the previous example, if  $P$  is not valid at time  $t = 2$ , we already know at this time that the operator will not be satisfied at the end of the time interval. Existing works regarding online semantics for STL [8] try to optimize the runtime evaluation of the predicate monitoring, detecting when one can conclude, positively or negatively.

Rather than optimizing execution time based on binary logic, Łukasiewicz proposed a three-valued logic [15]. This logic introduces a third truth-value, Unknown (**U**), describing values for which we are not yet able to conclude if the property is satisfied or violated. Later, Kleene proposes a strong logic of indeterminacy [14] similar to Łukasiewicz logic. The main difference lies in the returned value for implication. Kleene’s approach states that  $\mathbf{U} \Rightarrow \mathbf{U}$  is Unknown while Łukasiewicz considers that  $\mathbf{U} \Rightarrow \mathbf{U}$  should be True. In our use of three-valued logic, we base ourselves on Kleene strong logic.

Table 2 presents the truth tables showing the logical operations AND, OR, the logical implication, as well as the negation for Kleene’s strong logic.

### 3 Online semantics for STL

To evaluate STL properties online, we rely on Kleene strong three-valued logic introduced in Table 2. In this section, we first introduce a way to obtain a three-valued output as proposed by Kleene, from two-valued outputs. Then, we provide an online and three-valued semantics for STL properties.

**Definition 3** (Positive, negative and indeterminacy logics). *Each STL temporal operator can be expressed in a three-valued form. To implement it, we define three new concepts: 1. A Positive logic **T** returning True when the property is satisfied, and False when it is yet undetermined or negative. 2. A Negative logic **F** that acts like an alarm to underline a negative result, which means that a statement returns True when we are sure that the property is not satisfied, and it returns False otherwise (undetermined property or satisfied property situations). 3. An Indeterminacy logic **U** highlighting situations where it is not yet possible to conclude about the satisfaction or violation of the property.*

**Definition 4.** *Let  $\varphi$ ,  $\varphi_1$  and  $\varphi_2$  be STL properties. We denote by  $\mathbf{T}_\varphi^t$  (resp.  $\mathbf{U}_\varphi^t$  and  $\mathbf{F}_\varphi^t$ ) the evaluation of  $(\mathcal{X}, t) \models \varphi$  according to the positive (resp. indeterminacy and negative) logic. We denote by  $\mathbf{B}_\varphi$  the*

evaluation of  $\varphi$  according to the offline implementation introduced in Figure 1. Note that  $\mathbf{B}_\varphi$  does not depend on time instant  $t$ .

**Property 1** (Complete and pairwise distinct). *At any time instant  $t$ , exactly one of the three logics returns True for a given property.*

$$\mathbf{T}_\varphi^t \vee \mathbf{U}_\varphi^t \vee \mathbf{F}_\varphi^t \quad (\text{completeness}) \quad (10)$$

$$\neg((\mathbf{T}_\varphi^t \wedge \mathbf{F}_\varphi^t) \vee (\mathbf{T}_\varphi^t \wedge \mathbf{U}_\varphi^t) \vee (\mathbf{U}_\varphi^t \wedge \mathbf{F}_\varphi^t)) \quad (\text{disjointness}) \quad (11)$$

**Remark 2** (Deduction of the output of the third logic). *According to Property 1, we only need the output of a given property in two of these three logics to determine its output in the last one. For example, if a property  $\varphi$  returns False in positive and negative logic, it means that  $\varphi$  is still Unknown (True in the indeterminacy logic).*

**Remark 3** (Property determination). *There exists an instant  $t_d$  from which we cannot satisfy  $\mathbf{U}_\varphi^{t \geq t_d}$ . For a non-nested temporal operator evaluated on time interval  $[a, b]$ ,  $t_d$  corresponds at the latest to  $t + b$ .*

$$\exists t_d \leq t + b : \forall t' \geq t_d, \neg \mathbf{U}_\varphi^{t'} \quad (12)$$

**Property 2.** *From a specific time instant  $t_f$ , the offline and online results are similar. Thus, the outputs of the offline  $\mathbf{B}_\varphi$  and online  $\mathbf{T}_\varphi^{t_f}$  versions are equivalent. In the same way, the negation of the offline operator is equivalent to the online negative version  $\mathbf{F}_\varphi^{t_f}$ . For a non-nested temporal operator evaluated on time interval  $[a, b]$ , this time instant corresponds at the latest to  $t + b$ :*

$$\tau \geq t + b \implies ((\mathbf{B}_\varphi \iff \mathbf{T}_\varphi^\tau) \wedge (\neg \mathbf{B}_\varphi \iff \mathbf{F}_\varphi^\tau)) \quad (13)$$

**Property 3** (Immutability: Positive and negative logics are final). *If a property is satisfied in the positive (resp. negative) logic, it will remain so in the future.*

$$\exists t \in \mathbb{T} : \mathbf{T}_\varphi^t \implies \forall t' \geq t, \mathbf{T}_\varphi^{t'} \quad (14)$$

$$\exists t \in \mathbb{T} : \mathbf{F}_\varphi^t \implies \forall t' \geq t, \mathbf{F}_\varphi^{t'} \quad (15)$$

From these three logics, we obtain easily a three-valued output. The property is: 1. True in three-valued logic if it is True in positive logic; 2. False in three-valued logic if it is True in negative logic; 3. Unknown in three-valued logic if it is True in indeterminacy logic.

Let us now characterize for each construct, the sufficient and necessary conditions to determine a positive, a negative, or a temporary indeterminate value.

In the case of a non-temporal property, the property validity can always be determined as either satisfied or violated. Let  $\mu$  be an atomic proposition, ie. non temporal, and  $t \in \mathbb{T}$ :

$$\forall \mathcal{X}, \forall t, (\mathcal{X}, t) \models \mu \iff \mathbf{T}_\mu^t \quad (\mathcal{X}, t) \models \neg \mu \iff \mathbf{F}_\mu^t \quad \mathbf{U}_\mu = \perp \quad (16)$$

If the property is a combination of multiple predicates based on logical operators ( $\wedge, \vee, \implies, \neg, \varphi, \dots$ ), the validity is obtained using Kleene's three-valued strong logic presented in Table 2.

In the case of STL temporal operators as described in Figure 1, we define a three-valued semantics describing when each operator is True, False or Unknown. For positive and negative logics, we also provide an explicit version obtained by enumerating all the terms in the time horizon  $t + a$  and  $t + b$ . The unknown explicit version for a property  $P$  can be obtained by combining the positive and negative explicit versions :

$$\mathbf{U}_P^t \text{ (explicit)} \iff (\neg \mathbf{T}_P^t \text{ (explicit)}) \wedge (\neg \mathbf{F}_P^t \text{ (explicit)}) \quad (17)$$

Let us describe the positive, negative and indeterminate versions of each STL operator:

$$\mathbf{T}_P^\tau \quad \tau \geq t+a \wedge \exists t' \in [t+a, \min(\tau, t+b)] : (\mathcal{X}, t') \models \varphi \quad (18)$$

$$\mathbf{F}_P^\tau \quad \tau \geq t+b \wedge \forall t' \in [t+a, t+b], (\mathcal{X}, t') \models \neg\varphi \quad (19)$$

$$\mathbf{U}_P^\tau \quad (\tau < t+a) \vee (\tau < t+b \wedge \forall t' \in [t+a, \tau], (\mathcal{X}, t') \models \neg\varphi) \quad (20)$$

$$\mathbf{T}_P^\tau \text{ (explicit)} \quad ((\mathcal{X}, t+a) \models \varphi) \vee ((\mathcal{X}, t+a+1) \models \varphi) \vee \dots \vee ((\mathcal{X}, t+b-1) \models \varphi) \vee ((\mathcal{X}, t+b) \models \varphi) \quad (21)$$

$$\mathbf{F}_P^\tau \text{ (explicit)} \quad ((\mathcal{X}, t+a) \models \neg\varphi) \wedge ((\mathcal{X}, t+a+1) \models \neg\varphi) \wedge \dots \wedge \\ ((\mathcal{X}, t+b-1) \models \neg\varphi) \wedge ((\mathcal{X}, t+b) \models \neg\varphi) \quad (22)$$

**Figure 3:** Three-valued semantics of Eventually operator:  $P = \diamond_{[a,b]}\varphi$

$$\mathbf{T}_P^\tau \quad \tau \geq t+b \wedge \forall t' \in [t+a, t+b], (\mathcal{X}, t') \models \varphi \quad (23)$$

$$\mathbf{F}_P^\tau \quad \tau \geq t+a \wedge \exists t' \in [t+a, \min(\tau, t+b)] : (\mathcal{X}, t') \models \neg\varphi \quad (24)$$

$$\mathbf{U}_P^\tau \quad (\tau < t+a) \vee (\tau < t+b \wedge \forall t' \in [t+a, \tau], (\mathcal{X}, t') \models \varphi) \quad (25)$$

$$\mathbf{T}_P^\tau \text{ (explicit)} \quad ((\mathcal{X}, t+a) \models \varphi) \wedge ((\mathcal{X}, t+a+1) \models \varphi) \wedge \dots \wedge ((\mathcal{X}, t+b-1) \models \varphi) \wedge ((\mathcal{X}, t+b) \models \varphi) \quad (26)$$

$$\mathbf{F}_P^\tau \text{ (explicit)} \quad ((\mathcal{X}, t+a) \models \neg\varphi) \vee ((\mathcal{X}, t+a+1) \models \neg\varphi) \vee \dots \vee \\ ((\mathcal{X}, t+b-1) \models \neg\varphi) \vee ((\mathcal{X}, t+b) \models \neg\varphi) \quad (27)$$

**Figure 4:** Three-valued semantics of Always operator:  $P = \square_{[a,b]}\varphi$

**Eventually**  $\diamond_{[a,b]}\varphi$  (Fig. 3) In the positive logic, Eq. (18), we need to wait for time  $t+a$ , the beginning of the time interval, to have a chance to conclude positively if a valid condition has been observed. From time  $t+b$ , if the condition was not yet valid, the positive eventually operator always returns `False`. For the negative logic, Eq. (19), invalidity requires to wait until the end of the time interval, otherwise one cannot conclude. Finally, the validity is unknown if we have not yet reached the end of the time interval but have not yet observed a suitable time, Eq. (20).

The explicit positive version Equation (21) is obtained by considering each instant between  $t+a$  and  $t+b$ . One of these instants is supposed to satisfy the property. We use the disjunction between all the terms to check it. At the opposite, explicit negative version Equation (22) returns `True` if all the terms between  $t+a$  and  $t+b$  satisfy  $\neg\varphi$ . We therefore rely on the conjunction between all the terms.

**Always**  $\square_{[a,b]}\varphi$  (Fig. 4) In the positive logic, Eq. (23), similarly to the negative case of the eventually operator, one needs to wait until the end of the interval to claim validity. For the negative logic, Eq. (24), we detect invalidity as soon as we observe an invalid time, within the proper time interval. Unknown cases are either before the time interval or within it, if the property  $\varphi$  is valid, up to now, Eq. (25).

The Always explicit positive version Equation (26) returns `True` if each instant between  $t+a$  and  $t+b$  satisfies  $\varphi$ . Similarly to the explicit negative version of Eventually, we use the conjunction to verify this point. For the explicit negative version to return `True`, it suffices that at one instant between  $t+a$  and  $t+b$ , the property  $\varphi$  is not satisfied. Thus, we check the disjunction of all the terms, searching if one of them violates  $\varphi$ .

$$\mathbf{T}_P^\tau \quad (\tau \geq t+a) \wedge (\exists t_1 \in [t+a, \min(\tau, t+b)] : (\mathcal{X}, t_1) \models \varphi_2 \wedge \forall t_2 \in [t, t_1], (\mathcal{X}, t_2) \models \varphi_1) \quad (28)$$

$$\begin{aligned} \mathbf{F}_P^\tau & (\exists t_6 \in [t, \min(\tau, t+a)] : (\mathcal{X}, t_6) \models \neg \varphi_1) \quad \vee \\ & (\tau \geq t+a \wedge \tau < t+b \wedge \exists t_7 \in [t+a, \tau] : (\mathcal{X}, t_7) \models \neg \varphi_1 \wedge \\ & \quad \neg(\exists t_8 \in [t+a, \tau] : (\mathcal{X}, t_8) \models \varphi_2 \wedge \forall t_9 \in [t, t_8], (\mathcal{X}, t_9) \models \varphi_1)) \quad \vee \\ & (\tau \geq t+b \wedge \neg(\exists t_{10} \in [t+a, t+b] : (\mathcal{X}, t_{10}) \models \varphi_2 \wedge \forall t_{11} \in [t, t_{10}], (\mathcal{X}, t_{11}) \models \varphi_1)) \quad (29) \end{aligned}$$

$$\begin{aligned} \mathbf{U}_P^\tau & (\tau < t+a \wedge \forall t_3 \in [t, \tau], (\mathcal{X}, t_3) \models \varphi_1) \quad \vee \\ & (\tau \geq t+a \wedge \tau < t+b \wedge \forall t_4 \in [t, \tau], (\mathcal{X}, t_4) \models \varphi_1 \quad \wedge \forall t_5 \in [t+a, \tau], (\mathcal{X}, t_5) \models \neg \varphi_2) \quad (30) \end{aligned}$$

$$\begin{aligned} \mathbf{T}_P^\tau \text{ (explicit)} & \left( \left( \bigwedge_{n=0}^a (\mathcal{X}, n) \models \varphi_1 \right) \wedge ((\mathcal{X}, t+a) \models \varphi_2) \right) \vee \left( \left( \bigwedge_{n=0}^{a+1} (\mathcal{X}, n) \models \varphi_1 \right) \wedge ((\mathcal{X}, t+a+1) \models \varphi_2) \right) \vee \dots \vee \\ & \left( \left( \bigwedge_{n=0}^{b-1} (\mathcal{X}, n) \models \varphi_1 \right) \wedge ((\mathcal{X}, t+b-1) \models \varphi_2) \right) \vee \left( \left( \bigwedge_{n=0}^b (\mathcal{X}, n) \models \varphi_1 \right) \wedge ((\mathcal{X}, t+b) \models \varphi_2) \right) \quad (31) \end{aligned}$$

$$\begin{aligned} \mathbf{F}_P^\tau \text{ (explicit)} & \left( \bigvee_{n=0}^a (\mathcal{X}, n) \models \neg \varphi_1 \right) \vee \left( \bigvee_{n_1=a+1}^b ((\mathcal{X}, n_1) \models \neg \varphi_1 \wedge \left( \bigwedge_{n_2=a}^{n_1} (\mathcal{X}, n_2-1) \models \neg \varphi_2 \right)) \right) \vee \\ & \left( \bigwedge_{n=a}^b (\mathcal{X}, n) \models \neg \varphi_2 \right) \quad (32) \end{aligned}$$

**Figure 5:** Three-valued semantics of Until operator:  $P = \varphi_1 \mathcal{U}_{[a,b]} \varphi_2$

**Until**  $\varphi_1 \mathcal{U}_{[a,b]} \varphi_2$  (Fig. 5) Until operator is the most complex. We conclude positively when an event  $\varphi_2$  occurred within the proper time interval, and until this moment  $\varphi_1$  was always satisfied, Eq. (28). For the negative logic, there are multiple conditions that can lead to a violation of the property. First before the time interval, if  $\varphi_1$  is not satisfied. Then inside the time interval, if  $(\mathcal{X}, \tau) \models \neg \varphi_1$  before the moment when  $(\mathcal{X}, \tau) \models \varphi_2$ , or if it was false before. Finally from  $t+b$ , if  $\varphi_2$  is never reached inside the time interval or if it was false before, Eq. (29). About indeterminacy, we cannot yet conclude on the validity of the formula, if, for the moment the formula is neither validated nor violated. A first condition is that  $\varphi_1$  holds from time  $t$  until now. A second is that, at the current time  $\tau$ , we have not reach yet  $t+a$  or we always have  $\neg \varphi_2$ . These condition only apply before reaching the end of the time interval  $t+b$ , Eq. (30).

As the Until operator depends at the same time to the satisfaction of a property inside the time interval, and the satisfaction of another one before and inside the time interval, the explicit versions are less trivial to obtain than for others operators. For the explicit positive version to return **True**, we need to satisfy the Until property at least once between  $t+a$  and  $t+b$ , so we proceed by disjunction. Each term of the disjunction is satisfied only if  $\varphi_1$  is satisfied from time  $t$  until this time instant included (conjunction between all the terms between instant  $t$  and this time instant) and  $\varphi_2$  is satisfied at this moment. Note that time  $t$  is represented by the 0 value in the Until temporal referential, in the same way that instants  $t+a$  or  $t+b$  correspond to time  $a$  or  $b$  inside Until. As there are several ways of violated the Until operator, explicit false version of Until is built differently as others explicit versions. Indeed, we have a disjunction between the three possibilities to not satisfy the Until operator, as described above. We verify if  $\varphi_1$  is not satisfied before or at time  $t+a$  by relying on the disjunction between all the  $\varphi_1$  terms from  $t$  until  $t+a$ . Then, inside the time interval after time  $t+a$ , the property is violated with certainty if there exist a moment where  $\varphi_1$  is not **True**, and until the previous time,  $\varphi_2$  was never satisfied. Indeed, if  $\varphi_2$

was satisfied previously, either the property is satisfied, which means that explicit negative version must return `False`; either the property was already violated before, so there exist another anterior time where  $\varphi_1$  was not satisfied before  $\varphi_2$  was satisfied. Finally, explicit negative version must return `True` if  $\varphi_2$  is never satisfy inside the time interval, which is studied by examining the conjunction of all the  $\neg\varphi_2$  terms between  $t+a$  and  $t+b$ .

## 4 Operators implementation strategy

Based on this online semantics, we propose an implementation of *Eventually*, *Always* and *Until* in discrete time. We use the synchronous language Lustre. We recall that all the nodes are available at [https://garoche.net/publication/2023\\_fmas\\_submission/](https://garoche.net/publication/2023_fmas_submission/).

**Useful constructs for the implementation** First, we define the basic nodes needed to implement the temporal operators. Node `min` returns the minimum value between two variables. Node `exist (time : bool; prop: bool)` returns `True` as soon as a property `prop` has been satisfied during the time interval represented by `time`. Node `forall_a (time: bool ; prop: bool)` returns `True` if a property `prop` has always been `True` during the time interval. All these nodes can easily be implemented in Lustre.

Regarding the implementation of the nodes detecting whether or not we are in the time slot  $t+[a,b]$ , and since we work with finite intervals, we can optimize our clock, preventing it from incrementing to infinity. We can limit the counter until value  $b$ , ensuring the absence of overflow. We implement the node `timeab`, that returns `True` if the current time instant is inside the time interval, based on this bounded internal counter. As counter stops at  $b$ , end of the time interval is intercepted looking at the counter previous value. If it was already  $b$ , we know that we exceeded the end of the time interval.

We are now able to implement our nodes for each version of each operator. In the case of the *Positive* and *Negative* versions, we want to stay as close as possible to the definition proposed in the Section 3. We take two liberties in order to optimise the memory management. First, for each operator, we define a bounded internal clock as described above. The same strategy as for the counter is used to determine the end of the time interval, comparing the previous value of the internal node counter with  $b$ . Secondly, we want to have a bounded number of memories, not dependent on the trace-length or on the length of the time interval. We proceed as described in the literature [9], by reusing the outputs obtained at the previous time instant to obtain the outputs at the current one. For example, here is the implementation of the *Until False* node: Figure 6. Others *Positive* and *Negative* versions are obtained based on the same principle. Last, we deduce the *Unknown* version from the *Positive* and *Negative* versions, as described in the Property 1.

**Remark 4.** *The case where both Positive and Negative versions of the operator are True at the same time is never supposed to happen and would result in an error. Indeed, it would mean that the property is both satisfied and violated, which is impossible. This result comes directly from Property 1.*

Note that these implementations can only represent non-nested STL operators. That is to say that we only consider  $\diamond_{[a,b]}\varphi$ ,  $\square_{[a,b]}\varphi$  and  $\varphi_1\mathcal{U}_{[a,b]}\varphi_2$  with  $\varphi$ ,  $\varphi_1$  and  $\varphi_2$  being non-temporal predicates.

## 5 Formal verification of STL operators

In this section, we demonstrate by model checking that the operators implementation described in Section 4 corresponds to the given specification, as presented in Section 3. We first introduce the formalizing

```

1      node until_false (a,b: int ; phil, phi2: bool)
2          returns (result_until_false: bool);
3      var until_time: int;
4      let
5          -- internal clock
6          until_time = min(0 -> pre until_time + 1, b);
7
8          -- init t=0 : until is violated if phil is false
9          result_until_false = not phil ->
10
11         -- violated if phil false before a
12         ((until_time <= a) and (not phil)) or
13
14         -- violated if we are in the time interval,...
15         (until_time > a and until_time <= b and
16             exist(timeab(a,b), not phil)
17             -- and before this moment we never had
18             and not (exist(timeab(a,b),
19                 -- phi2 is true and until this moment phil is true.
20                 ((phi2) and forall_a(timeab(0,b),phil)))))) or
21
22         -- violated if there is no instant in the time interval
23         ((until_time >= b) and not (exist(timeab(a,b),
24             -- where phi2 is true and until this moment
25             -- phil is true
26             ((phi2) and forall_a(timeab(0,b),phil)))))) or
27
28         -- still violated if it was violated once in the past
29         pre result_until_false;
30     tel

```

**Figure 6:** Until Lustre node for the False version of the operator

```

1      node P_at_k (const k: int; clk:int; P:bool)
2          returns (ok: bool);
3      let
4          ok = if clk = k then P else (false -> pre ok);
5      tel

```

Figure 7:  $P\_at\_k$  Lustre node

of each STL operator proof node for positive, negative and three-valued versions. Then, we present the use of Kind2 to concretely verify these proof nodes.

### 5.1 Induction on time interval size

To demonstrate the correctness of the positive and negative versions of the operators, we compare the outputs of our implementation proposition for each operator and an explicit equivalent, as provided in Equations (21), (22), (26), (27), (31) and (32). We remind that the explicit version is obtained by enumerating all the terms in the time horizon  $t + [a, b]$ . This allows to check directly the value of each term of the operator, and hence, to be sure to understand the obtained output. We have to show that our implementation and the explicit one are equivalent for any time interval. We prove this property by strong structural induction on the time interval size.

We proceed as follows. In a first time, we demonstrate that a statement is true for the smallest possible STL time interval, cf. base case of Eq. (33). Then, we demonstrate that if the statement is true for a given time interval size, it is also true when we increase the size interval by 1, cf. Eq. (35). By verifying these two properties, we demonstrate the correctness of our operators for all intervals  $[a, b]$  such that  $a, b \in \mathbb{T} \wedge a < b$ .

**Base case:**  $[a, a+1]$  Let  $\mathbf{Op}$  be a version of a temporal operator, and  $\mathbf{Op\_exp}$  its explicit representation as described in Section 3.

$$\forall a \in \mathbb{T}, \mathbf{Op}_{[a, a+1]} \varphi \iff \mathbf{Op\_exp}_{[a, a+1]} \varphi \quad (33)$$

For the base case proof, we create a new Lustre node  $P\_at\_k$ , cf. Fig. 7 that checks if a property is satisfied at a specific time or was satisfied before. This allows us to implement the explicit case. Let us take the example of the *Positive* version of the Eventually. For the  $[a, a+1]$  time interval, its explicit version is Eq. (34) and its implementation corresponds to the lines 13 and 14 of the Figure 8

$$\diamond_{[a, a+1]} \varphi \equiv ((\mathcal{X}, a) \models \varphi) \vee ((\mathcal{X}, a+1) \models \varphi) \quad (34)$$

**Inductive case:**  $[a, b+1]$  Let  $\mathbf{Op}$  be a version of a temporal operator and  $\mathbf{Op\_exp}$  its explicit representation as described in Section 3.

$$(\mathbf{Op}_{[a, b]} \varphi \iff \mathbf{Op\_exp}_{[a, b]} \varphi) \implies (\mathbf{Op}_{[a, b+1]} \varphi \iff \mathbf{Op\_exp}_{[a, b+1]} \varphi) \quad (35)$$

In our implementation,  $\mathbf{Op\_exp}_{[a, b+1]}$  is obtained thanks to the previous value of  $\mathbf{Op}_{[a, b]}$ , that we assume equivalent to  $\mathbf{Op\_exp}_{[a, b]}$ . For example, inductive case of the explicit version of Eventually



```

1      returns (base_case, ind_case: bool);
2      (*@contract
3         assume a<b and a>=0;
4         guarantee base_case;
5         guarantee ind_case;
6      *)
7      var clk : int;
8         output_ev_true, output_ev_true_bp1: bool;
9      let
10     clk = 0 -> 1 + pre clk;
11     output_ev_true = eventually_true (a,b,phi);
12     output_ev_true_bp1 = eventually_true (a,b+1,phi);
13     base_case = (b=a+1) =>
14         (output_ev_true = P_at_k(a,clk,phi) or P_at_k(a+1,clk,phi));
15
16     ind_case =
17         (output_ev_true_bp1 = (output_ev_true or P_at_k(b+1,clk,phi)));
18 tel

```

Figure 8: Eventually True proof node

```

1      node always_3v(const a,b:int; phi: bool)
2      returns (output_ev_true, output_ev_false: bool);
3      (*@contract
4         assume a<b and a>=0;
5         -- their are mutually exclusive
6         guarantee not (output_ev_true and output_ev_false);
7      *)
8      let
9         output_ev_true = eventually_true(a,b,phi);
10        output_ev_false = eventually_false(a,b,phi);
11 tel

```

Figure 9: Three-valued Eventually node in Lustre

True operator is obtained as described in Eq. (36) and its implementation corresponds to lines 16 and 17 of Figure 8

$$\diamond_{[a,b+1]}\varphi \equiv (\diamond_{[a,b]}\varphi) \vee ((\mathcal{X}, b+1) \models \varphi) \quad (36)$$

To concretely check these basic and inductive cases, we use the Kind2 model checker, cf Section 2.3. For each positive and negative version of STL operators, we express the base and inductive case as two properties, ie. two lemmas inside a contract to guarantee basic and inductive case, cf. Figure 8.

Finally, to obtain a three-valued output, we need to encode the result on two booleans. We combine the positive and negative outputs - previously verified - to determine the state of the operator. According to Property 1, Unknown is obtained if Positive and Negative outputs return `False` at the same time. As a complementary check, we ensure inside a contract that Positive and Negative versions are mutually exclusive as mentioned in the Remark 4. Figure 9 summarizes the implementation of this final node in Lustre.

Node name	# Property	Method	Proof time
timeab	assume	PDR	0.339s
	assume	induction	0.351s
	guarantee	PDR	2.618s
eventually_true	-	-	-
proof_ev_true	assume	PDR	0.653s
	assume	2-induction	0.713s
	guarantee	2-induction	26.778s
	guarantee	2-induction	29.631s
eventually_false	-	-	-
proof_ev_false	guarantee	PDR	37.215s
	guarantee	PDR	37.215s
eventually_3v	assume	PDR	0.677s
	assume	2-induction	0.688s
	guarantee	2-induction	0.688s
	guarantee	2-induction	13.216s
	guarantee	2-induction	18.855s

**Table 3:** Experiments for operator Eventually.

## 5.2 Using Kind2 as a theorem prover

Each of the three temporal operators is defined in a separate file. They all rely on basic nodes mentioned in Sect. 4, in which only `timeab` is fitted with a contract. The following tables summarize all contract elements automatically proved by Kind2 model-checker. We recall that Kind2 relies on different model-checking algorithms that are executed in parallel. The method that succeeds first interrupt the proof process. In the table, *PDR* stands for Property-Direct-Reachability [4] while *k-induction* specifies the number of steps of the k-induction process used to conclude. In both methods, Kind2 produces subproblems that are solved using Z3 [18].

As mentioned above, each operator `op` is defined using two underlying nodes `op_false` and `op_true` as well as a node `op_3v` that reconstruct the three-valued output. The nodes `op_false` and `op_true` are not directly associated to a contract but their soundness is expressed through the validity of another node: respectively `proof_op_false` and `proof_op_true`. These nodes are defining the base and inductive cases and associated to the main contract (cf. Fig. 8). Last, the final node `op_3v` is only fitted with an extra contract guarantying disjunctiveness of the output (cf. Fig. 9 encoding Eq. (11)).

Experiments were run with `kind2 v2.0.0-7-gdcc7f6f` on a 1,2 GHz Quad-Core Intel Core i7 with 16 GB of RAM. To build the table, each node is analyzed independently, but a quicker analysis of each file can be performed with all nodes analyzed at once. Note also results with the same execution time such as the elements of the node `eventually_3v`. Typically, in this case, they denote properties that were proved together k-inductive by the algorithm. We observe something similar with PDR for node `proof_ev_false`.

As a last remark, we have to say that, because of the parallel architecture of Kind2, it is difficult to obtain perfect reproductibility of the results. For example, one can observe that the runtime of the validity proof of the simple node `timeab_tmp` varies slightly between experiments while it is the exact same node. The difference can also appear in the number of unrolling of the k-induction engine.

Node name	# Property	Method	Proof time
timeab	assume	PDR	0.432s
	assume	2-induction	0.454s
	guarantee	PDR	2.855s
until_true	-	-	-
proof_until_true	assume	PDR	0.537s
	assume	2-induction	0.595s
	guarantee	2-induction	6.347s
	guarantee	2-induction	93.526s
	guarantee	2-induction	161.614s
until_false	-	-	-
proof_until_false	assume	induction	0.898s
	assume	induction	0.898s
	assume	induction	0.898s
	guarantee	2-induction	606.067s
	guarantee	PDR	1605.403s
until_3v	guarantee	PDR	34.530s

**Table 4:** Experiments for operator Until.

Node name	# Property	Method	Proof time
timeab	assume	PDR	0.400s
	assume	induction	0.425s
	guarantee	PDR	3.375s
always_true	-	-	-
proof_alw_true	guarantee	2-induction	76.633s
	guarantee	2-induction	96.686s
always_false	-	-	-
proof_alw_false	guarantee	2-induction	19.407s
	guarantee	PDR	27.540s
always_3v	guarantee	PDR	12.854s

**Table 5:** Experiments for operator Always.

## 6 Discussions and conclusion

**Related Works.** The use of three-valued logic has already been explored in the context of temporal logic, particularly in LTL, with the same division used in this paper: one value indicating the certainty of satisfaction of a property, another indicating the certainty of violation of a property, and a final value representing indeterminacy [3, 11].

Formal verification of STL properties has also been studied. Roehm et al. [19] propose to check STL properties on reach sequences, using hybrid model checking algorithms such as Cora [1] or SpaceEx [10]. A first step consists in the transformation of STL properties into their *reachset temporal logic* (RTL) equivalent. This transformation comes close to the explicit development of each operator that we described in Section 3, requiring potentially a large set of memories.

Moreover, several examples of algorithms and online implementation of STL properties have been produced, using a finite number of memories, cf [17, 9]. Thus, [9] proposes an algorithm for quantitative online STL implementation, and show on different examples the time-saving benefits of using their online method compared to the offline one. Balsini et al. [2] propose a qualitative online implementation of STL in Simulink, which nevertheless has some limitations. In particular, since three-valued logic is not used in this implementation, we cannot be sure whether a property has been satisfied or violated until the end of execution. These proposals go further than ours, allowing operators to be nested, sometime with some limitations like [2] that can only contain one operator inside another. However the soundness of the encoding is not formally proven.

**Conclusion.** In this paper, we propose an online discrete implementation of the STL semantics, in the continuity of Balsini's work [2]. Our contribution is twofold. First, we proposed an implementation based on Kleene's three-valued logic in order to be able to represent indeterminacy. Second, we formally demonstrated the soundness of our implementation, proving the validity of each operator with respect to its semantics.

Our approach was the following: we first defined the online STL semantics, and used it to build each STL operator as a synchronous observer in the Lustre language. Finally, we formally demonstrated the correctness of their implementation, using the Kind2 model-checker. We proceed by induction on the size of temporal intervals. We succeed to demonstrate all the proof objectives for each temporal operator implemented.

**Future Work.** They are mainly two directions to continue this work. A first one is to apply these operators to models and see how model-checkers such as Kind2 can verify properties or produce counter-examples. For example revisiting the use case of Roehm et al. [19]. The other direction is to extend the set of STL formulas that can be encoded in our framework. While Balsini et al. [2] proposed a similar encoding (but without proof) of nested operators with restricted form and up to two levels, we would like to lift the restrictions and deal with more general formulas. The notion of propagation delays introduced in Kempa et al. [13] could also lead to an efficient encoding with memories, also associated with proof of the implementation.

## 7 Acknowledgment

The authors would like to thank the Institute for Cybersecurity in Occitania (ICO) for partially funding this work.

## References

- [1] Matthias Althoff (2015): *An Introduction to CORA 2015*. In: *EPiC Series in Computing*, 34, EasyChair, pp. 120–151, doi:10.29007/zbkv. Available at <https://easychair.org/publications/paper/xMm>. ISSN: 2398-7340.
- [2] Alessio Balsini, Marco Di Natale, Marco Celia & Vassilios Tsachouridis (2017): *Generation of simulink monitors for control applications from formal requirements*. In: *2017 12th IEEE International Symposium on Industrial Embedded Systems (SIES)*, IEEE, Toulouse, pp. 1–9, doi:10.1109/SIES.2017.7993389. Available at <https://ieeexplore.ieee.org/document/7993389/>.
- [3] Andreas Bauer, Martin Leucker & Christian Schallhart (2006): *Monitoring of Real-Time Properties*. In S. Arun-Kumar & Naveen Garg, editors: *FSTTCS 2006: Foundations of Software Technology and Theoretical Computer Science*, Lecture Notes in Computer Science, Springer, Berlin, Heidelberg, pp. 260–272, doi:10.1007/11944836\_25.
- [4] Aaron R. Bradley (2012): *IC3 and beyond: Incremental, Inductive Verification*. In P. Madhusudan & Sanjit A. Seshia, editors: *Computer Aided Verification - 24th International Conference, CAV 2012, Berkeley, CA, USA, July 7-13, 2012 Proceedings, Lecture Notes in Computer Science 7358*, Springer, p. 4, doi:10.1007/978-3-642-31424-7\_4.
- [5] Paul Caspi, Daniel Pilaud, Nicolas Halbwachs & John Plaice (1987): *Lustre: A Declarative Language for Programming Synchronous Systems*. In: *POPL'87*, pp. 178–188, doi:10.1145/41625.41641.
- [6] Adrien Champion, Arie Gurfinkel, Temesghen Kahsai & Cesare Tinelli (2016): *CoCoSpec: A Mode-Aware Contract Language for Reactive Systems*. In: *SEFM'16*, pp. 347–366, doi:10.1007/978-3-319-41591-8\_24.
- [7] Adrien Champion, Alain Mebsout, Christoph Stickseel & Cesare Tinelli (2016): *The Kind 2 Model Checker*. In Swarat Chaudhuri & Azadeh Farzan, editors: *Computer Aided Verification - 28th International Conference, CAV 2016, Toronto, ON, Canada, July 17-23, 2016, Proceedings, Part II, Lecture Notes in Computer Science 9780*, Springer, pp. 510–517, doi:10.1007/978-3-319-41540-6\_29.
- [8] Jyotirmoy V. Deshmukh, Alexandre Donzé, Shromona Ghosh, Xiaoqing Jin, Garvit Juniwal & Sanjit A. Seshia (2017): *Robust online monitoring of signal temporal logic*. *Formal Methods in System Design* 51(1), pp. 5–30, doi:10.1007/s10703-017-0286-7.
- [9] Alexandre Donzé & Oded Maler (2010): *Robust Satisfaction of Temporal Logic over Real-Valued Signals*. In Krishnendu Chatterjee & Thomas A. Henzinger, editors: *Formal Modeling and Analysis of Timed Systems*, Lecture Notes in Computer Science, Springer, Berlin, Heidelberg, pp. 92–106, doi:10.1007/978-3-642-15297-9\_9.
- [10] Goran Frehse, Colas Le Guernic, Alexandre Donzé, Scott Cotton, Rajarshi Ray, Olivier Lebeltel, Rodolfo Ripado, Antoine Girard, Thao Dang & Oded Maler (2011): *SpaceEx: Scalable Verification of Hybrid Systems*. In Shaz Qadeer Ganesh Gopalakrishnan, editor: *Proc. 23rd International Conference on Computer Aided Verification (CAV)*, LNCS, Springer, pp. 379–395, doi:10.1007/978-3-642-22110-1\_30.
- [11] Hsi-Ming Ho, Joël Ouaknine & James Worrell (2014): *Online Monitoring of Metric Temporal Logic*. In Borzoo Bonakdarpour & Scott A. Smolka, editors: *Runtime Verification*, Lecture Notes in Computer Science, Springer International Publishing, Cham, pp. 178–192, doi:10.1007/978-3-319-11164-3\_15.
- [12] C. A. R. Hoare (1969): *An Axiomatic Basis for Computer Programming*. *Commun. ACM* 12(10), pp. 576–580, doi:10.1145/363235.363259.
- [13] Brian Kempa, Pei Zhang, Phillip H. Jones, Joseph Zambreno & Kristin Yvonne Rozier (2020): *Embedding Online Runtime Verification for Fault Disambiguation on Robonaut2*. In Nathalie Bertrand & Nils Jansen, editors: *Formal Modeling and Analysis of Timed Systems - 18th International Conference, FORMATS 2020, Vienna, Austria, September 1-3, 2020, Proceedings, Lecture Notes in Computer Science 12288*, Springer, pp. 196–214, doi:10.1007/978-3-030-57628-8\_12.
- [14] Stephen Cole Kleene (1952): *Introduction to Metamathematics*. North-Holland, Amsterdam.

- [15] J. Lukasiewicz (1970): *Selected Works*. Available at <https://www.scribd.com/document/359602256/J-Lukasiewicz-Selected-Works-L-Borkowski-Editor>.
- [16] Oded Maler & Dejan Nickovic (2004): *Monitoring Temporal Properties of Continuous Signals*. In Yassine Lakhnech & Sergio Yovine, editors: *Formal Techniques, Modelling and Analysis of Timed and Fault-Tolerant Systems*, Lecture Notes in Computer Science, Springer, Berlin, Heidelberg, pp. 152–166, doi:10.1007/978-3-540-30206-3\_12.
- [17] Oded Maler & Dejan Ničković (2013): *Monitoring properties of analog and mixed-signal circuits*. *International Journal on Software Tools for Technology Transfer* 15(3), pp. 247–268, doi:10.1007/s10009-012-0247-9.
- [18] Leonardo Mendonça de Moura & Nikolaj S. Bjørner (2008): *Z3: An Efficient SMT Solver*. In C. R. Ramkrishnan & Jakob Rehof, editors: *Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings, Lecture Notes in Computer Science* 4963, Springer, pp. 337–340, doi:10.1007/978-3-540-78800-3\_24.
- [19] Hendrik Roehm, Jens Oehlerking, Thomas Heinz & Matthias Althoff (2016): *STL Model Checking of Continuous and Hybrid Systems*. In Cyrille Artho, Axel Legay & Doron Peled, editors: *Automated Technology for Verification and Analysis*, Springer International Publishing, Cham, pp. 412–427, doi:10.1007/978-3-319-46520-3\_26.
- [20] Mary Sheeran, Satnam Singh & Gunnar Stålmarck (2000): *Checking Safety Properties Using Induction and a SAT-Solver*. In Warren A. Hunt Jr. & Steven D. Johnson, editors: *Formal Methods in Computer-Aided Design, Third International Conference, FMCAD 2000, Austin, Texas, USA, November 1-3, 2000, Proceedings, Lecture Notes in Computer Science* 1954, Springer, pp. 108–125, doi:10.1007/3-540-40922-X\_8.

# Runtime Verification of Learning Properties for Reinforcement Learning Algorithms

Tommaso Mannucci

Julio de Oliveira Filho

Intelligent Autonomous Systems

TNO – Netherlands Organisation for Applied Scientific Research  
The Hague, The Netherlands

tommaso.mannucci@tno.nl

julio.deoliveirafilho@tno.nl

Reinforcement learning (RL) algorithms interact with their environment in a trial-and-error fashion. Such interactions can be expensive, inefficient, and timely when learning on a physical system rather than in a simulation. This work develops new runtime verification techniques to predict when the learning phase has not met or will not meet qualitative and timely expectations. This paper presents three verification properties concerning the quality and timeliness of learning in RL algorithms. With each property, we propose design steps for monitoring and assessing the properties during the system’s operation.

## 1 Introduction

Reinforcement learning (RL) [16] is a bio-inspired approach to machine learning which formalizes the notion of “trial-and-error” and “learn-by-doing”. RL enables systems to learn during operation based on sequential interactions with the environment. During their learning phase, RL algorithms encourage decisions that led to good results in the past while avoiding detrimental choices. This simple concept is at the base of some stunning results in robotics automation [14], natural language processing [7], and computerised gaming, such as the Atari [10], StarCraft [19] video games, and the ancient tabletop games of Chess, Shogi, and Go [13].

Due to the runtime and interactive nature of RL algorithms, there has been an increasing demand for guarantees about their learning; e.g., that it will be concluded within a certain amount of time or interactions when done in an operational environment. It is also necessary to guarantee the agent learned its solution space well enough during the learning phase. Offering such guarantees for RL algorithms is a challenging task. Traditional testing is often not possible due to the difficulty of acquiring a representative set of operating conditions [6]. Formal testing methods do not yet scale well, and many RL algorithms use “black box” components [4], such as artificial neural networks. The underlying models for these algorithms are uninformative and highly dimensional; and the individual effect of their many parameters on the overall performance is not apparent nor easy to assess.

This work proposes new *Runtime Verification* (RV) techniques for checking properties of the learning phase of RL algorithms. RV is an engineering discipline concerned with checking a system behaviour during its execution [2]. Many RL algorithms’ properties require such a runtime verification approach. Safety, timeliness, and robustness properties, for example, can become invalid when RL algorithms engage in learning during operation. This happens because properties observed in the system during design time might no longer hold as the system changes by learning from new data. Timeliness properties, such as the duration of the learning phase, depend on the order and variety of interactions presented to the RL agent.

The specific contributions of this paper are:

M. Farrell, M. Luckcuck, M. Schwammberger, and  
M. Gleirscher (Eds): Fifth International Workshop on  
Formal Methods for Autonomous Systems (FMAS 2023)  
EPTCS 395, 2023, pp. 205–219, doi:10.4204/EPTCS.395.15

© T. Mannucci, J. de Oliveira Filho

- We propose formal specifications for three verification properties related to the learning phase of RL algorithms:

**Quality of learning** This property measures *how well has the agent learned its environment*. It is related to the variety and frequency of experiences presented to an agent during the learning phase.

**Distance to optimal policy** This property assesses *how far the current learned policy is from the optimal policy*.

**Time to learn** A property that estimates *the amount of interactions the learning process will need to evaluate a (new) policy*.

- Along with each property, we propose the design of an RV monitor able to assess the property from observations and during the learning phase. We discuss which information should be observable from the learning phase, how to collect it systematically, and how to use observations to assess the properties. We propose ideas for the monitor’s implementation and how it can be efficiently instrumented in an RL-based system.

The paper is organized as follows. Section 2 provides a short review of related and relevant work. Section 3 introduces basic concepts of RL and RV we need to derive the verification properties. Section 4 derives formal specifications for the properties, proposes monitoring techniques and examples. Section 5 discusses our conclusions and further work.

## 2 Related Work

Verification of RL safety properties has received the main priority in the literature because RL algorithms obviously need to explore in a safe way [3, 11, 22] if they are to learn in real-world setups. Pathak et al. [11] and Zhu et al. [22] go beyond system checks and specify how the result of their verification procedure can be used to enforce the safety constraints after a system re-design. Other specific RL frameworks [1, 5, 9] use monitors to guide the system preventing the agent from violating the properties specified. Safety is an important aspect and has received significant attention in research. In this work however, we focus on other two important runtime properties of RL algorithms: learning quality and timeliness.

Verification of properties for quality of learning has received less attention than safety properties in research. There are guarantees of convergence for specific algorithms, such as Q-learning [20] and SARSA [15]. But this only means that such an algorithm will eventually learn. We differ in which we provide explicit ways to assess how much a system already learned after a set of experiences. Like us, Van Wesel and Goodloe [18] propose off-line and online verification techniques for quality-of-learning properties. However, their approach does not leverage from knowledge of the inner structure of the algorithm. Xin et al. [21] introduce the concept of *exploration entropy* to guide the learning until the final policy is of sufficient quality; their approach differs from proper verification in that it steers, and thus interferes with, the learning process.

Verification of timeliness of the learning phase is even less prominent. Szepesvari [17] investigates the rate of convergence of Q-learning, and Potapov and Ali [12] analyze the influence of learning parameters on the convergence speed. But none of them provide an approach to verify if an RL algorithm will be able to learn within a desired number of interactions. This work differs from the aforementioned previous work by (1) providing formal specifications for quality and timeliness of the RL learning process and (2) providing monitoring techniques to check such properties at runtime. The monitors we propose do not modify the behaviour of the learning phase.



This work derives from the analysis of Markov decision processes in Mannor et al. [8]. We use their approach on the calculation of estimates for the RL value function and its bias and variance estimates. We extend many of their results to define formal verification properties. And we show how to monitor the RL algorithm during the learning phase to check these properties on-line.

### 3 Fundamentals

#### 3.1 Reinforcement learning

Reinforcement Learning is a class of machine learning (ML) algorithms that solves control and decision problems. This work targets RL variants which can be modelled as finite Markov decision processes (MDP) such as Q-learning [20]. We use the finite MDP problem structure to formally derive verification properties, and later, to design the verification monitor. An MDP problem is defined by a tuple  $\{S, A, T, R, \gamma\}$ , where  $S$  is a set possible of environment states.  $A$  is a set of actions an agent can take at each state and  $T := S \times A \times S \rightarrow [0, 1]$  is a probabilistic state transition function.  $R := S \times A \times S \rightarrow \mathbb{R}$  is a reward function attributing a payoff for each state transition and action.  $\gamma \in [0, 1)$  is a discount factor over past rewards.

During the execution of the RL algorithm, an agent operates in a sequence of distinct steps. During the  $n^{th}$  step, the agent observes the current state  $s_n \in S$ , chooses and performs an action  $a \in A$ . The action is chosen according to a (probabilistic) policy  $\pi := S \times A \rightarrow [0, 1]$ . This causes the environment to transition to a subsequent state  $s_{n+1} \in S$  according to  $T$ . For its action and the new state achieved, the agent receives an instantaneous reward  $r$  according to  $R$ . Rewards obtained from state  $s$  following policy  $\pi$  are accrued into the so-called value function:

$$V^\pi(s) := E\left[\sum_{n=0}^N \gamma^n R(s_n, \pi(s_n), T(s_n, a))\right] \quad (1)$$

where  $\gamma^n$  is the  $n^{th}$  power of  $\gamma$ , and  $s_n$  is the  $n^{th}$  state encountered after starting in  $s = s_0$ . A policy is optimal (indicated as  $\pi^*$ ) if it maximizes  $V^{\pi^*}(s)$  for all states. Note that the expectation operator in Eq. 1 is due to the potential stochasticity of  $\pi$  and  $T$ . In many cases, it is convenient to define the action value function

$$Q^\pi(s, a) = E[R(s, a, T(s, a)) + \gamma V^\pi(T(s, a))] \quad (2)$$

indicating the value obtainable in  $s$  by taking action  $a$  and following the policy thereafter.

Temporal difference(TD) [16] is a method to solve RL problems when functions  $T$  and  $R$  are unknown. In this method, the value function is randomly initialized and a policy  $\pi$  is followed. The reward observed at every transition is used to correct the value function, with a chosen learning rate  $\alpha$  dictating the speed of correction. TD learning is proven to converge to a fixed value function  $V^\pi$ , given the agent has had enough and representative interactions with the environment. In section 4, we will use this fact as an intuitive notion for the quality of learning.

To define RV properties and monitors, we will use two results from Mannor et al. [8]. First, estimates  $\hat{T}$  and  $\hat{R}$  (of transition and reward functions  $T$  and  $R$ , respectively) can be reconstructed from observing transitions during the learning phase. As a consequence, it is also possible to produce a value function estimate  $\hat{V}^\pi$  directly via Eq. 1. Second, estimators for the bias and variance of the this value function

estimate can be obtained as follows. Under the assumption that all state action combinations are visited at least once:

$$\text{bias}(\hat{V}^\pi) = \gamma^2 X Q V^\pi + \gamma X B + o\left(\frac{1}{\min_{s,a} N(s,a)}\right) \approx \gamma^2 X Q \hat{V}^\pi + \gamma X B; \quad (3)$$

$$\text{cov}(\hat{V}^\pi) = X W X^T + o\left(\frac{1}{\min_{s,a} N(s,a)}\right) \approx X W X^T \quad (4)$$

where  $X$ ,  $Q$ ,  $W$  and  $B$  are matrices computed from transitions and from  $\hat{T}$  and  $\hat{R}$ . The derivation and interpretation of these matrices is out of scope for this paper; the interested reader is referred to [8]. That being said, this result provides the bias and variance of the value function, which reflect the uncertainty of the agent due to lack of data. This is confirmed by the fact that  $\text{bias}(\hat{V}^\pi)$  and  $\text{cov}(\hat{V}^\pi)$  tend to zero by construction [8] if  $\min_{s,a} N(s,a) \rightarrow \infty$ . Thus estimates  $\hat{T}$ ,  $\hat{R}$  and  $\hat{V}^\pi$  will converge to their true value with more transitions.

The value function  $V^\pi$  can be used to iteratively improve the agent policy via the so-called *policy improvement*:

$$\pi_{k+1}(s) = \operatorname{argmax}_a E[R(s,a, T(s,a)) + V^{\pi_k}(T(s,a))] = \operatorname{argmax}_a Q^{\pi_k}(s,a) \quad (5)$$

which converges during a proper learning experience to the optimal policy  $\pi^*$  yielding the optimal value function  $V^{\pi^*}$ .

Replacing  $V^{\pi_k}$  with  $\hat{V}^{\pi_k}$  will yield an estimate  $\hat{V}^{\pi^*}$  of the optimal value function  $V^{\pi^*}$ . However, such an estimate will be biased, as optimization will favor actions for which the expected cumulative reward is overestimated. [8] recognizes the problem and proposes to divide the set of all transitions into a calibration set and a validation set to mitigate this inconvenience. With this, compute calibrated estimates  $\hat{T}_{\text{cal}}$ ,  $\hat{R}_{\text{cal}}$  and  $\hat{V}^{\pi_{\text{cal}}}$ , and apply policy improvement to obtain  $\pi_{\text{cal}}^*$ . From the validation set, obtain the transition and reward function estimates  $\hat{T}_{\text{val}}$  and  $\hat{R}_{\text{val}}$ . Finally, compute  $\hat{V}^{\pi_{\text{cal}}}$  via Eq. 1.

### 3.2 Runtime verification

Runtime Verification (RV) is an engineering discipline that combines (semi-)formal methods and monitoring of the system operation to check if a system's behaviour conforms to requirements. *Monitors* assess the system based on carefully collected observations of the system behaviour – called *traces*. Traces must conform to formally specified *properties*. Therefore, Bartocci et al. [2] indicate three steps to define an RV technique:

1. First, it is necessary to describe the property under verification using an unambiguous specification. Mathematical or logical formulations which can be assessed on system traces are the most common.
2. Second, it is necessary to design a monitor, which is a component able to collect and to assess traces of the system. Assessment here means any analysis steps necessary to evaluate the trace against the specified property.
3. Third, it is necessary to *instrument* the system. That is, insert observation mechanisms for correctly collecting the system traces. Good instrumentation minimally interferes with the system behaviour and performance.

This work follows these three steps for each of the proposed properties. For each property, we derive a property specification, and provide monitoring steps to observe the system and calculate the property. The monitor can be implemented to assess all the three properties concurrently and based on the same observed traces.

### 3.3 Use case: police patrol scheduling

We sketch a fictional but typical example for an RL-based learning system. In the remainder of the paper, we will use this example to illustrate the defined properties and discuss aspects relevant to the RV monitors and instrumentation.

A police department wants to use a new scheduling system for police night shifts in a city with frequent crime. This system will use an RL module that learns the most effective patrol schedules between three risk areas: the docks, the slums, and the bus station. Specifically, every hour between 00:00 am and 06:00 am, a patrol car is assigned to one of the three areas, for a total of six shifts per night. The RL algorithm for such a system has a state set  $S := \{(t, loc) | loc \in \{\text{docks, slums, station}\}, t \in [0, 5]\}$  and an action set  $A := \{\text{docks, slums, station}\}$ .

For our approach, the underlying problem structure ( $S$  and  $A$ ) must be known to the designer of the RV monitor. Neither the system nor the monitor designer knows the transition function  $T$  and the reward function  $R$  to be used. The transition function is not known because when a patrol car is sent to a location, it may take more or less time to complete the patrol. As a consequence it may miss a shift or terminate a shift early. The reward function cannot be estimated in advance as it is unknown which criminal activities can be prevented and where. However, it is decided to assign a reward between 0 and 3 in proportion to the severity of the spotted criminal activity, with 0 corresponding to no crime and 3 corresponding to a very severe crime.

## 4 Runtime Verification for quality and timeliness of RL Learning

In this section, we propose formal specifications for three verification properties related to the learning phase of RL-algorithms: *quality of learning*, *distance to optimal policy*, and *time to learn*.

### 4.1 How well has the agent learned its environment?

The first question is how to estimate if the agent has learned “enough” from its environment. Intuitively, a TD learning agent has learned enough if the current value function  $V^\pi(s)$  is close to its converged value (for all states). This choice is justified by the fact that the value function is related to both environmental stochastic functions  $T$  and  $R$ , as well as to the fact that correctly estimating  $V^\pi$  means correctly estimating the performance of the agent’s policy as well<sup>1</sup>. Unfortunately, due to the stochasticity of both policy and environment, analyzing the value function error in time can lead to premature convergence assessments. Instead, we propose using bias and covariance of the value function estimate  $\hat{V}$ . Since these reflect the lack of gathered data of the agent, they can be used to assess when enough transitions have been accumulated by the agent to learn from, even if the agent does not make direct use of the estimates  $\hat{T}$  and  $\hat{R}$ , but relies on another method to solve the MDP problem, e.g., TD learning. The procedure, based on Eq. 3 and Eq. 4, is as follows:

---

<sup>1</sup>Assuming that  $\pi$  and  $T$  act ergodically concerning  $S$  and  $A$ , i.e., that all state-action combinations are visited with non-zero probability.

1. query the policy of the system  $\pi$ ;
2. read traces, assumed in the form  $\{s, a, r, s'\}$ , i.e., the MDP transitions;
3. compute off-policy estimates  $\hat{T}$  and  $\hat{R}$  based on observed transitions, as well as on-policy estimates  $\hat{T}^\pi(s, s') := \hat{T}(s, \pi(s), s')$  and  $\hat{R}^\pi(s) := \hat{R}(s, \pi(s), \hat{T}(s, \pi(s)))$ ;
4. compute value  $\hat{V}^\pi$  via Eq. 1.
5. compute matrices  $X$ ,  $W$ ,  $B$  and  $Q$  from  $\hat{T}$  and  $\hat{R}$ ;
6. compute bias and covariance, ignoring in first approximation the  $o(\frac{1}{\min_{s,a} N(s,a)})$  term;
7. compute relative bias  $bias_{rel} := bias(\hat{V}^\pi(s))/\hat{V}^\pi(s)$  and relative variance  $\sigma_{rel} := \sigma(\hat{V}^\pi(s))/\hat{V}^\pi(s)$ ;
8. compare relative bias and variance with predefined upper thresholds; if for all states the two are below their respective thresholds, the property is satisfied.

The procedure is straightforward to follow and implement but presents a few limitations as well. First, the monitor must have the memory to store past traces to be able to recompute the estimates  $\hat{T}^\pi$  and  $\hat{R}^\pi$ . Second, the monitor must have access to the policy of the system. This is different than observing a signal trace in that the policy is not a signal, but a function utilized internally by the system. In case the policy is not observable within the system, it would be recommendable to use an estimate  $\hat{\pi}$  from the observed actions in the trace. In this case, however, some error in the estimated bias and covariance can be expected given that the policy is an estimate in itself. Third, the procedure described here is not incremental, i.e., it does not provide for a method to update the estimates of bias and covariance at time  $k+1$  given the trace at time  $k+1$  and previous estimates of bias and covariance at time  $k$ ; however, the calculations can be repeated by storing the previous traces. Finally, the method as presented requires that all state-action combinations are visited at least once (i.e.,  $\min_{s,a} N(s,a) > 0$ ). If this condition is not verified, then the bias and covariance cannot be computed.

### Example

Looking back at the example sketched in Sec.3.3, imagine the RL scheduling system has been provided with an initial exploratory policy assumed to perform decently, so as not to waste the patrolling effort while the RL system gathers information. How long should information be gathered with this policy? To answer this question, a monitor is designed, following the given procedure, to verify the property

$$\max_s bias_{rel}(s) < 0.05 \bigwedge \max_s \sigma_{rel}(s) < 0.02 \quad (6)$$

which indicates that the bias and variance are small, respectively 2% and 5% of the estimated value, and therefore the epistemic uncertainty on the value function is low (note that these thresholds are for illustration purposes). The monitor shall inform whether this property is violated, or unverified<sup>2</sup>, or satisfied at each moment. Note that the monitor will not be able to predict when the property will be satisfied, it can only say if it is so at the current time. Furthermore, the fact that the property is initially unsatisfied does not mean that it cannot be satisfied eventually.

After a sufficient amount of traces is collected, so that  $\min_{s,a} N(s,a) > 0$ , the initial estimates of bias and variance can be generated. However,  $\hat{T}$  and  $\hat{R}$  will initially be poor estimates of  $T$  and  $R$ , so the

<sup>2</sup>It might appear unsound that the monitor shall be able to report that the property is unverified since the inequality formulation of the property can in theory be always verified. However, at the start of the exploration, the condition  $\min_{s,a} N(s,a) > 0$  will not be verified for the applicability of the method. Therefore, the monitor will produce a “property unverified” response.

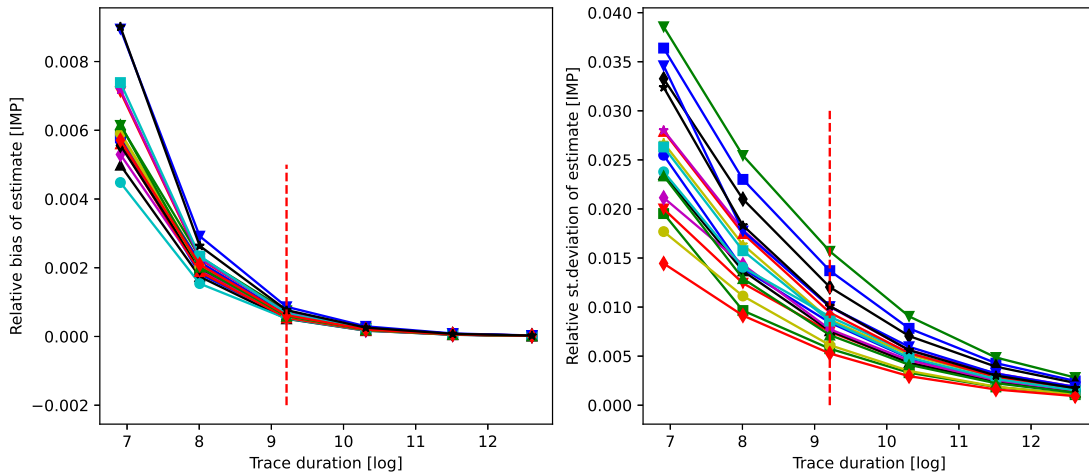


Figure 1: Evolution of relative bias and variance in the RL value function (estimates). The vertical dashed line indicates the moment of convergence after which all  $bias_{rel}$  and  $\sigma_{rel}$  are below the limits set by the property, and the property is satisfied.

corresponding bias and variance are likely to be outside of the ranges provided by Eq. 6. Therefore, the monitor will produce a “property violated” result. The longer the trace, however, the more  $\hat{T}$  and  $\hat{R}$  will resemble the true matrices. Accordingly, bias and variance will reduce, until eventually Eq. 6 will be true. The monitor will then produce a “property satisfied” response.

Figure 1 shows the relative bias and standard deviation of Eq. 6 (plotted in logarithmic scale). It can be seen how both bias and standard deviation decay to zero, in agreement with the theory. The vertical dashed line indicates the iteration at which the property is satisfied. It is possible to empirically verify the correctness of the monitor response by confronting the value function obtainable from the actual matrices  $T$  and  $R$  versus the one that can be computed from estimates  $\hat{T}$  and  $\hat{R}$  at different iterations of the monitoring. Figure 2 shows the relative error  $\frac{V(s) - \hat{V}(s)}{\hat{V}(s)}$  for all 18 states. It can be seen that this error is initially very high, indicating that  $\hat{T}$  and  $\hat{R}$  are bad estimates. However, the error reduces sensibly with the increase in iterations. At the iteration for which the property is positively validated, it can be seen that the absolute error at such iteration is lower than the estimated bias, confirming the indication of the monitor that both  $T$  and  $R$  are reasonably learned.

## 4.2 How far from the optimum is the current policy?

If the optimal value function  $V^{\pi^*}$  was known, one could compute how well  $\pi$  is faring compared to  $\pi^*$ . Unfortunately, the optimal value function is not available before learning is concluded. However, it is possible to use the estimate  $\hat{V}^{\pi^*}$  as given in Sec. 3 in first approximation.

To simplify the exposition, assume the case of positive definite reward:  $R(s, a, s') \geq 0$ . In this case, both  $\hat{V}^{\pi}$  and  $V^{\pi^*}$  are positive by construction, so that an *optimality ratio*  $\eta(s) := \frac{\hat{V}^{\pi}(s)}{V^{\pi^*}(s)}$  can be defined: if the ratio is sufficiently high for all states, this indicates that the policy is “almost optimal”.

Under the assumption that both value functions  $\hat{V}^{\pi}$  and  $\hat{V}^{\pi_{cal}} \approx \hat{V}^{\pi^*}$  are normally distributed, it is possible to bound the optimality ratio of  $\pi$ . Given that the  $i^{th}$  diagonal elements  $\sigma^2$  of the covariance

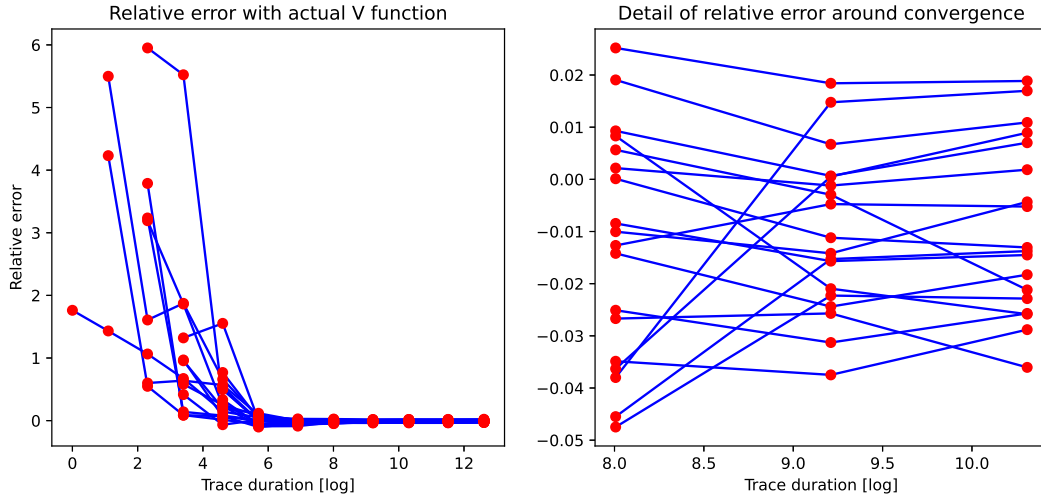


Figure 2: Relative error between  $V(s)$  and  $\hat{V}(s)$  (left); detail when the property is verified (right). This relative error is an indicator of how well the agent has learned the transition function  $T$  and the reward function  $R$ .

matrix  $cov(\hat{V}^\pi)$  coincide with the variance of the value  $\hat{V}^\pi(s_i)$  for the  $i^{th}$  state  $s_i$ , a 95% confidence interval in  $V^\pi$  can be computed.

$$\lfloor V^\pi \rfloor := \max(0, \hat{V}^\pi - bias(\hat{V}^\pi) - 2\sigma(\hat{V}^\pi)); \lceil V^\pi \rceil := \max(0, \hat{V}^\pi - bias(\hat{V}^\pi) + 2\sigma(\hat{V}^\pi)), \quad (7)$$

and similarly for  $\hat{V}^{\pi^*}$ :

$$\lfloor V^{\pi^*} \rfloor := \max(0, \hat{V}^{\pi^*} - bias(\hat{V}^{\pi^*}) - 2\sigma(\hat{V}^{\pi^*})); \lceil V^{\pi^*} \rceil := \max(0, \hat{V}^{\pi^*} - bias(\hat{V}^{\pi^*}) + 2\sigma(\hat{V}^{\pi^*})), \quad (8)$$

and thus  $\eta$  is bounded as  $\underline{\eta} \leq \eta \leq \bar{\eta}$ , with

$$\underline{\eta} = \lfloor V^\pi \rfloor / \lceil V^{\pi^*} \rceil; \bar{\eta} = \min(1, \lceil V^\pi \rceil / \lfloor V^{\pi^*} \rfloor), \quad (9)$$

again within a 95% confidence interval. Omitting the dependency from  $s$  for legibility, the procedure is as follows.

1. divide the trace into a calibration set and a validation set;
2. utilize the calibration set to obtain the transition and reward function estimates  $\hat{T}_{cal}$  and  $\hat{R}_{cal}$ ;
3. obtain the optimal policy  $\pi_{cal}^*$  via iterated policy improvement;
4. utilize the validation set to obtain the transition and reward function estimates  $\hat{T}_{val}$  and  $\hat{R}_{val}$ ;
5. compute the value function  $\hat{V}^{\pi_{cal}}$ , as well as the bias  $bias(\hat{V}^{\pi_{cal}})$  and covariance matrix  $cov(\hat{V}^{\pi_{cal}})$  substituting  $\hat{T}_{val}$  and  $\hat{R}_{val}$  for  $\hat{T}$  and  $\hat{R}$ ;
6. compute upper and lower bounds  $\lfloor V^\pi \rfloor$ ,  $\lceil V^\pi \rceil$  and  $\lfloor V^{\pi^*} \rfloor$ ,  $\lceil V^{\pi^*} \rceil$ , for all states;
7. compute upper and lower bounds on the optimality ratio  $\underline{\eta}$  and  $\bar{\eta}$ , again for all states;
8. compare  $\underline{\eta}$  and  $\bar{\eta}$  with predefined lower thresholds; if for all states the two are above their respective thresholds, the property is satisfied.

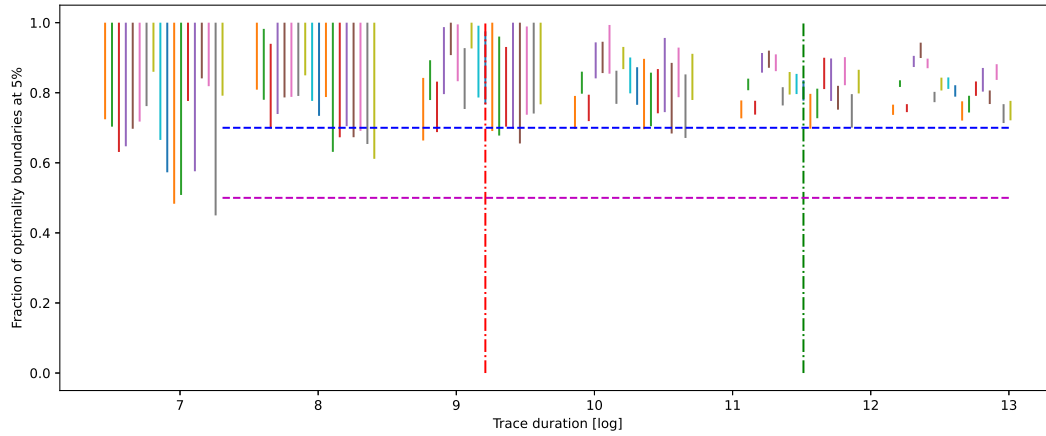


Figure 3: Optimality boundaries  $\underline{\eta}$  and  $\bar{\eta}$ , including property thresholds (horizontal dashed lines) and convergence iterations (vertical dot-dashed lines).

### Example

Consider once more the recurring use case example of police patrol. Initially, the system has been provided with a sensible exploratory policy. This is not likely to be the optimal one, but it could be close enough to optimality to not be worth changing. Conversely, it could be so suboptimal to warrant a change of policy. This loose intuition of “distance from optimum” is encoded in the property:

$$\underline{\eta}(s) \geq 50\% \wedge \bar{\eta}(s) \geq 70\% \forall s \quad (10)$$

which guarantees that  $\pi$  is not too far from the optimum (based on  $\underline{\eta}$ ) as well as indicating that  $\pi$  is potentially close to the optimum (based on  $\bar{\eta}$ ). To this must be added the condition

$$\forall s \max bias_{rel}(s) < 0.05 \wedge \max \sigma_{rel}(s) < 0.02, \quad (11)$$

on both  $V^\pi$  and  $V^{\pi^*}$ , to ensure that all estimates of  $T$  and  $R$ , which are necessary to compute the bounds on  $\eta$ , are accurate.

To verify this property (under this condition), it is necessary to first estimate the range of the optimum  $V^{\pi^*}$ . In this example, the monitor utilizes a random calibration set equal to 5% of the total traces to estimate the optimum policy. The remainder of 95% of the traces are utilized to estimate the value function  $V^{\pi^*}$  as well as the bias and variance following the procedure introduced in this section. After that, it is possible to compute the optimality ratio boundaries  $\underline{\eta}$  and  $\bar{\eta}$ .

Figure 3 shows the optimality ratio boundaries computed at different times during a sample run of the system. The optimality boundaries are indicated via vertical error bars. The dashed horizontal lines indicate the 50% and 70% optimality thresholds of in Eq. 10. The vertical lines indicate at which iteration the value function bias and standard deviation satisfy the property in Eq. 6 for the policy value function  $V^\pi$  (left) and for the optimal value function  $V^{\pi^*}$  (right) respectively. The example shows that Eq. 10 is punctually verified after Eq. 11 is satisfied.

As can be seen from Figure 3, bounds on  $\eta$  shrink with the trace duration; this is due both to the reduction in bias and variance of the estimated optimal value function, as well as due to an actual im-

provement of the calibrated policy  $\pi_{cal}$  due to additional learning. This continuous improvement also explains why error bars at a given iteration are not always contained within the error bars of the previous iterations.

### 4.3 How long will it take to evaluate a (new) policy?

We consider three scenarios when an RL system must learn or update a policy:

1. the system attempts to introduce a new policy  $\pi'$  for the same functionality;
2. the system functionality change, so that a new unknown reward function  $R'$  takes effect;
3. the system environment changes, so that a new unknown transition function  $T'$  takes effect.

Estimating how long this will take might be relevant, especially if there is a finite amount of resources to do so. This section determines an upper bound to the required number of iterations. For all cases, we assume that both policy and learning rate  $\alpha$  are stationary. We also assume that the reward function  $R$  is bounded and deterministic. To derive this property, we assume an on-policy RL algorithm is used. And our monitor design needs an estimate of the transition function  $\hat{T}$ . We discuss this assumption in details when illustrating the third case.

#### 4.3.1 New policy

We analyze first the case when the system desires to implement a different policy, but neither the functionality nor the environment is any different. Estimates  $\hat{T}$  and  $\hat{R}$  are still valid and refer to a well-learned environment. To predict the learning time, it is necessary to predict how long  $V^\pi$  will take to converge to its new value. Consider first the case in which the entire value function  $V^\pi$  or action value function  $Q^\pi$  can be re-estimated in *updates*, with each update covering the full state space or state-action space, until convergence. This can be done if both  $R$  and  $T$  are known. In this case, it can be demonstrated that for an on-policy value update, such as a SARSA [16], the expected consecutive difference  $\Delta := Q_{k+1}^\pi - Q_k^\pi$  before and after an update decays by a factor  $\hat{\gamma} = 1 - \alpha(1 - \gamma)$  at each update iteration (see Appendix for details). Furthermore, assuming  $R$  is not stochastic, the decay bound holds at every iteration. Defining convergence as  $\|\Delta\| < \varepsilon$ , it can be seen that for an on-policy update,  $Q$  will converge in  $M_u$  updates, if

$$\hat{\gamma}^{M_u} \|\bar{\Delta}\| < \varepsilon \Rightarrow M_u = \text{ceil}\left(\frac{\log(\frac{1}{\varepsilon} \|\bar{\Delta}\|)}{\log \frac{1}{\hat{\gamma}}}\right), \quad (12)$$

where  $\text{ceil}(x)$  is the ceiling function of  $x$ , and  $\bar{\Delta}$  is an upper estimate of the initial consecutive difference. This is maximum in case  $V_0^\pi(s) = Q_0^\pi(s, a) = R_{min}/(1 - \gamma)$  and  $V_1^\pi(s) = Q_1^\pi(s, a) = R_{max}/(1 - \gamma)$ ,  $\forall s, a$ . Therefore, a rigorous upper bound on convergence can be found when this is re-estimated via full state or state-action updates<sup>3</sup>.

However, such updates seldom apply in practice, with state visits determined by the stochasticity of  $\pi$  and  $T$ . In this case, it is necessary to reduce the above estimation in terms of  $M_u$  updates to a different boundary  $M_t$  in terms of *state transitions*. During updates all states or state-action combinations are visited equally often, which is not true for state transitions. However, one could make the rough assumption that if  $M_u$  updates are necessary to converge then the same convergence can be reached when

---

<sup>3</sup>Note that, even though the derivation of  $M_u$  makes use of both  $T$  and  $R$ , we do not imply that the RL agent being monitored will make direct use of these quantities.



each state is visited  $M_u$  times (whether or not this assumption is valid will be discussed in due time). Following this first assumption, it is possible to convert the previous bound into

$$M_t = M_u \max_{s,a} \tau(s,a|T, \pi, s_{\text{init}}), \quad (13)$$

where  $\tau(s,a)$  indicates the period, i.e., the number of transitions between two consecutive visits of state  $s$  in which action  $a$  is selected (also known as the *revisit time*). It is immediate to see that in the most general case, the revisit time is an unbounded stochastic variable; a first approximation of  $\tau$  can nonetheless be obtained from the steady-state state probabilities  $p$  of the stochastic transition matrix  $T$ . This can be found as the solution to the following eigenvector problem:

$$\mathbf{p} := \{p(s_0), \dots, p(s_N)\} \Rightarrow \mathbf{p} = T\mathbf{p}, \sum_{i=0}^N p(s_i) = 1, p(s_i) \geq 0. \quad (14)$$

Note that to compute such steady-state probability it is necessary to know  $T$ . That being said, Eq. 13 can be rewritten in first approximation as:

$$M_t = M_u \max_{s,a} [p(s)\pi(s,a)]^{-1}. \quad (15)$$

Note that the naive assumption of Eq. 13 that  $M_u$  visits at each state-action pair are equivalent to  $M_u$  updates does not hold in general. Indeed, updates not only guarantee an equal visit count among all states but also an equal “mixing” of visits, so that the entire value or action value function is re-estimated evenly. That being said, in the absence of pathological graphs induced by  $T$  (e.g., with absorbing states), one can utilize Eq. 15 as a first upper bound estimate.

Also note that Eq. 15 yields  $M_t \rightarrow \infty$  in case either  $p(s) \rightarrow 0$  or  $\pi(s,a) \rightarrow 0$ , i.e., if a state is not reachable or an action for a reachable state is never selected by the policy. In both cases, the limitation is more theoretical than practical. Consider the case  $p(\tilde{s}) = 0$  and assume for simplicity that the environment is not in  $\tilde{s}$  at the start of exploration: since  $\tilde{s}$  (or  $(\tilde{s}, a)$ ) cannot be reached, this never affects the value computation in Eq. 1 except for  $V(\tilde{s})$  (or  $Q(\tilde{s}, a)$ ), and  $\tilde{s}$  can be effectively ignored. The same considerations apply for  $\pi(s,a) = 0$ . By extension, a state-action couple for which  $p(s)\pi(s,a)$  is smaller than an arbitrarily small negligibility threshold  $\omega \geq 0$  will also have a negligible effect on the policy evaluation and can also be ignored. Hence, the new upper bound for convergence in terms of state transitions is

$$M_t = M_u \max_{(s,a) \in \Omega} [p(s)\pi(s,a)]^{-1}, \text{ where } \Omega := \{(s,a) | p(s)\pi(s,a) \geq \omega\}. \quad (16)$$

Finally, Eq. 16 estimates the transitions required to evaluate the new policy. Selecting a proper  $\omega > 0$  is not an immediate choice but requires some insight into the problem at hand. In the absence of this, a conservative estimate can be made by selecting  $\omega = 0$ , thus ensuring that  $M_t$  will be finite.

Summarizing, the procedure is as follows:

1. select arbitrarily small coefficients  $\varepsilon$  and  $\omega$ ;
2. compute bound  $M_u$ ;
3. solve eigenvalue problem of Eq. 14 given  $T$ ;
4. compute bound  $M_t$  from  $M_u$ ,  $\mathbf{p}$  and  $\pi$ ;
5. if  $M_t$  is lower than a predefined upper threshold, the property is verified.

### 4.3.2 New reward

The case of a new reward is identical to the case of a new policy. Indeed, Eq. 16 can be directly reutilized, with the caveat that  $M_u$  must be estimated according to the new values  $R_{min}$  and  $R_{max}$  if these differ from the previous.

### 4.3.3 New environment

When considering a new environment,  $T$  can change from what was previously estimated by the monitor. In principle, this means that the approach investigated in this section is no longer valid due to the unknown state probability  $p(s)$ . Therefore, the monitor will first need to reestimate  $\hat{T}$  (as per the first property discussed). Before that, no estimation can be provided.

That being said, a very conservative, worst-case estimation can still be given if one were to provide a positive negligibility  $\omega$  on state-action visits. In that case, a worst-case estimation can be provided in the form:

$$\overline{M}_t = M_u \frac{1}{\omega} \geq M_t. \quad (17)$$

To apply the previous estimate  $M_t$  for runtime monitoring, it is important to make a few observations. First, the convergence and negligibility constants  $\varepsilon$  and  $\omega$  must be provided to the monitor (to use in Eq. 12 and Eq. 16). Second, the monitor computes the steady-state probability  $\mathbf{p}$  of Eq. 14 by replacing  $T$  with  $\hat{T}$ . If such an estimate is missing, the monitor can either provide an immediate response using the worst-case estimate of Eq. 17 or return an unverified response until the property in Eq. 6 is verified.

### Example

Consider once more the police patrol case, and assume that informants notify police authorities of rumoured changes in the location and time of criminal activities. As a result, the initial policy  $\pi$  might no longer be satisfactory, due to the unspecified change of the reward function  $R$ , and reestimating its value will take time. Police officials want to know if the time it will take for the RL system to reestimate correctly the value of  $\pi$  is acceptable, i.e., if the following property holds:

$$M_t \leq M_t^{\max}, \quad (18)$$

where  $M_t^{\max}$  can be derived from, e.g., a time constraint. A monitor is then deployed to evaluate the property: this must have access to the estimate  $\hat{T}$  (since  $T$  is unchanged, the estimate is still valid); as well as to the potentially new quantities  $R_{min}$  and  $R_{max}$ .

Thus, the monitor can employ Eq. 16 to estimate  $M_t$  and therefore to verify the property of Eq. 18 before each policy improvement step. As an example, consider the case in which  $\pi$  dictates a fixed patrol schedule in the form

$$\pi(t, loc) = \begin{cases} \text{slums} & \text{if } t \leq 1 \\ \text{station} & \text{if } t = 2 \\ \text{docks} & \text{else.} \end{cases} \quad (19)$$

The monitor can now estimate the number of transitions necessary to evaluate this policy. Given learning parameters  $\alpha = 0.75$ ,  $\gamma = 0.5$ , constants  $\varepsilon = 0.05$  and  $\omega = 0$ , and assuming  $R_{min}$  and  $R_{max}$  to be 0 and 3 once more, the monitor estimates 62 transitions necessary for convergence.

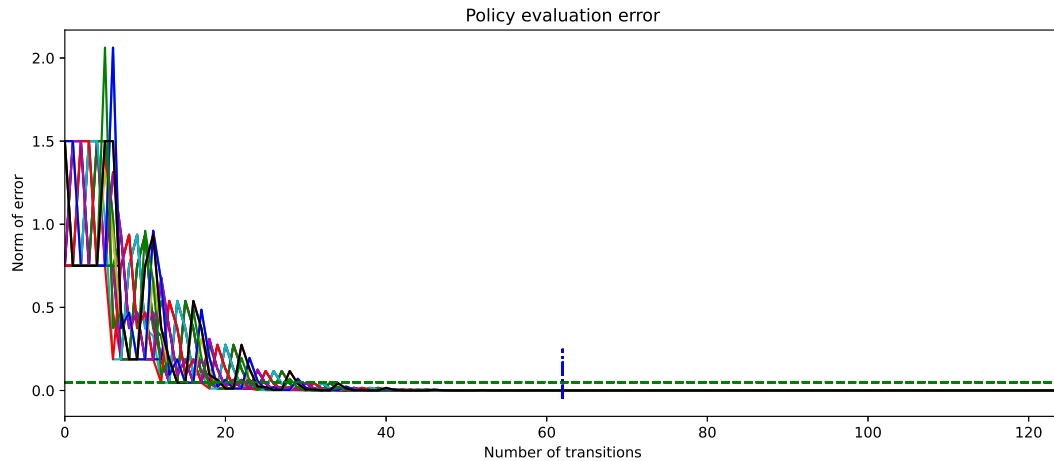


Figure 4: Norm of value function error  $\|\Delta\|$  for all initial states. The horizontal dotted line indicates the convergence threshold  $\varepsilon$ , while the vertical dot-dashed line indicates the expected convergence time.

Figure 4 shows the actual amount of transitions necessary to estimate the policy value. The continuous lines indicate the change in value function error  $\|\Delta\|$  for each of the 18 initial states. It can be seen that  $\|\Delta\|$  reduces as the number of transitions increases, albeit not monotonically. It can also be seen that approximately at 30 transitions, the error for all initial states is below the threshold of  $\varepsilon = 0.05$ , indicated by the dashed horizontal line. Thus, the value  $M_t$  for this example was indeed a conservative estimate.

## 5 Conclusions

In this paper, we propose three verification properties for the learning phase of reinforcement learning (RL) agents. The first property relates to the quality of the learning phase. It expresses whether or not the agent has learned from sufficiently enough and sufficiently varied experiences, to have a suitable representation of its environment. We devise a runtime verification monitoring technique to assess this property by estimating the variance and bias of the learned value function. This property enables the verification of the second and third properties.

The second property measures the actual learned policy relative to the ideal optimum. We extend our monitoring techniques to derive an optimality ratio  $\eta$ . The verification monitor uses confidence intervals of  $\eta$  to indicate upper and lower bounds on optimality. Our RV monitor can check whether the current policy guarantees minimal satisfactory behaviour and whether the policy could potentially be close to optimal.

The third property checks if the learning time falls within a desired number of interactions. It can be used when a new policy is put in place, or when a known policy is applied to a new transition function  $T$  or a new reward function  $R$ . Verification of such property is relevant to systems that require estimating the value induced by a new policy or new environment within a limited time. In this case, the evaluation time is a rough estimate, due to the assumption of Eq. 13 on the equivalence between updates and state transitions. In future work, we will improve this estimate by examining the mixing properties of the graph induced by  $\pi$  on the transition matrix  $T$ . For this property, we discuss the monitoring techniques necessary to assess the property using the system's runtime observations.

Our future work includes incorporating these techniques into a runtime verification tool for autonomous robots. We will also expand the formulation of our properties and monitoring techniques towards Deep Reinforcement Learning algorithms (DRL). DRL algorithms typically use continuous state- and action spaces, which our approach cannot yet cover.

## References

- [1] Greg Anderson, Abhinav Verma, Isil Dillig & Swarat Chaudhuri (2020): *Neurosymbolic Reinforcement Learning with Formally Verified Exploration*. In: *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS'20*, Curran Associates Inc., Red Hook, NY, USA, doi:10.48550/arXiv.2009.12612.
- [2] Ezio Bartocci, Yliès Falcone, Adrian Francalanza & Giles Reger (2018): *Introduction to runtime verification*. In: *Lectures on Runtime Verification*, Springer, pp. 1–33, doi:10.1007/978-3-319-75632-5.
- [3] Davide Corsi, Enrico Marchesini & Alessandro Farinelli (2021): *Formal verification of neural networks for safety-critical tasks in deep reinforcement learning*. In: *Uncertainty in Artificial Intelligence*, PMLR, pp. 333–343, doi:10.48448/tj1d-sk77.
- [4] Rudiger Ehlers (2017): *Formal verification of piece-wise linear feed-forward neural networks*. In: *International Symposium on Automated Technology for Verification and Analysis*, Springer, pp. 269–286, doi:10.1007/978-3-319-68167-2\_19.
- [5] Nathan Hunt, Nathan Fulton, Sara Magliacane, Trong Nghia Hoang, Subhro Das & Armando Solar-Lezama (2021): *Verifiably safe exploration for end-to-end reinforcement learning*. In: *HSCC '21: 24th ACM International Conference on Hybrid Systems: Computation and Control, Nashville, Tennessee, May 19-21, 2021*, ACM, pp. 14:1–14:11, doi:10.1145/3447928.3456653.
- [6] Zachary Kenton, Angelos Filos, Owain Evans & Yarin Gal (2019): *Generalizing from a few environments in safety-critical reinforcement learning*. doi:10.48550/arXiv.1907.01475.
- [7] Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley & Jianfeng Gao (2016): *Deep Reinforcement Learning for Dialogue Generation*. In: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Austin, Texas, pp. 1192–1202, doi:10.18653/v1/D16-1127. Available at <https://aclanthology.org/D16-1127>.
- [8] Shie Mannor, Duncan Simester, Peng Sun & John N Tsitsiklis (2007): *Bias and variance approximation in value function estimates*. *Management Science* 53(2), pp. 308–322, doi:10.1287/mnsc.1060.0614.
- [9] George Mason, Radu Calinescu, Daniel Kudenko & Alec Banks (2017): *Assured Reinforcement Learning with Formally Verified Abstract Policies*. In: *Proceedings of the 9th International Conference on Agents and Artificial Intelligence, ICAART 2017, Volume 2, Porto, Portugal, February 24-26, 2017*, SciTePress, pp. 105–117, doi:10.5220/0006156001050117.
- [10] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra & Martin Riedmiller (2013): *Playing atari with deep reinforcement learning*. doi:10.48550/arXiv.1312.5602.
- [11] Shashank Pathak, Luca Pulina & Armando Tacchella (2018): *Verification and repair of control policies for safe reinforcement learning*. *Appl. Intell.* 48(4), pp. 886–908, doi:10.1007/s10489-017-0999-8.
- [12] Alex Potapov & MK Ali (2003): *Convergence of reinforcement learning algorithms and acceleration of learning*. *Physical Review E* 67(2), p. 026706, doi:10.1103/PhysRevE.67.026706.
- [13] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy P. Lillicrap, Karen Simonyan & Demis Hassabis (2017): *Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm*. *CoRR* abs/1712.01815, doi:10.48550/arXiv.1712.01815.
- [14] Bharat Singh, Rajesh Kumar & Vinay Pratap Singh (2021): *Reinforcement learning in robotic applications: a comprehensive survey*. *Artificial Intelligence Review* 55(2), pp. 945–990, doi:10.1007/s10462-021-09997-9.

- [15] Satinder Singh, Tommi S. Jaakkola, Michael L. Littman & Csaba Szepesvári (2000): *Convergence Results for Single-Step On-Policy Reinforcement-Learning Algorithms*. *Mach. Learn.* 38(3), pp. 287–308, doi:10.1023/A:1007678930559.
- [16] Richard S Sutton & Andrew G Barto (2018): *Reinforcement learning: An introduction*. MIT press, doi:10.1109/TNN.1998.712192.
- [17] Csaba Szepesvári (1997): *The Asymptotic Convergence-Rate of Q-learning*. In: *Advances in Neural Information Processing Systems 10, [NIPS Conference]*, The MIT Press, pp. 1064–1070, doi:10.5555/3008904.3009053.
- [18] Perry Van Wesel & Alwyn E Goodloe (2017): *Challenges in the verification of reinforcement learning algorithms*. Technical Report. Available at <https://ntrs.nasa.gov/citations/20170007190>.
- [19] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev et al. (2019): *Grandmaster level in StarCraft II using multi-agent reinforcement learning*. *Nature* 575(7782), pp. 350–354, doi:10.1038/s41586-019-1724-z.
- [20] Christopher J.C.H. Watkins & Peter Dayan (1992): *Machine Learning* 8(3/4), pp. 279–292, doi:10.1023/a:1022676722315.
- [21] Bo Xin, Haixu Yu, You Qin, Qing Tang & Zhangqing Zhu (2020): *Exploration entropy for reinforcement learning*. *Mathematical Problems in Engineering* 2020, doi:10.1155/2020/2672537.
- [22] He Zhu, Zikang Xiong, Stephen Magill & Suresh Jagannathan (2019): *An inductive synthesis framework for verifiable reinforcement learning*. In: *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019, Phoenix, AZ, USA, June 22-26, 2019*, ACM, pp. 686–701, doi:10.1145/3314221.3314638.

## Appendix

First, observe that the action value update can be rewritten as

$$Q_{k+1}^\pi(s, a) = (1 - \alpha)Q_k^\pi(s, a) + \alpha \sum_{s' \in S} T(s, a, s') [r(s, a, s') + \gamma \sum_{a'} \pi(s', a') Q_k^\pi(s', a)].$$

Define now the matrix  $\Pi \in \mathbb{R}^{|S| \times |A| \times |S|}$ , whose entries are defined as:

$$\Pi_{i,j} = \begin{cases} \pi(s_j, a_{i-(j-1)|A|}) & \text{if } (j-1)|S||A| < i \leq j|S||A| \\ 0 & \text{otherwise.} \end{cases}$$

$\Pi$  is the matrix obtained by staking  $|S|$  copies of the vector-wise policy  $\pi \in \mathbb{R}^{|S| \times |A|}$ , and setting to zero each element of the resulting  $j^{\text{th}}$  column which does not contain the probabilities for the  $j^{\text{th}}$  state. Denote as  $\mathbf{Q} \in \mathbb{R}^{|S| \times |A|}$  the column vector of values  $Q^\pi(s, a)$ , as  $\mathbf{R} \in \mathbb{R}^{|S| \times |A|}$  the column vector of the expected rewards  $E_{s' \in S}[r(s, a, s')]$ , and as  $\mathbf{T} \in \mathbb{R}^{|S| \times |A| \times |S|}$  the matrix corresponding to the transition function  $T$ . Then the value update can be written in matrix form as

$$\mathbf{Q}_{k+1} = \alpha \mathbf{R} + [\alpha \gamma \mathbf{T} \Pi^T + (1 - \alpha) \mathbf{I}] \mathbf{Q}_k = \alpha \mathbf{R} + \mathbf{B} \mathbf{Q}_k,$$

where  $\mathbf{I}$  is the identity matrix. Note now that the square matrix  $\mathbf{B} := [\alpha \gamma \mathbf{T} \Pi^T + (1 - \alpha) \mathbf{I}]$  has norm  $\|\mathbf{B}\| \leq \alpha \gamma + 1 - \alpha = 1 - \alpha(1 - \gamma)$  by construction, because  $\alpha, \gamma \leq 1$ , and  $\mathbf{T}, \Pi$  are stochastic matrices.

If we define the consecutive error vector as  $\Delta_k := \mathbf{Q}_{k+1} - \mathbf{Q}_k$ . Then it follows that  $\Delta_{k+1} = \mathbf{B} \Delta_k$  indicates a contraction since  $1 - \alpha(1 - \gamma) \leq 1$ .