# Formal Methods for Autonomous Systems

**Suggested Citation:** Tichakorn Wongpiromsarn, Mahsa Ghasemi, Murat Cubuktepe, Georgios Bakirtzis, Steven Carr, Mustafa O. Karabag, Cyrus Neary, Parham Gohari and Ufuk Topcu (2022), "Formal Methods for Autonomous Systems", : Vol. xx, No. xx, pp 1–18. DOI: 10.1561/XXXXXXXXXX.

# **Tichakorn Wongpiromsarn**

Iowa State University nok@iastate.edu

Mahsa Ghasemi

Purdue University mahsa@purdue.edu

Murat Cubuktepe Georgios Bakirtzis Steven Carr Mustafa O. Karabag Cyrus Neary Parham Gohari Ufuk Topcu

The University of Texas at Austin utopcu@utexas.edu

This article may be used only for the purpose of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval.



# Contents

1	Intr	oduction	3		
	1.1	Autonomous Systems	5		
	1.2	Overview of Verification and Synthesis	8		
	1.3	What Makes It Hard?	10		
	1.4	Organizational Outline	11		
	1.5	A Note on the Coverage of the Article	11		
2	Models and Specifications				
	2.1	Transition Systems	15		
	2.2	Formal Specification Languages	20		
3	Verification and Model Checking				
	3.1	Model Checking	29		
	3.2	Computational Complexity	32		
	3.3	Tools for Model Checking	33		
	3.4	Model Checking for Autonomous Vehicles	34		
4	Closed-System Synthesis				
	4.1	Control Protocol Synthesis	39		
	4.2	Model Checking-Based Synthesis			
	4.3	A Case Study	44		

	4.4	Minimum-Violation Synthesis	46
5	Read	ctive Synthesis	59
	5.1	Synthesis of Reactive Control Protocol	60
	5.2	Receding Horizon Temporal Logic Planning	67
	5.3	Reactive Synthesis with Maximum Realizability	74
6	Prob	pabilistic Synthesis and Verification	89
	6.1	Probabilistic Model Checking in Markov Decision Processes	92
	6.2	A Case Study	98
	<ul><li>6.3</li><li>6.4</li></ul>	Parameter Synthesis for Markov Decision Processes Scenario-Based Verification in Uncertain Markov Decision	105
		Processes	116
7	Deal	ling with Information Limitations	129
	7.1	Partially Observable Markov Decision Processes	130
	7.2	Uncertain Partially Observable Markov Decision Processes	133
	7.3	Task-Oriented Active Perception and Planning	142
8	Runt	time Assurance via Shielding	157
	8.1	Shielding a UAV mission	160
	8.2	Synthesizing Admissible Shields	165
9	Veri	fiable Learning-Based Synthesis	172
	9.1	Verifiable Recurrent Neural Network-Based Policies for Par-	
		tially Observable Markov Decision Processes	174
	9.2	Case Study: Autonomous Driving in Traffic Lights	187
10	Som	e Future Directions	192
	••••	e ratare Birections	
		Formal Methods for Reinforcement Learning	192
	10.1	Formal Methods for Reinforcement Learning Operating Under Limited Information and Concealing Infor-	192
	10.1 10.2	Formal Methods for Reinforcement Learning	
	10.1 10.2	Formal Methods for Reinforcement Learning	192 195
	10.1 10.2 10.3	Formal Methods for Reinforcement Learning	192 195 196
	<ul><li>10.1</li><li>10.2</li><li>10.3</li><li>10.4</li></ul>	Formal Methods for Reinforcement Learning Operating Under Limited Information and Concealing Information	192 195 196 197
	10.1 10.2 10.3 10.4 10.5	Formal Methods for Reinforcement Learning	192 195 196

Acknowledgements	202
References	203

# Formal Methods for Autonomous Systems

Tichakorn Wongpiromsarn<sup>1</sup>, Mahsa Ghasemi<sup>2</sup>, Murat Cubuktepe<sup>3</sup>, Georgios Bakirtzis<sup>3</sup>, Steven Carr<sup>3</sup>, Mustafa O. Karabag<sup>3</sup>, Cyrus Neary<sup>3</sup>, Parham Gohari<sup>3</sup> and Ufuk Topcu<sup>3</sup>

#### ABSTRACT

Formal methods refer to rigorous, mathematical approaches to system development and have played a key role in establishing the correctness of safety-critical systems. The main building blocks of formal methods are models and specifications, which are analogous to behaviors and requirements in system design and give us the means to verify and synthesize system behaviors with *formal guarantees*.

This monograph provides a survey of the current state of the art on applications of formal methods in the autonomous systems domain. We consider correct-by-construction synthesis under various formulations, including closed systems, reactive, and probabilistic settings. Beyond synthesizing systems in known environments, we address the concept of uncertainty and bound the behavior of systems that employ learning using formal methods. Further, we examine the synthesis of systems with monitoring, a mitigation technique

Tichakorn Wongpiromsarn, Mahsa Ghasemi, Murat Cubuktepe, Georgios Bakirtzis, Steven Carr, Mustafa O. Karabag, Cyrus Neary, Parham Gohari and Ufuk Topcu (2022), "Formal Methods for Autonomous Systems", : Vol. xx, No. xx, pp 1–18. DOI: 10.1561/XXXXXXXXX.

<sup>&</sup>lt;sup>1</sup> Iowa State University; nok@iastate.edu

 $<sup>^2</sup> Purdue\ University;\ mahs a@purdue.edu$ 

<sup>&</sup>lt;sup>3</sup> The University of Texas at Austin; utopcu@utexas.edu

for ensuring that once a system deviates from expected behavior, it knows a way of returning to normalcy. We also show how to overcome some limitations of formal methods themselves with learning. We conclude with future directions for formal methods in reinforcement learning, uncertainty, privacy, explainability of formal methods, and regulation and certification.

# 1

#### Introduction

This monograph is about a class of formal methods that verify systems properties we care about, such as "bad things never happen" (safety) and "good things will eventually happen" (liveness), for autonomous system analysis and synthesis. Unlike traditionally engineered systems, autonomous systems need to readily react to changing environments and operational situations, often by coordinating and adapting. This ability of autonomous systems makes them both valuable and challenging to certify simultaneously. For example, compared to traditional industrial control systems where the operational contexts are static, autonomous systems must behave acceptably in various environments, often partially unknown (Bakirtzis et al., 2022b).

Finding design bugs that occur at the interaction of subsystems is particularly challenging (Bakirtzis et al., 2022a). In traditionally engineered systems testing and simulation are often forms of sufficient certification. However, these techniques are inadequate in ensuring the absence of design bugs in autonomous systems. Formal methods applied to autonomous systems bound the uncertainty arising from the unknown physical environments they deploy in and, therefore, assure that they

4 Introduction

will not misbehave.

Standard systems engineering goes through multiple steps of iteration to produce a system. Late in the lifecycle, system designs go through verification & validation. It is at verification & validation where formal methods provide proof of correct behavior for autonomous systems. The most significant barrier to certifying systems is ensuring that the assumptions made at the previous stages of systems engineering agree with the formal model. In subsequent sections, we introduce some of these models in increasing levels of expressivity. From simple to complex, system designers use these different models based on the system's expected behavior and deployment context. Therefore, applying formal methods to autonomous system design is not merely picking the most expressive model but, rather, the most appropriate one for a given application.

The impact of formal methods can be twofold. First, they provide provable guarantees that the system satisfies the desired properties (Baier and Katoen, 2008; Holzmann, 2004; Holzmann, 1994). This is the usual way we use formal methods, and it requires that both the system and the properties we are trying to prove are well-defined within a formalism. In addition, provable guarantees relate to the model and not the actual system. In the past, the congruence between the actual system and its model was done manually in an ad-hoc manner. Today, synthesis methods exist that automatically output the assured behavior of the system under examination from a formal model with provable guarantees. Second, we use formal methods to interrogate our assumptions for the system design (Lamport, 2002; Newcombe et al., 2015). Often, it is more useful to think about a design problem formally and use formal models as decision guides for the rest of the development of the system. This way of using formal methods is done from the early lifecycle all the way up to deployment. The expectation is not that the system contains provable guarantees via formal methods but that the system developers are informed about potential issues they need to address at the design phase, e.g., clashing requirements and beyond, modeling the open world environments we deploy autonomous systems.

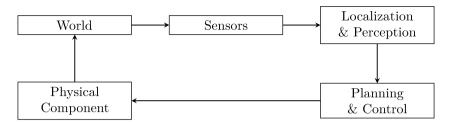


Figure 1.1: A typical architecture of autonomous vehicles.

#### 1.1 Autonomous Systems

While the public<sup>1</sup> is most familiar with misbehavior in the context of autonomous vehicles, similar inconsistencies arise in a range of applications of autonomous systems.

#### Historical roadmap

The modern development roadmap of autonomous systems can be illustrated by that of modern autonomous vehicles revealing that there are several design challenges to overcome. In particular, building and deploying reliable autonomy requires designers and engineers to consider the interactive nature between subsystems. While individually subsystems may be well understood, their composition can give rise to emergent behavior (Campbell et al., 2010). Briefly, most autonomous vehicles can be modeled via two main subsystems: perception and localization, and planning and control (Figure 1.1). The essential sensors can include GPS, IMU, odometry, lidar, radar, and cameras. The perception and localization subsystem use the measurements from these sensors to localize the vehicle within the map, detect and track objects and relevant features (e.g., road, lane marking, stop lines, etc.), and create a parsable representation of the world around the vehicle. The planning component uses this information to compute a vehicle's trajectory. Finally, the control component sends actuation commands, including brake, throttle, steering, gear shifting, to keep the vehicle within this calculated trajectory.

<sup>&</sup>lt;sup>1</sup>This section partially follows Wongpiromsarn (2020).

6 Introduction

One of the main key technological developments to steer the improvement of autonomous system was the introduction Velodyne's HDL-64E sensor. This spinning lidar and its successors, including the more affordable models with 16 and 32 beams and the high-end models with 128 beams, are still a key component of many autonomous vehicles today.

The planning component is typically decomposed into three levels—the mission, the behavioral, and the trajectory planners—although naming and detail of responsibilities and algorithms varies between implementations. Roughly, the mission planner computes a high-level route for the vehicle to complete its mission. The behavioral planner is responsible for making local decisions (e.g., whether to stay in lane, proceed through an intersection, etc.) and typically is implemented as a finite-state machine. The trajectory planner then translates the decision into a trajectory for the vehicle to follow, using variations of optimization-based (e.g., model predictive control (MPC)) and graph-based (e.g., rapidly-exploring random trees (RRT) and probabilistic roadmap (PRM)) approaches. The early controllers of this era were typically based on pure pursuit and proportional-integral-derivative (PID) control. More details about planning and control algorithms can be found in Paden et al. (2016).

Similarly, to the development of modern autonomous vehicles, another example is that of unmanned aerial vehicles, which can be perceived as having the same components, except that the calculated dynamics concern with flight rather than driving. Unmanned aerial vehicle software design has improved from manual definition of state machines and ad hoc tuning of control parameters in a fly-fix-fly fashion to automated synthesis from verified and validated formal models, such as state machine representations. The hardware side has also seen both an increase in capabilities and a reduction in size, giving rise to, for example, networked clusters of very small unmanned aerial vehicles cooperating for common goals. Larger unmanned aerial vehicles have become more precise, less perturbed by uncertainties in the environmental parameters, and with faster and smaller motors.

One useful delineation for understanding modern autonomous systems

is, therefore, the following: (1) there is increasing development in highprecision sensors and controllers with reduced cost and (2) there are new software synthesis techniques that can make those controllers highly reliable and modular. In this work we will consider the second part of the modern development of autonomous systems as it related from the translation of formal specification to the synthesis of behavior on fabric, abstracting away from particular implementation details in the systems and microcontrollers.

#### Major Technological Challenges

The challenges associated with individual components (e.g., developing scalable algorithms for perception, planning, control, and contingency management) and their integration into a holistic system are further intensified by the safety-critical nature of autonomous systems. In particular, subtle design bugs may arise from the unforeseen interactions among different components and manifest as undesirable behavior only under a specific set of conditions, making them very hard to catch using simulation and testing. For example, consider an implementation of an autonomous vehicle that composes of the following components.

- The trajectory planner, which generated a path for the vehicle to follow.
- The safety system, which rapidly decelerated the vehicle when it deviated too much from the planned path and got too close to an obstacle.
- The low-level steering controller, which limited the steering rate at low speeds to protect the vehicle steering system.

Each of these functionalities is straightforward to implement in isolation. However, when combined, they can leed to unsafe behavior under specific circumstances and contexts, e.g., when the vehicle has to make a tight turn while merging into traffic facing a concrete barrier next to the major road. In this case, the vehicles planned path contains a sharp turn. Accelerating from a low speed, the controller cannot execute the turn closely due to the limited steering rate. As a result, the vehicle may

8 Introduction

deviate from the path and head instead towards the concrete barrier.

Further safety systems are necessary to avoid hazardous situations like the above, but even then the safety system may be a cause of concern as well. By activating it may slow the vehicle down as it is taking the sharp turn, leading to an even stricter limit on the steering rate. This cycle can cause the vehicle to be stuck at the corner of a sharp turn, dangerously stuttering in the middle of an intersection. The analysis presented in Wongpiromsarn et al. (2009) reveals that the software design was not inherently flawed. The undesirable behavior was caused by an unfortunate choice of certain parameters relating to the geometric properties of the planner-generated paths, compounding the challenge of safety when deploying autonomy.

In short, the key technical challenges in autonomous systems evolved around the following factors: uncertainties, complex tasks, and interconnection of computing, communication, and physical components. The uncertain and unstructured nature of environments lead to an unreasonably large number of test scenarios and give rise to the question of how to address edge cases. The complex interaction of different components can cause any change in one component to engender interaction faults once integrated with others. Finally, complex tasks are primarily handled by handcrafted implementations, e.g., via increasingly complicated finite state machines, which ends up hosting several hacks to handle corner cases encountered during testing. In particular, a naive behavioral planner architecture can be implemented as a finite state machine with less than five states. Still, to address corner cases such state machines at a minimum might need three interacting finite state machines, each containing more than 20 states, making it almost impossible to analyze or debug.

# 1.2 Overview of Verification and Synthesis

This article aims to explain frameworks for formal verification and synthesis of autonomous systems to provide a formal, mathematical guarantee of the correctness of such a system with respect to its desired properties. These systems typically consist of both low-level (continuous)

dynamics associated with the physical hardware and the high-level (discrete) logics that govern the overall behavior of the systems. Synthesis and verification of these systems thus require integration of reasoning about discrete and continuous behaviors within a single framework. A common approach to enable such integration is to construct a finite state model that serves as an abstract model of the physical system (which typically has infinitely many states) and apply formal verification and synthesis to the resulting finite state model. Several abstraction methods have been proposed based on a fixed abstraction (Kress-Gazit et al., 2007; Conner et al., 2007; Kloetzer and Belta, 2008a; Wongpiromsarn et al., 2012; Tabuada and Pappas, 2006; Girard and Pappas, 2009) and sampling (Karaman and Frazzoli, 2009; Castro et al., 2013; Wongpiromsarn et al., 2021). The focus of this article, however, is on the high-level logics and we assume that the system can be abstracted using a finite state model.

The research presented here consists of three key components: specification, synthesis, and verification. Specification refers to a precise description of the system and its desired properties. However, this precise description of the system does not need to capture all the details of the actual implementation itself. To simplify the analysis of the system, one may want to capture only the essential aspects and abstract the actual implementation in this description. In this article, we use the term "desired properties", "desired behaviors", and "requirements" interchangably to refer to what we want the system to do.

Verification is the process of checking the correctness of the system. Here, correctness is only defined relative to the desired properties. Specifications of both the system and its desired properties are essential in this process. As previously discussed, verifying the correctness of complex systems such as autonomous vehicles can be very difficult due to the interleaving between their continuous and their discrete components. Although much work exists in this domain, verifying such systems remains time-consuming and requires some level of expertise.

There has been a growing interest in the automatic synthesis of autonomous systems that provide formal system correctness guarantees

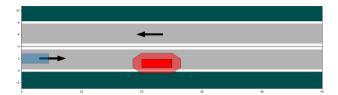


Figure 1.2: The autonomous vehicle (blue rectangle) encounters a stationary vehicle (red rectangle) on a two-lane road with a double white lane divider. The red octagon represents the clearance zone around the stationary vehicle.

to complement system verification. This type of automated synthesis can potentially reduce the time and cost of the system development cycle. It will ultimately help reduce the number of iterations between redesigning the system and verifying the new design.

#### 1.3 What Makes It Hard?

Autonomous systems typically feature a tight interaction between the computing and the physical components. A combination of model-based and data-driven approaches may be employed to design these components. Furthermore, autonomous systems are composed of multiple heterogeneous components whose complex interactions may cause any change in one component to affect others in unexpected ways.

Many autonomous systems need to perform complex tasks and are safety-critical. As a result, they are subject to strict regulations. A software bug in these systems can lead to a violation of law and morality. The complex tasks autonomous systems have to perform and regulatory requirements form a set of complex rules that the system needs to satisfy. Under certain situations, these rules may be conflicting, i.e., they cannot be simultaneously satisfied. That means that there can be requirement misalignment with our values that then clash in the eventual implementation of the system. For example, item 221 of Singapore's final theory of driving (Final Theory of Driving: The Official Handbook 2017) suggests keeping a safe gap of one meter when passing by a parked vehicle. In contrast, item 52 of Singapore's basic theory of driving (Basic Theory of Driving: The Official Handbook 2018) prohibits crossing a

solid double white lane divider. As a result, when encountering a vehicle that is improperly parked in a lane with a solid double white lane divider (Figure 1.2), an autonomous vehicle may need to violate either of these rules unless the lane is wide enough to laterally accommodate two cars with a buffer of one meter.

The uncertain and unstructured nature of environments in which the systems operate further amplifies the complexity of the verification and synthesis of these systems. In particular, the sociotechnical environments may change abruptly, drastically, and unexpectedly and may include adversaries. Such an open-world challenge leads to an unreasonably large number of test scenarios. It drives the question of edge cases from the safety and verification and the certification and regulation aspects.

#### 1.4 Organizational Outline

This monograph is segmented to incrementally introduce different types of formalisms for the analysis and synthesis of autonomous systems. In particular, we first introduce basic concepts of models—what they are, why the exist, and how they are used—and specifications—formal semantics for modeling autonomous system behavior (Section 2). Then, we show how those formal specifications can verify certain system requirements that are capturable within a logic (Section 3). We continue by showing how those specifications can synthesize behavior equipped with formal guarantees in varying mathematical settings that capture different systems scenarios (Sections 4, 5 and 6). Using these problem formulation models we show how to deal with partiality in the information the system can perceive (Section 7). We extend the synthesis problem to monitoring for runtime assurance of correct behavior (Section 8). We address the addition of learning in the verification and synthesis problem (Section 9). We conclude with the open problems in the intersection of formal methods and autonomy (Section 10).

## 1.5 A Note on the Coverage of the Article

In the rest of this article, we follow the exposition of several earlier publications by the authors. We list such publications at the beginning 12 Introduction

of the associated sections. Additionally, while the upcoming sections strive to give an objective coverage of the existing work, they do not provide a complete literature survey and, at times, are biased toward the references that had influenced the authors' own work.

# Models and Specifications

We present a series of formal verification and synthesis problems and their applications in the context of autonomous systems. This section provides a high-level view on formal verification and synthesis and their main building blocks: *models* and *specifications*.

Models are representations of our knowledge—and the limitations in it—on, e.g., the capabilities of the system of interest, the environment in which it is to operate, the uncertainties to which it is subject to, and the resources it has access. Specifications are representations of the requirements that we impose on the system's operation and the assumptions we place on the aspects that influence the system's behavior but cannot be controlled by the system.

Models and specifications are used to formally verify systems (Figure 2.1). Roughly speaking, verification is concerned with whether all (or most) executions of the system generated by implementing a given control strategy will respect a specification. The answer to a verification question is either affirmative, i.e., all possible executions of the system respect the specification, or negative, i.e., counter-example traces that

At times, the exposition in this section follows Baier and Katoen (2008).

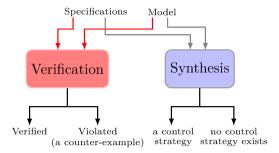


Figure 2.1: Inputs and potential outcomes of formal verification and synthesis.

don't respect the specification are output. Synthesis is concerned with the possibly more ambitious question of whether and how a system, e.g., by choice of a control strategy, can ensure the satisfaction of the specifications in all system executions. The outcome of synthesis is a strategy that, when implemented, ensures that the system satisfies the specifications or evidence that no such strategy exists.

The choice of mathematical representations and specification languages for verification and synthesis depends on several factors. Factors include the objective of verification and synthesis, the type of knowledge about the underlying system, and the uncertainties in this knowledge. A hierarchy of mathematical models and specification languages exists to address these factors. For example, one sequence of models that we will encounter in the upcoming sections is from finite transition systems to Markov decision processes to partially observable Markov decision processes to uncertain partially observable Markov decision processes. This evolution will help account for the introduction of stochastic uncertainties, limitations in the availability of run-time information, and limitations in the knowledge about the stochastic uncertainties in the way the system behaves or perceives its environment.

We introduce the first type of model in the sequence, (finite) transition systems as the basis for modeling as well as linear temporal logic as an example of a formal specification language. As they become relevant, we will introduce extensions of such models and variants of temporal logic.

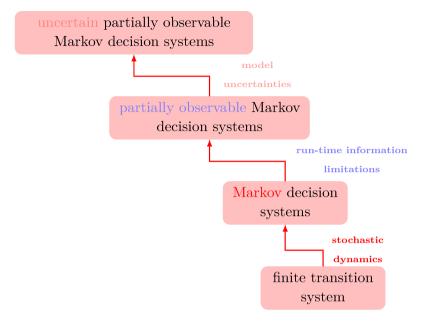


Figure 2.2: A sequence of models with increasing expressivity.

It is impossible to present all models and specification languages exhaustively. For example, one of the aspects of formally verifying autonomous systems that we do not consider is verification and synthesis using models with continuous (or hybrid) state and action spaces, which is illustrated elsewhere (Kress-Gazit *et al.*, 2007; Conner *et al.*, 2007; Kloetzer and Belta, 2008a; Wongpiromsarn *et al.*, 2012; Tabuada and Pappas, 2006; Girard and Pappas, 2009; Karaman and Frazzoli, 2009; Castro *et al.*, 2013; Wongpiromsarn *et al.*, 2021; Srinivasan *et al.*, 2018).

## 2.1 Transition Systems

We first introduce some notation. Given a set X, let  $X^*$ ,  $X^{\omega}$  and  $X^+$  denote the set of finite, infinite and nonempty finite strings, respectively, of X and let |X| denote the cardinality of X. For sequences  $\pi$ ,  $\pi_1$  and  $\pi_2$ , let  $\pi_1\pi_2$  denote a sequence obtained by concatenating  $\pi_1$  and  $\pi_2$  and let  $\pi^{\omega}$  denote an infinite sequence obtained by concatenating  $\pi$  infinitely many times.

A transition system is a mathematical description of the behavior of systems with discrete inputs, outputs, internal states, and transitions between states. Atomic propositions that express essential characteristics of individual states of the system formalize the behavior of transition systems. Roughly, a proposition is a statement that can be either true or false, but not both. An atomic proposition is a proposition whose truth or falsity does not depend on the truth or falsity of any other proposition. For example, a statement "traffic light is green" is an atomic proposition, whereas a statement "traffic light is either green or red" is not an atomic proposition.

**Definition 2.1.** A (labeled) transition system TS is a tuple

$$TS = (S, Act, \rightarrow, I, AP, L)$$

is composed of the following data.

- A set of states, S.
- A set of actions, Act.
- A transition relation,  $\rightarrow \subseteq S \times Act \times S$ .
- A set of initial states,  $I \subseteq S$ .
- A set of atomic propositions, AP.
- A labeling function  $L: S \to 2^{AP}$ .

We use the relation notation,  $s \stackrel{\alpha}{\to} s'$ , to denote  $(s, \alpha, s') \in \to$ . The transition system TS is called *finite* if S, Act and AP are finite. Note that a transition system TS may consist of a subset of these elements

The following example is borrowed from (Baier and Katoen, 2008).

**Example 2.1.** Consider a traffic light that can be either red or green. Let g denote an atomic proposition stating that the light is green. This traffic light can be modeled by a transition system

$$T = (S, Act, \rightarrow, I, AP, L),$$

where 
$$S = \{s_1, s_2\}$$
,  $Act = \{\alpha\}$ ,  $\rightarrow = \{(s_1, \alpha, s_2), (s_2, \alpha, s_1)\}$ ,  $I = \{s_1\}$ ,  $AP = \{g\}$  and  $L: S \rightarrow 2^{AP}$  is defined by  $L(s_1) = \emptyset$  and  $L(s_2) = \{g\}$ .

Given a transition system  $TS = (S, Act, \rightarrow, I, AP, L), s \in S$  and  $\alpha \in Act$ , we let

$$Act(s) = \{ \alpha \in Act : \exists s' \in S \text{ such that } s \stackrel{\alpha}{\to} s' \}$$

denote the set of enabled actions in s,

$$Post(s, \alpha) = \{s' \in S : s \xrightarrow{\alpha} s'\} \text{ and } Post(s) = \bigcup_{\alpha \in Act} Post(s, \alpha)$$

denote the set of direct successors of s. We say that TS is action-deterministic if and only if  $|I| \leq 1$  and  $|Post(s,\alpha)| \leq 1$  for all  $s \in S$  and  $\alpha \in Act$ . A sequence of states, either finite  $\pi = s_0s_1 \cdots s_n$ , or infinite  $\pi = s_0s_1 \cdots$ , is a path fragment if  $s_{i+1} \in Post(s_i)$  for all  $i \geq 0$ . A path is a path fragment such that  $s_0 \in I$  and it is either a finite path fragment that ends in a state s with  $Post(s) = \emptyset$  or an infinite path fragment. We denote the set of paths in TS by Path(TS). The trace of an infinite path fragment  $\pi = s_0s_1 \ldots$  is defined by  $trace(\pi) = L(s_0)L(s_1) \cdots$ . The set of traces of TS is defined by

$$Trace(TS) = \{trace(\pi) : \pi \in Path(TS)\}.$$

**Remark 2.1.** Transition systems that are not action-deterministic are those in which some action, when applied in some state, leads to several possible next states. Hence, they can be used to capture uncertainties in the system, especially those that arise from difference choices of valid environment behaviors over which the system does not have control.

**Example 2.2.** Consider a traffic light T (Example 2.1). An infinite sequence  $\pi = (s_1 s_2)^{\omega}$  is a path of T with the corresponding trace  $trace(\pi) = (\emptyset\{g\})^{\omega}$ . The set of paths and the set of traces of T are given by  $Paths(T) = \{\pi\}$  and  $Traces(T) = \{trace(\pi)\}$ , respectively.

Complex systems are typically composed of multiple components that can be executed at the same time. Suppose a transition system can model each component of a system. Composing the complete system amounts to composing the finite transition systems representing individual components. There are several composition techniques, depending on how the components interact, e.g., no communication to synchronous and asynchronous message transfer (Baier and Katoen, 2008). For example, hand-shaking is a mode of communication between the components that leads to synchrony. The composition of finite transition systems by hand-shaking is defined as follows.

#### **Definition 2.2.** Let two transition systems

$$TS_1 = (S_1, Act_1, \rightarrow_1, I_1, AP_1, L_1)$$
 and  $TS_2 = (S_2, Act_2, \rightarrow_2, I_2, AP_2, L_2)$ 

be given. Their composition (by hand-shaking), denoted by  $TS_1||TS_2$ , is the transition system defined by

$$TS_1||TS_2 = (S_1 \times S_2, Act_1 \cup Act_2, \rightarrow, I_1 \times I_2, AP_1 \cup AP_2, L),$$

where  $L(\langle s_1, s_2 \rangle) = L_1(s_1) \cup L_2(s_2)$  and  $\rightarrow$  is defined by the following rules:

- If  $\alpha \in Act_1 \cap Act_2$ ,  $s_1 \stackrel{\alpha}{\to} s'_1$  and  $s_2 \stackrel{\alpha}{\to} s'_2$ , then  $\langle s_1, s_2 \rangle \stackrel{\alpha}{\to} \langle s'_1, s'_2 \rangle$ .
- If  $\alpha \in Act_1 \setminus Act_2$  and  $s_1 \stackrel{\alpha}{\to} s'_1$ , then  $\langle s_1, s_2 \rangle \stackrel{\alpha}{\to} \langle s'_1, s_2 \rangle$ .
- If  $\alpha \in Act_2 \setminus Act_1$  and  $s_2 \stackrel{\alpha}{\to} s_2'$ , then  $\langle s_1, s_2 \rangle \stackrel{\alpha}{\to} \langle s_1, s_2' \rangle$ .

**Example 2.3.** Consider the system composed of 2 sets of traffic lights (Figure 2.3). For  $i \in \{1, 2\}$ , we let  $g_i$  denote an atomic proposition stating that light  $T_i$  is green. Then,

$$T_1 = (S_1, Act_1, \rightarrow_1, I_1, AP_1, L_1)$$

is a finite transition system where  $S_1 = \{s_{1,1}, s_{1,2}\}, Act_1 = \{\alpha_1\}, \\ \rightarrow_1 = \{(s_{1,1}, \alpha_1, s_{1,2}), (s_{1,2}, \alpha_1, s_{1,1})\}, I_1 = \{s_{1,1}\}, AP_1 = \{g_1\} \text{ and } L: \\ S_1 \rightarrow 2^{AP_1} \text{ is defined by } L(s_{1,1}) = \emptyset \text{ and } L(s_{1,2}) = \{g_1\}. \text{ Also,}$ 

$$T_2 = (S_2, Act_2, \rightarrow_2, I_2, AP_2, L_2)$$

is a finite transition system where  $S_2 = \{s_{2,1}, s_{2,2}\}, Act_2 = \{\alpha_2\}, \rightarrow_2 = \{(s_{2,1}, \alpha_2, s_{2,2}), (s_{2,2}, \alpha_2, s_{2,1})\}, I_2 = \{s_{2,1}\}, AP_2 = \{g_2\} \text{ and } L:$ 

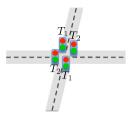
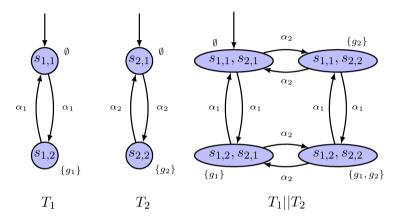


Figure 2.3: A system of traffic lights considered in Example 2.3.



**Figure 2.4:** The transition systems representing the models of traffic lights in Example 2.3.

 $S_2 \to 2^{AP_2}$  is defined by  $L(s_{2,1}) = \emptyset$  and  $L(s_{1,2}) = \{g_2\}$ . Figure 2.4 shows the graphical representation of  $T_1$  and  $T_2$  and their composition  $T_1||T_2$ . Note that  $T_1||T_2$ , is action-deterministic because at every state, the actions uniquely determines the next state.

In settings where two players, one representing the controllable system and one representing the non-controllable environment, determine the properties characterized by the atomic propositions independently, one may resort to an alternative definition of transition system. An input-output transition system, different from the transition system in Definition 2.1, treats the input atomic propositions controlled by the environment and the output atomic propositions controlled by the system separately. The overall set of atomic propositions is  $\mathcal{AP} = \mathcal{I} \cup \mathcal{O}$ , where  $\mathcal{I}$  and  $\mathcal{O}$  are disjoint sets, denoting input and output atomic propositions,

respectively.

**Definition 2.3.** An input-output transition system over a set of input propositions  $\Im$  and a set of output propositions  $\Im$  is a tuple  $\Im = (S, s_0, \tau)$ , where S is a set of states,  $s_0$  is the initial state, and the transition function  $\tau: S \times 2^{\Im} \to S \times 2^{\Im}$  maps a state s and a valuation  $\sigma_I \in 2^{\Im}$  of the input propositions to a successor state s' and a valuation  $\sigma_O \in 2^{\Im}$  of the output propositions. For any letter  $\sigma$ , the projection to input propositions is  $\sigma_I \stackrel{\text{def}}{=} \sigma \cap \Im$  and to output propositions is  $\sigma_O \stackrel{\text{def}}{=} \sigma \cap \Im$ .

An execution of  $\mathcal{T}$  is an infinite sequence  $s_0, (\sigma_{I0} \cup \sigma_{O0}), s_1, (\sigma_{I1} \cup \sigma_{O1}), s_2 \dots$  such that  $s_0$  is the initial state, and  $(s_{i+1}, \sigma_{Oi}) = \tau(s_i, \sigma_{Ii})$  for every  $i \geq 0$ . The corresponding sequence  $(\sigma_{I0} \cup \sigma_{O0}), (\sigma_{I1} \cup \sigma_{O1}), \dots \in \Sigma^{\omega}$  is considered a trace.

The input-output transition system is applicable in modeling of reactive systems, as discussed in Section 5.3.

## 2.2 Formal Specification Languages

A first step toward providing guarantees on the behavior of the design artifacts is formally expressing what operational requirements and safety constraints the system ought to satisfy. These specifications shall be in a language that is (i) mathematically based to allow automated tool development, (ii) rich enough to capture diverse set of properties, and (iii) relatively natural to allow the designers to express their intent at a high level, potentially with domain-specific terms. One such family of possible specification languages is temporal logic which is used for specifying properties of infinite sequences of states. For example, for an autonomous system, snapshots of the relevant variables needed to model the behavior of the system including its own position, availability of its actuation and motion capabilities, quality of its sensors, position and status of the other systems can all be modeled with temporal logic.

#### **Linear Temporal Logic**

The broad family of temporal-logic-based specifications provides a basis for expressing specifications for autonomous systems in a formal language. Temporal logics is a branch of logic that implicitly incorporates temporal aspects and can be used to reason about a time line (Baier and Katoen, 2008; Emerson, 1990; Huth and Ryan, 2004; Manna and Pnueli, 1992). Its use as a specification language was introduced by Pnueli (1977). Since then, temporal logic has been demonstrated to be an appropriate specification formalism for reasoning about various kinds of systems, especially those of concurrent programs. It has been used to formally specify and verify behavioral properties in various applications, including concurrent systems, reactive systems, discrete event systems, robotics and aerospace (Clarke et al., 1986; Pnueli, 1986; Galton, 1987; Lin, 1993; Holzmann, 1994; Bouma et al., 1994; Seow and Devanathan, 1994; Gabbay et al., 1995; Jagadeesan et al., 1996; Cerrito and Mayer, 1998; Jiang and Kumar, 2001; Schneider et al., 1998; Havelund et al., 2001; Holzmann, 2014).

We consider a version of temporal logic, namely linear temporal logic (LTL). An LTL formula is built up from a set of atomic propositions and two kinds of operators: logical connectives and temporal modal operators. The logic connectives are those used in propositional logic:  $negation (\neg)$ ,  $disjunction (\lor)$ ,  $conjunction (\land)$  and  $material implication (\Longrightarrow)$ . The temporal modal operators include  $next (\bigcirc)$ ,  $always (\square)$ ,  $eventually (\diamondsuit)$  and until (U). Specifically, an LTL formula over a set AP of atomic propositions is inductively defined as follows.

- (1) True is an LTL formula,
- (2) any atomic proposition  $p \in AP$  is an LTL formula and
- (3) given LTL formulas  $\varphi$ ,  $\varphi_1$  and  $\varphi_2$ ,  $\neg \varphi$ ,  $\varphi_1 \lor \varphi_2$ ,  $\bigcirc \varphi$  and  $\varphi_1 U \varphi_2$  are also LTL formulas.

Additional operators can be derived from the logical connectives  $\vee$  and  $\neg$  and the temporal modal operator  $\mathcal{U}$ . For example,  $\varphi_1 \wedge \varphi_2 = \neg(\neg \varphi_1 \vee \neg \varphi_2)$ ,  $\varphi_1 \implies \varphi_2 = \neg \varphi_1 \vee \varphi_2$ ,  $\diamond \varphi = True \ \mathcal{U} \ \varphi$  and  $\Box \varphi = \neg \diamond \neg \varphi$ .

LTL formulas are interpreted on infinite strings  $\sigma = \sigma_0 \sigma_1 \sigma_2 \cdots$  where  $\sigma_i \in 2^{AP}$  for all  $i \geq 0$ . Such infinite strings are referred to as words. The satisfaction relation is denoted by  $\models$ , i.e., for a word  $\sigma$  and an LTL formula  $\varphi$ , we write  $\sigma \models \varphi$  if and only if  $\sigma$  satisfies  $\varphi$ . The satisfaction relation is defined inductively as follows.

- $\sigma \models True$ ,
- for an atomic proposition  $p \in AP$ ,  $\sigma \models p$  if and only if  $p \in \sigma_0$ ,
- $\sigma \models \neg \varphi$  if and only if  $\sigma \not\models \varphi$ ,
- $\sigma \models \varphi_1 \land \varphi_2$  if and only if  $\sigma \models \varphi_1$  and  $\sigma \models \varphi_2$ ,
- $\sigma \models \bigcirc \varphi$  if and only if  $\sigma_1 \sigma_2 \cdots \models \varphi$  and
- $\sigma \models \varphi_1 \mathcal{U} \varphi_2$  if and only if there exists  $j \geq 0$  such that  $\sigma_j \sigma_{j+1} \cdots \models \varphi_2$  and for all i such all  $0 \leq i < j$ ,  $\sigma_i \sigma_{i+1} \cdots \models \varphi_1$ .

Let  $\varphi$  be an LTL formula over AP. The linear-time property induced by  $\varphi$  is defined as  $Words(\varphi) = \{\sigma \in (2^{AP})^{\omega} : \sigma \models \varphi\}$ . Given a transition system TS, its infinite path fragment  $\pi$  and an LTL formula  $\varphi$  over AP, we say that  $\pi$  satisfies  $\varphi$ , denoted  $\pi \models \varphi$ , if  $trace(\pi) \models \varphi$ . Finally, we say that TS satisfies  $\varphi$ , denoted  $TS \models \varphi$ , if  $Trace(TS) \subseteq Words(\varphi)$ .

The decision problem of determining whether there exists a transition system that satisfies an LTL formula is called the *realizability problem* for LTL. If an LTL formula  $\varphi$  is realizable, the goal of LTL synthesis problem is to construct a transition system TS such that  $TS \models \varphi$ .

## **Examples of LTL formulas**

Given propositional formulas p and q, important and widely used properties can be defined in terms of their corresponding LTL formulas as follows.

**Safety (invariance)** A safety formula is of the form  $\Box p$ , which asserts that the property p remains invariantly true throughout an execution. Typically, a safety property ensures that nothing bad happens. A classic example of safety property frequently used in the robot motion planning domain is obstacle avoidance.

- **Guarantee (reachability)** A guarantee formula is of the form  $\Diamond p$ , which guarantees that the property p becomes true at least once in an execution. Reaching a goal state is an example of a guarantee property.
- **Obligation** An obligation formula is a disjunction of safety and guarantee formulas,  $\Box p \lor \Diamond q$ . It can be shown that any safety and progress property can be expressed using an obligation formula. (By letting  $q \equiv False$ , we obtain a safety formula and by letting  $p \equiv False$ , we obtain a guarantee formula.)
- **Progress (recurrence)** A progress formula is of the form  $\Box \Diamond p$ , which essentially states that the property p holds infinitely often in an execution. As the name suggests, a progress property typically ensures that the system makes progress throughout an execution.
- **Response** A response formula is of the form  $\Box(p \Longrightarrow \Diamond q)$ , which states that following any point in an execution where the property p is true, there exists a point where the property q is true. A response property can be used, for example, to describe how the system should react to changes in the operating conditions.
- **Stability (persistence)** A stability formula is of the form  $\Diamond \Box p$ , which asserts that there is a point in an execution where the property p becomes invariantly true for the remainder of the execution. This definition corresponds to the definition of stability in the controls domain since it ensures that eventually, the system converges to a desired operating point and remains there for the remainder of the execution.
- **Example 2.4.** Consider the system of traffic lights (Example 2.3). Its desired properties might include the following.
  - "At least one of the lights is always on" is a safety property and can be expressed in LTL as  $\Box(g_1 \vee g_2)$ .
  - "Two lights are never green at the same time" is also a safety property and can be expressed in LTL as  $\Box(\neg g_1 \lor \neg g_2)$ .

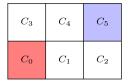


Figure 2.5: The robot environment of Example 2.5.

• " $T_1$  will turn green infinitely often" is known as a "progress" property and can be expressed in LTL as  $\Box \Diamond g_1$ .

**Example 2.5.** Consider a robot motion planning problem where the robot is moving in an environment that is partitioned into six regions (Figure 2.5).

Let s represent the position of the robot and  $C_0, \dots, C_5$  represent the polygonal regions in the robot environment. Suppose the robot receives an externally triggered PARK signal. Consider the following desired behaviors.

- (a) Visit region  $C_5$  infinitely often.
- (b) Eventually go to region  $C_0$  when a PARK signal is received.

Assuming that infinitely often, a PARK signal is not received, the desired properties of the system can be expressed in LTL as

$$\Box \Diamond (\neg park) \implies (\Box \Diamond (s \in C_5) \land \Box (park \implies \Diamond (s \in C_0))). \tag{2.1}$$

Here, park is a Boolean variable that indicates whether a PARK signal is received.

Let  $S_1 = \{s_0, s_1, \dots, s_5\}$  be a finite set of the (discretized) positions of the robot such that  $s_i \in C_i$  for all i, Let  $S_2 = \{p, p'\}$ , where p and p' indicates that the PARK signal is received and is not received, respectively,  $Act = \{l, r, u, d\}$ , where l, r, u, d represent the left, right, up, and down movement of the robot, respectively,  $I = \{s_0\}$ , and  $AP = \{C_1, C_2, \dots, C_5, park\}$ . We can model the complete system as a finite transition system  $TS = (S, Act, \rightarrow, I, AP, L)$  where  $S = S_1 \times S_2$ ,  $((s_i, \tilde{p}), \alpha, (s_i, \hat{p})) \in \rightarrow$  if and only if  $C_i$  and  $C_j$  are adjacent,  $\alpha$ 

corresponds to the location of  $C_j$  relative to  $C_i$ , and  $\tilde{p}, \hat{p} \in \{p, p'\}$ . The labeling function  $L: S \to 2^{AP}$  is defined by  $L(s_i, p) = \{C_i, park\}$  and  $L(s_i, p') = \{C_i\}$  for all i.

A path  $\pi = ((s_0, p), (s_1, p), (s_2, p'), (s_5, p), (s_4, p), (s_3, p'))^{\omega}$  of TS satisfies the LTL formula 2.1 while  $\pi' = ((s_0, p), (s_1, p'))^{\omega}$  does not satisfy this formula.

**Example 2.6.** An autonomous vehicle competing in the DARPA Urban Challenge is required to follow traffic rules as well as completing a task specified by a sequence of checkpoints that the vehicle has to cross (Figure 2.6). We define the state of the autonomous vehicle as  $(x, \theta, v)$  where  $x \in \mathbb{R}^2$  represents the center of its front bumper,  $\theta \in \mathbb{R}$  represents its heading, and  $v \in \mathbb{R}$  represents its speed. Additionally, let  $\mathsf{Obs} \subset \mathbb{R}^2$  represent the union of the footprints of all the vehicles and obstacles in the environment. The state of the complete system (autonomous vehicle and the environment) is then given by  $(x, \theta, v, \mathsf{Obs})$ . Examples of some desired properties and LTL formulas expressing those properties are given below.

Traffic rule 1. No collision:

$$\Box(\mathsf{FP}(x,\theta) \cap \mathsf{Obs} = \emptyset),$$

where for any  $x, \theta \in \mathbb{R}$ ,  $\mathsf{FP}(x, \theta) \subset \mathbb{R}^2$  is the footprint of the autonomous vehicle when the center of its front bumper is at x and its heading is  $\theta$ .

Traffic rule 2. Obey speed limits:

$$\Box((x \in \text{Reduced\_Speed\_Zone}) \implies (v \leq v_{reduced})),$$

where  $v_{reduced}$  is a pre-specified parameter for the maximum speed in Reduced\_Speed\_Zone.

**Goal.** Eventually visit the checkpoint:  $\Diamond(x = \text{ck\_pt})$  where ck\_pt denote the position of the checkpoint.

The following are considered to be atomic propositions in this example as they are evaluated to be either true or false, given the state  $(x, \theta, v, \mathsf{Obs})$ 

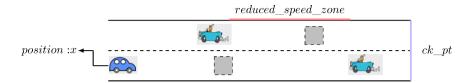


Figure 2.6: A simplified autonomous driving problem considered in Example 2.6.

of the system:

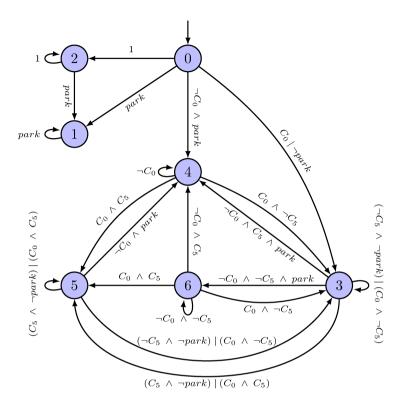
$$\begin{aligned} & \mathsf{FP}(x,\theta) \ \cap \ \mathsf{Obs} = \emptyset, \\ & x \in \mathsf{Reduced\_Speed\_Zone}, \\ & v \leq v_{reduced}, \ \mathsf{and} \\ & x = \mathsf{ck\_pt}. \end{aligned}$$

**Remark 2.2.** LTL offers an extension to the properties, e.g., safety (in the form of constraints on the system state) and reachability (in the form of convergence to a desired state) that have typically been used in the controls literature.

#### Automata Representation of LTL Formulas

There are several ways for checking whether the behaviors (i.e., traces) that can be generated by a transition system satisfy a given temporal logic specification. One notion that will appear frequently in subsequent sections is that of an automaton that witnesses the satisfaction of a temporal logic formula.

The language accepted by an LTL formula can equivalently be represented by a nondeterministic (or universal) Büchi (or co-Büchi) automaton (Büchi, 1990). A Büchi automaton over a finite alphabet  $\Sigma$  is a tuple  $\mathcal{A} = (Q, Q_0, \delta, F)$ , where Q is a finite set of states,  $Q_0$  is the set of initial states,  $\delta \subseteq Q \times \Sigma \times Q$  is the transition relation, and  $F \subseteq Q$  is a subset of states. A run of  $\mathcal{A}$  on an infinite word  $w = \sigma_0 \sigma_1 \cdots \in \Sigma^{\omega}$  is an infinite sequence  $q_0, q_1, \ldots$  of states, where  $q_0 \in Q_0$  is an initial state and for every  $i \geq 0$  it holds that  $(q_i, \sigma_i, q_{i+1}) \in \delta$ .



**Figure 2.7:** Büchi automaton for the LTL formula  $\Box \Diamond (\neg park) \implies (\Box \Diamond (s \in C_5) \land \Box (park) \implies \Diamond (s \in C_0))$  in Example 2.5.

A run of a Büchi automaton is accepting if it contains infinitely many occurrences of states in F. A co-Büchi automaton  $\mathcal{A} = (Q, Q_0, \delta, F)$  differs from a Büchi automaton in the accepting condition: a run of a co-Büchi automaton is accepting if it contains only finitely many occurrences of states in F. For a Büchi automaton the states in F are called accepting states, while for a co-Büchi automaton they are called rejecting states. A nondeterministic automaton  $\mathcal{A}$  accepts a word  $w \in \Sigma^{\omega}$  if some run of  $\mathcal{A}$  on w is accepting. A universal automaton  $\mathcal{A}$  accepts a word  $w \in \Sigma^{\omega}$  if every run of  $\mathcal{A}$  on w is accepting.

Figure 2.7 and Figure 2.8 represent examples of Büchi automata for the LTL formulas in Example 2.5 and Example 2.6, respectively.

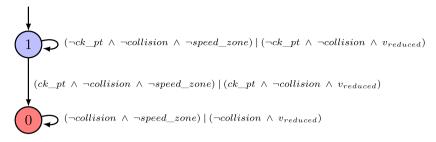


Figure 2.8: Büchi automaton for the LTL formula  $\Box (\mathsf{FP}(x,\theta) \cap \mathsf{Obs} = \emptyset) \land \Box ((x \in \mathsf{Reduced\_Speed\_Zone}) \Longrightarrow (v \leq v_{reduced})) \land \Diamond (x = \mathsf{ck\_pt})$  in Example 2.6. Here,  $ck\_pt$ , collision,  $speed\_zone$ , and  $v_{reduced}$  represent  $x = \mathsf{ck\_pt}$ ,  $\mathsf{FP}(x,\theta) \cap \mathsf{Obs} \neq \emptyset$ ,  $x \in \mathsf{Reduced\_Speed\_Zone}$ , and  $v \leq v_{reduced}$ , respectively.

While the class of nondeterministic Büchi automata are sufficient to represent any LTL formula, there are other classes of automata that are widely used. The major variations in the definition of these different automata relate to their input, states, transition function, and acceptance condition. Other automata types include deterministic and nondeterministic finite-state automata (McCulloch and Pitts, 1943), deterministic Büchi automata (Büchi, 1990), Rabin automata (Rabin and Scott, 1959), and parity automata (Gradel and Thomas, 2002). For instance, deterministic and nondeterministic finite-state automata are expressive enough to represent regular languages. Therefore, if an application requires only finite-horizon specifications, its LTL formulas can be represented by finite-state automata that are simpler than Büchi automata.

# **Verification and Model Checking**

This section reviews existing approaches to system verification. These methods perform automated analysis of the correctness of the system's abstract mathematical model relative to the system's requirements. As a result, these approaches provide a formal guarantee that the desired system properties hold over all of its possible executions, provided that the actual execution of the system respects its model.

## 3.1 Model Checking

Model checking is a well-established technique for system verification based on exhaustive exploration of the state space. The key requirement of this technique is that the description of the system and its requirements be formulated in some precise mathematical language. From the description of the system, all of its possible behaviors can be derived. In addition, all the valid and invalid behaviors can be obtained from the system requirements. A model checker then checks whether an intersection of all the possible behaviors of the system and all the invalid behaviors is empty. It terminates with a yes/no answer and provides an error trace in case of a negative result. This technique is

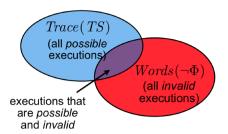


Figure 3.1: The basic idea behind model checking.

very attractive because it is automatic, fast and requires no human interaction. However, as it is based on exhaustive exploration of the state space, model checking is limited to finite state systems. It also faces a combinatorial blow up of the state space, commonly known as the state explosion problem (Holzmann, 2004; Baier and Katoen, 2008).

A model checking problem is to find a path in the finite transition system TS that violates the specification  $\varphi$ . All the possible behaviors of TS can be captured by its trace, Trace(TS) whereas all the invalid behaviors of the system can be captured by the linear-time property  $Words(\neg \varphi)$  (Chapter 2). The satisfiability problem – determining whether TS satisfies  $\varphi$  – can then be solved by claiming that  $Trace(TS) \cap Words(\neg \varphi) = \emptyset$  as shown in Figure 3.1. In case of negative result, a word in  $Trace(TS) \cap Words(\neg \varphi)$  is a counterexample. A positive result means  $Trace(TS) \cap Words(\neg \varphi) = \emptyset$ , i.e., a path  $\pi$  of TS that violates  $\varphi$  does not exist; hence, we can conclude that  $\varphi$  is satisfied.

For simplicity of the presentation, we assume that TS only has one valid initial state. For TS with multiple valid initial states, the procedure described below can be applied to each initial state separately. To check whether  $Trace(TS) \cap Words(\neg \varphi) = \emptyset$ , we first compute a non-deterministic Büchi automaton  $\mathcal{A} = (Q, \delta, Q_0, F)$  over  $2^{AP}$  that accepts all and only words over AP that satisfy  $\neg \varphi$ . The product  $TS_p = TS \otimes \mathcal{A}$  can then be constructed based on the following definition.

**Definition 3.1.** Let  $TS = (S, Act, \rightarrow, I, AP, L)$  be a transition system

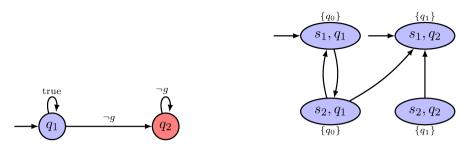


Figure 3.2: (Left) A non-deterministic Büchi automaton  $\mathcal{A}$  that accepts all and only words over  $AP = \{g\}$  that satisfy  $\neg \Box \Diamond g$ . (Right) The product  $T \otimes \mathcal{A}$  for Example 3.1

and  $\mathcal{A} = (Q, \delta, Q_0, F)$  be a non-deterministic Büchi automaton over  $2^{AP}$ . The product of TS and  $\mathcal{A}$  is the transition system  $TS_p = TS \otimes \mathcal{A}$  defined by  $TS_p = (S \times Q, Act, \rightarrow_p, I_p, Q, L_p)$  where

- (i) for any  $s, t \in S$ ,  $\alpha \in Act$  and  $p, q \in Q$ ,  $\langle s, p \rangle \xrightarrow{\alpha}_{p} \langle t, q \rangle$  if and only if  $s \xrightarrow{\alpha} t$  and  $p \xrightarrow{L(t)} q$ ,
- (ii)  $I_p = \{\langle s_0, q_0 \rangle : s_0 \in I \text{ and } \exists q \in Q_0 \text{ such that } q \stackrel{L(s_0)}{\longrightarrow} q_0 \}$  and
- (iii)  $L_p: S \times Q \to 2^Q$  is given by  $L_p(\langle s, q \rangle) = \{q\}.$

**Example 3.1.** Consider the traffic light T in Example 2.1. Suppose we want to ensure that the light is green infinitely often. This requirement can be expressed in LTL as  $\varphi = \Box \diamondsuit g$ . A non-deterministic Büchi automaton  $\mathcal{A}$  that accepts all and only words over  $AP = \{g\}$  that satisfy  $\neg \varphi$  as well as the product  $T \otimes \mathcal{A}$  are shown in Figure 3.2.

Consider a path  $\pi_p = \langle s_0, q_0 \rangle \langle s_1, q_1 \rangle \dots$  on  $TS_p$ . We say that  $\pi_p$  is accepting if and only if there exist infinitely many  $j \geq 0$  such that  $q_j \in F$ . Stepping through Definition 3.1 shows that given a path  $\pi_p$  on  $TS_p$ , the corresponding path  $\pi = s_0 s_1 \dots$  on TS generates a word  $L(s_0)L(s_1)\dots$  that satisfies  $\neg \varphi$  if and only if  $\pi_p$  is accepting. Hence, an accepting path of  $TS_p$  uniquely corresponds to a path of TS that violates  $\varphi$ . As a result, model checking can be reduced to a graph search problem to find a state  $\langle s, q \rangle$  in  $TS_p$  satisfying the following conditions:

- (MC1)  $L_p(\langle s, q \rangle) \in F$ .
- (MC2)  $\langle s, q \rangle$  is reachable, i.e., there exists a finite path fragment  $\pi_p^p$  from some  $\langle s_0, q_0 \rangle \in I_p$  to  $\langle s, q \rangle$  in  $TS_p$ .
- (MC3)  $\langle s, q \rangle$  is on a direct cycle, i.e., there exists a finite path fragment  $\pi_p^c$  from some  $\langle s', q' \rangle \in Post(\langle s, q \rangle)$  to  $\langle s, q \rangle$  in  $TS_p$ .

If such  $\langle s,q \rangle$  does not exists, we can conclude that  $Trace(TS) \cap Words(\neg \varphi) = \emptyset$  and therefore  $TS \models \varphi$ . Otherwise, an accepting path  $\pi_p$  on  $TS_p$  can be simply defined by  $\pi_p = \pi_p^p(\pi_p^c)^\omega$ . A path  $\pi = s_0 s_1 \dots$  on TS corresponding to  $\pi_p$  is a counterexample of a run on TS that violates  $\varphi$ .

**Example 3.2.** Let us revisit Example 3.1. Recall that  $F = q_2$  and the set of states of  $TS_p$  is given by  $\{s_1, s_2\} \times \{q_1, q_2\}$  as shown in Figure 3.2. First, consider the states  $(s_1, q_1)$  and  $(s_2, q_1)$ . Here,  $L((s_1, q_1)) = L((s_2, q_1)) = \{q_1\} \not\subset F$ ; thus, condition (MC1) fails. Next, the state  $(s_1, q_2)$  is not on a direct cycle; thus, condition (MC3) fails. Finally, the state  $(s_2, q_2)$  is not reachable; thus, condition (MC2) fails. We thus conclude that there is no accepting path  $\pi_p$  on  $TS_p$  and so the transition system TS modeling the traffic light T satisfies the specification  $\varphi = \Box \diamondsuit q$ .

## 3.2 Computational Complexity

The number of states of  $\mathcal{A}$  is exponential in the length of  $\varphi$ , i.e., the number of operators in  $\varphi$ . Hence, the number of states of the product transition system  $TS_p$  is  $O(|S|)2^{O(|\varphi|)}$  where |S| is the number of states in TS and  $|\varphi|$  is the length of  $\varphi$ . A nested depth-first search algorithm (Baier and Katoen, 2008) can be used to detect accepting cycles efficiently, with the worst-case time complexity that is linear in the number of states and transitions of  $TS_p$ .

Several reduction techniques have been proposed to allow model checkers to handle large state spaces. One example is state compression, including lossy compression (e.g., hash-compact and bitstate hashing), lossless compression, and alternate state representation methods, help

reduce memory requirements by reducing the amount of memory required to store each state. In contrast, partial order reduction strategies avoid computing equivalent paths, which helps reduce the number of states that needs to be explored. Symbolic model checkers such as SMV and NuSMV use compressed representation of the state space known as binary decision diagram (BDD). On the other hand, Spin avoids constructing the complete state space by employing on-the-fly construction of the finite transition system, the non-deterministic Büchi automaton, and the product automaton. We refer the reader to (Holzmann, 2004) for more details on these reduction techniques.

### 3.3 Tools for Model Checking

There are various model checkers for different specification languages. TLC (Yu et al., 1999) is a model checker for specifications written in TLA+, which is a specification language based on temporal logic of actions (TLA) (Abadi and Lamport, 1994; Lamport, 1983; Lamport, 1994). TLA introduces new kinds of temporal assertions (Lamport, 1994) to traditional linear temporal logic (LTL) to make it practical to describe a system by a single formula and to make the specifications simpler and easier to understand.

The Spin model checker deals with specifications written in process meta-language (PROMELA) (Holzmann, 2004). This language was influenced by Dijkstra, Hoare's CSP language and C. It emphasizes the modeling of process synchronization and coordination, not computation and is not meant to be analyzed by a human. Spin can be run in two modes—simulation and verification. The simulation mode performs random or iterative simulations of the modeled system's execution while the verification mode generates a C program that performs a fast exhaustive verification of the system state space. Spin is mainly used for checking for deadlocks, livelocks, unspecified receptions, unexecutable code, correctness of system invariants and non-progress execution cycles. It also supports the verification of linear time temporal constraints. Spin has been used in many applications, especially in proving correctness of safety-critical software (Havelund et al., 2001; Gluck and Holzmann,

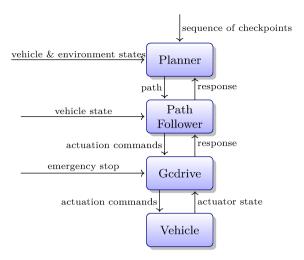


Figure 3.3: The embedded control component of Alice.

2002). Other popular model checkers include Symbolic Model Verifier (SMV) (McMillan, 1993) and its successor NuSMV (Cimatti *et al.*, 2002).

## 3.4 Model Checking for Autonomous Vehicles

This section illustrates the applications of model checking (Section 3.1) to the embedded control component of Alice, an autonomous vehicle built at the California Institute of Technology for the DARPA Urban Challenge. Alice was equipped with 25 CPUs and utilized a networked control system architecture to provide high performance and modular design. The embedded control component of Alice is shown in Figure 3.3. We refer the reader to (Burdick *et al.*, 2007; DuToit *et al.*, 2008; Wongpiromsarn and Murray, 2008) for more details on this hierarchical control architecture.

The case studies presented in this section focus on the low-level module, namely Gcdrive, which is the overall driving software for Alice. It takes independent commands from Path Follower and DARPA and sends appropriate commands to the actuators. Commands from Path Follower include control signals to throttle, brake and transmission. Commands

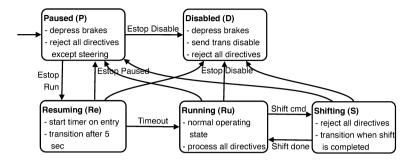


Figure 3.4: Finite state machine implemented in Actuation Interface.

from DARPA include estop pause, estop run and estop disable. An estop pause command should cause the vehicle to be brought quickly and safely to a rolling stop and reject commands to any actuator. An estop run command resumes the operation of the vehicle. An estop disable command is used to stop the vehicle and put it in the disable mode. A vehicle that is in the disable mode may not restart in response to an estop run command.

The logic in Gcdrive to handle these concurrent commands can be described by a finite state machine shown in Figure 3.4, which is implemented in the Actuation Interface component of Gcdrive (See Figure 3.5). This example illustrates the use of model checking in proving the correctness of the implementation of this finite state machine. We model Follower, Gcdrive, and DARPA (see Figure 3.5) in the Spin model checker with the following global variables.

- state ∈ {DISABLED (D), PAUSED (P), RUNNING (Ru), RE-SUMING (Re), SHIFTING (S)} is the state of the finite state machine as described in Figure 3.4.
- $estop \in \{ DISABLE (0), PAUSE (1), RUN (2) \}$  is the emergency stop command sent by DARPA.
- $acc \in [-1, 1]$  and  $acc\_cmd \in [-1, 1]$  are the acceleration commands sent from Actuation Interface to Acceleration Module and from Follower to Actuation Interface, respectively.

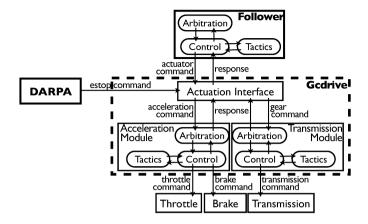


Figure 3.5: The components involved in the Gcdrive FSM example.

- $gear \in \{-1,0,1\}$  and  $gear\_cmd \in \{-1,0,1\}$  are the gear commands sent from Actuation Interface to Transmission Module and from Follower to Actuation Interface, respectively.
- $timer \in \{0, 1, 2, 3, 4, 5\}$  keeps track of the time after which the latest estop run command is received.

The desired properties can be expressed in LTL as follows.

(1) If DARPA sends an estop disable command, Gcdrive state will eventually stay at DISABLED and Acceleration Module will eventually command full brake forever.

$$\Box \big( (estop = 0) \implies \Diamond \Box (state = D \land acc = -1) \big) \tag{3.1}$$

(2) If DARPA sends an estop pause command while the vehicle is not disabled, eventually Gcdrive state will be PAUSED.

$$\Box((estop = 1 \land state \neq D) \implies \Diamond(state = P))$$
 (3.2)

(3) If DARPA sends an estop run command while the vehicle is not disabled, eventually Gcdrive state will be RUNNING or RESUMING or DARPA will send an estop disable or estop pause command.

$$\Box((estop = 2 \land state \neq D) \implies \Diamond(state \in \{Ru, Re\} \lor estop \neq 2))$$
(3.3)

(4) If the current state is RESUMING, eventually the state will be RUNNING or DARPA will send an estop disable or pause command.

$$\Box((state = Re) \implies \Diamond(state = Ru \lor estop \in \{0, 1\}))$$
 (3.4)

(5) The vehicle is disabled only after it receives an estop disable command.

$$((state \neq D)U(estop = 0)) \vee \Box(state \neq D)$$
 (3.5)

(6) Actuation Interface sends a full brake command to the Acceleration Module if the current state is DISABLED, PAUSED, RESUMING or SHIFTING. In addition, if the vehicle is disabled, then the gear is shifted to 0.

$$\Box(state \in \{D, P, Re, S\} \implies acc = -1) \land \Box(state = D \implies gear = 0)$$
(3.6)

(7) After receiving an estop pause command, the vehicle may resume the operation 5 seconds after an estop run command is received.

$$\Box(state = \text{Ru} \implies timer \ge 5) \tag{3.7}$$

Model checking revealed that an early implementation of Gcdrive failed to satisfy the above properties. In particular, the counterexample showed that the variable *state* did not change as required when an estop command was sent. The counterexample suggested that we incorporate the following assumptions.

- (a) Actuation Interface gets executed at least once each time an estop command is sent.
- (b) Actuation Interface reads the current estop status at the beginning of each iteration. It then performs a computation based on this estop status for the rest of the iteration.
- (c) All the estop commands are eventually received.

We introduce a global variable enableEstop to incorporate assumption (a). Assumption (b) is enforced using atomic sequences (See (Holzmann, 2004)). Lastly, assumption (c) is enforced by letting the variable estop represent the estop command received by Gcdrive as well. With these assumptions, Spin can verify the correctness of the system with respect to the desired properties. The PROMELA models of the components involved in this example can be found in (Wongpiromsarn, 2010).

Realizing that these assumptions needed to be enforced, we then modified the implementation of Alice by having Gcdrive store all the estop commands in a queue and process all these commands one by one. If an estop command is not stored but only sampled at the beginning of each iteration, an estop disable or pause command may not be handled appropriately. Consider, for example, the case where an estop pause command is sent while Actuation Interface is in the middle of an iteration and an estop run command is sent immediately after. In this case, Alice will not stop because the estop pause command is not processed, leading to an incorrect, unsafe behavior.

4

# **Closed-System Synthesis**

This section provides an overview of correct-by-construction synthesis of control protocols for autonomous systems and discusses approaches that merge concepts from formal methods and controls. These concepts include but are not limited to formal specification languages, discrete protocol synthesis, and optimization-based control. It also points to approaches that deal with settings where the set of desired specifications is not realizable as a whole.

We consider a discrete system modeled as an action-deterministic finite transition system, i.e., at every state, the actions uniquely determine the next state. We refer to these systems as deterministic systems. In particular, we consider closed systems by referring to systems whose outputs are generated purely by themselves without any exogenous input. We assume that at any time instance, the state of the system is fully observable.

## 4.1 Control Protocol Synthesis

In this section, we are interested in synthesizing a control protocol for a transition system to ensure that a given linear temporal logic (LTL) specification is satisfied. We define a control protocol for a transition system as follows:

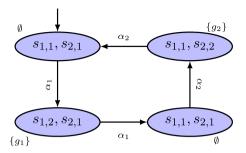
**Definition 4.1.** Let  $TS = (S, Act, \rightarrow, I, AP, L)$  be a transition system. A control protocol for TS is a function  $u: S^+ \rightarrow Act$  such that  $u(s_0s_1...s_n) \in Act(s_n)$  for all  $s_0s_1...s_n \in S^+$ , where  $S^+$  denotes the set of nonempty finite strings of S.

A control protocol u for a transition system TS essentially restricts the non-deterministic choices in TS by picking an action based on the path fragment that leads to the system's current state. Hence, u induces a transition system  $TS^u$  that formalizes the behavior of TS under the control protocol u.

In general,  $TS^u$  contains all of the states in  $S^+$ ; therefore, the induced transition system may not be finite even though TS is finite. However, for special cases where u is a memoryless or a finite-memory control protocol, it can be shown that  $TS^u$  can be identified with a finite transition system. Roughly, a memoryless control protocol picks an action based on the current state of TS, irrespective of the path fragment that led to that state. For example, a memoryless control protocol is sufficient for a car in a controlled intersection; the car proceeds if the light is green and stops if the light is red. A finite memory control protocol, on the other hand, maintains a "mode", picks an action based on the current mode and the current state of TS, and modifies the mode according to the next state. For example, an uncontrolled intersection environment with stop signs requires cars to have finite memory control protocols. The "mode" of the controller keeps the number of cars that arrived at the intersection before, and the car proceeds if and only if there are no cars that arrived before.

**Example 4.1.** Consider the transition system that represents the complete traffic light system in Example 2.3 (see Figure 2.4). Define a control protocol  $u: S^+ \to Act$  such that

• 
$$u(\langle s_{1,1}, s_{2,1} \rangle) = \alpha_1,$$



**Figure 4.1:**  $(T_1||T_2)^u$ , the transition systems induced by applying u defined in Example 4.1 on the traffic light system  $T_1||T_2$ .

- $u(\pi\langle s_{1,1}, s_{2,1}\rangle\langle s_{1,2}, s_{2,1}\rangle) = \alpha_1,$
- $u(\pi \langle s_{1,2}, s_{2,1} \rangle \langle s_{1,1}, s_{2,1} \rangle) = \alpha_2,$
- $u(\pi \langle s_{1,1}, s_{2,1} \rangle \langle s_{1,1}, s_{2,2} \rangle) = \alpha_2$ , and
- $u(\pi \langle s_{1,1}, s_{2,2} \rangle \langle s_{1,1}, s_{2,1} \rangle) = \alpha_1,$

for any  $\pi \in S^*$ . According to the transition system  $T_1||T_2|$  in Figure 2.4, the initial state is  $\langle s_{1,1}, s_{2,1} \rangle$ . At this state, the control protocol picks the action  $u(\langle s_{1,1}, s_{2,1} \rangle) = \alpha_1$ , causing the system to transition to the state  $\langle s_{1,2}, s_{2,1} \rangle$ . The path fragment so far is then given by  $\langle s_{1,1}, s_{2,1} \rangle \langle s_{1,2}, s_{2,1} \rangle$ . Thus, the control protocol picks the action  $u(\pi\langle s_{1,1}, s_{2,1}\rangle\langle s_{1,2}, s_{2,1}\rangle) = \alpha_1$ , with  $\pi = \emptyset$ , causing the system to transition back to the initial state  $\langle s_{1,1}, s_{2,1} \rangle$ . Following this procedure, the next action is  $u(\pi\langle s_{1,2}, s_{2,1}\rangle\langle s_{1,1}, s_{2,1}\rangle) = \alpha_2$ , with  $\pi = \langle s_{1,1}, s_{2,1}\rangle$ . The system then transition to  $\langle s_{1,1}, s_{2,2} \rangle$ , at which the control action  $u(\pi\langle s_{1,1}, s_{2,1}\rangle\langle s_{1,1}, s_{2,2}\rangle) = \alpha_2$  is applied, causing the system to transition back to the initial state  $\langle s_{1,1}, s_{2,1} \rangle$  once again. The transition system induced by u is shown in Figure 4.1. This transition system satisfies  $\Box(\neg g_1 \lor \neg g_2)$  since no state in the induced system has a label including  $g_1$  and  $g_2$ . It also satisfies  $\Box \Diamond g_1$  since state  $\langle s_{1,2}, s_{2,1} \rangle$  with label  $\{g_1\}$  is visited infinitely often. However, it violates  $\square(g_1 \vee g_2)$ since states  $\langle s_{1,1}, s_{2,1} \rangle$  and  $\langle s_{1,2}, s_{2,2} \rangle$  in the induced system do not have labels including  $g_1$  or  $g_2$ .

## Control Protocol Synthesis Problem: Given a finite transition

$C_{21}$	$C_{22}$	$C_{23}$	$C_{24}$	$C_{25}$
$C_{16}$	$C_{17}$	$C_{18}$	$C_{19}$	$C_{20}$
$C_{11}$	$C_{12}$	$C_{13}$	$C_{14}$	$C_{15}$
$C_6$	$C_7$	$C_8$	$C_9$	$C_{10}$
$C_1$	$C_2$	$C_3$	$C_4$	$C_5$

Figure 4.2: The robot motion planning problem in Example 4.2. Each cell represents a state of TS. The possible transitions are between the adjacent cells.

system TS and a specification  $\varphi$  expressed as an LTL formula, automatically synthesize a control protocol u such that the induced transition system  $TS^u$  satisfies  $\varphi$ .

**Example 4.2.** Consider the robot motion planning problem where the robot navigates in an area that is partitioned into cells as shown in Figure 4.2. The dynamics of the robot is abstracted to a finite transition system TS shown on the right of Figure 4.2. The state of TS represents the cell occupied by the robot. If TS is action-deterministic, then the system is deterministic. Otherwise, the system is non-deterministic. In this case, non-determinism potentially arises due to disturbances that affect the dynamics of the robot, leading to multiple possible next states when an action is taken. The desired property of the system is for a robot to visit cell  $C_8$ , then  $C_1$  and, subsequently, cover  $C_{10}$ ,  $C_{17}$ , and  $C_{25}$  in any order while always avoiding cells  $C_2$ ,  $C_{14}$ , and  $C_{18}$ . This property can be expressed in LTL as

$$\diamondsuit \big( C_8 \land \diamondsuit (C_1 \land \diamondsuit C_{10} \land \diamondsuit C_{17} \land \diamondsuit C_{25}) \big) \land \Box \neg (C_2 \lor C_{14} \lor C_{18}).$$

**Example 4.3.** Consider the simplified autonomous driving problem described in Example 2.6. The system consists of the autonomous

vehicle, obstacles (i.e., Obs in Traffic rule 1), and other vehicles (i.e., Veh in Traffic rule 1). If the obstacles and other vehicles are not stationary, and their motion is not known exactly, then the system is non-deterministic since the system does not have control over the motion of the obstacles and other vehicles. In this case, a control protocol for this system needs to ensure that the desired properties described in Example 2.6 are satisfied for all the possible motion (i.e., behavior) of the obstacles and other vehicles.

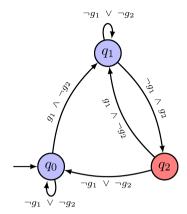
#### 4.2 Model Checking-Based Synthesis

Consider a closed system that is modeled as an action-deterministic finite transition system TS. The control-protocol synthesis problem can be formulated as finding a path in TS that satisfies a given specification  $\varphi$ , which is essentially a model checking problem as described in Section 3.1. Using model checking methods, we can start with the hypothesis that  $\varphi$  is not satisfiable; that is, we investigate whether or not there exist a path  $\pi$  of TS that satisfies  $\varphi$ . In case the hypothesis can be refuted with a counterexample, the counterexample can be used as a synthesized path  $\pi$  of TS that satisfies  $\varphi$ .

Let  $\mathcal{A}$  be a non-deterministic Büchi automaton over  $2^{\mathcal{AP}}$  that accepts all and only words over  $\mathcal{AP}$  that satisfy  $\varphi$ . We construct the product  $TS_p = TS \otimes \mathcal{A}$  as described in Definition 3.1. Let  $\pi_p = \langle s_0, q_0 \rangle \langle s_1, q_1 \rangle \dots$  be an accepting path on  $TS_p$  and  $\pi = s_0 s_1 \dots$  be the path on TS corresponding to  $\pi_p$ . We define a control protocol u for TS by

$$u(s_0's_1'\ldots s_i') = \begin{cases} \alpha_i & \text{if } s_0's_1'\ldots s_i' = s_0s_1\ldots s_i, \\ \alpha_i' & \text{otherwise,} \end{cases}$$

where  $\alpha_i$  satisfies  $s_i \stackrel{\alpha_i}{\to} s_{i+1}$  and  $\alpha'_i \in Act(s'_i)$  is any arbitrary action. Under the control protocol u, the system simply picks the next state according to the path  $\pi$ , which ensures that the resulting path satisfies  $\varphi$ . Note that this construction of control protocol works because we consider a closed system, which has a full control over the non-deterministic choices in TS and is not affected by exogenous inputs (e.g., from the environment).



**Figure 4.3:** Non-deterministic Büchi automaton  $\mathcal{A}$  that recognizes  $\varphi = \Box(\neg g_1 \lor \neg g_2) \land \Box \diamondsuit g_1 \land \Box \diamondsuit g_2$ . Accepted states are drawn with a double (red) circle.

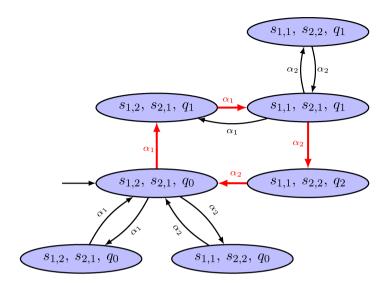
**Example 4.4.** Consider the traffic light system  $TS = T_1 || T_2$  shown in Figure 2.4 and the desired property  $\varphi = \Box(\neg g_1 \lor \neg g_2) \land \Box \diamondsuit g_1 \land \Box \diamondsuit g_2$ . In words, the property makes sure that the two lights are never green at the same time and each light turns green infinitely often. A non-deterministic Büchi automaton  $\mathcal{A}$  that recognizes  $\varphi$  and the product transition system  $TS_p = TS \otimes \mathcal{A}$  are shown in Figure 4.3 and Figure 4.4, respectively. In relation to Figure 4.1, projecting the path that is shown in Figure 4.4 onto the state of TS yields the same transition system that applying the control protocol in Example 4.1 induces.

## 4.3 A Case Study

We now take a closer look into a case study motivated by autonomous driving. Consider a scenario where an autonomous vehicle  $TS^a$  needs to make an unprotected left turn with an oncoming vehicle  $TS^h$  as shown in Figure 4.5. The complete system consists of the autonomous vehicle, the oncoming vehicle, and the traffic light.

First, consider the case where the complete system is deterministic, i.e., starting from any state, the actions of the autonomous vehicle lead to a unique state of the complete system. This implies that the behavior of the oncoming vehicle and the traffic light is deterministic. Figure

45



**Figure 4.4:** The product transition system  $TS_p = (T_1||T_2) \otimes A$ , showing only reachable states. An accepting path is highlighted by double (red) arrows.

4.6 shows the finite transition systems that describe the behavior of each element as well as the complete system TS. Here, a state  $s \in S$  is of the form  $s = (c_i, c_j, s_k)$ , where i and j are the labels of the cells occupied by the autonomous vehicle and the oncoming vehicle, respectively, and k indicates whether or not the light is green; k = 1 if the light is green and k = 2 if the light is red.  $Act = \{acc, brake\}$  where acc and brake represent accelerating and braking, respectively.  $AP = \{a_0, \ldots, a_9, h_0, \ldots, h_9, g\}$  and the labeling function L is defined such that for any state  $s = (c_i, c_j, s_k)$ ,  $a_i, h_j \in L(s)$  and  $g \in L(s)$  if and only if k = 1.

An LTL formula  $\varphi = \neg (h_4 \land a_4) U a_9$  specifies that the autonomous vehicle and the oncoming vehicle do not simultaneously occupy cell  $c_4$  until the autonomous vehicle reaches cell  $c_9$ . Figure 4.7 shows the corresponding nondeterministic Büchi automaton  $\mathcal{A}$ . The product  $TS \otimes \mathcal{A}$  as well as a solution is shown in Figure 4.8.

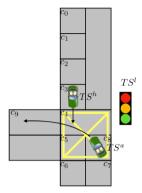


Figure 4.5: An unprotected left turn scenario.

#### 4.4 Minimum-Violation Synthesis

Autonomous systems are often subject to multiple regulatory requirements or safety rules that may not be equally important. In general, it is infeasible to guarantee the satisfaction of all the rules under all conditions. Thus, there is a need to allow for formally justifiable violation of these rules. In particular, we assume that each rule has a certain penalty associated with its violation. The goal of *minimum-violation synthesis* is to minimize such penalties.

The formulation of the minimum-violation synthesis requires some extension to the models and specifications considered in Section 2. First, we consider a system that can be modeled as a weighted, finite, and action-deterministic transition system.

**Definition 4.2.** A weighted transition system (WTS) is a tuple  $TS = (S, Act, \rightarrow, I, \mathcal{AP}, L, W)$  where  $S, Act, \rightarrow, I, AP$ , and L are the same as in Definition 2.1 and, for some fixed  $m \in \mathbb{N}$ ,  $W: S \times S \rightarrow \mathbb{R}^m_{\geq 0}$  is a weight function.

For states  $s_1, s_2 \in S$ , the weight  $W(s_1, s_2)$  typically represents the time or distance between  $s_1$  and  $s_2$ . Similar to conventional transition systems, a WTS is *finite* if S, Act and AP are all finite, and is action-deterministic

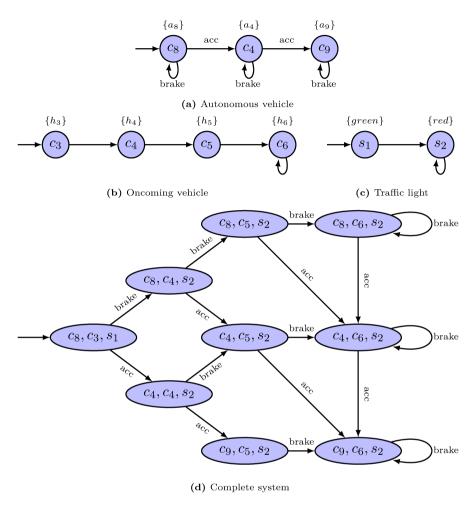


Figure 4.6: The finite transition systems that describe the behavior of the autonomous vehicle, the oncoming vehicle, the traffic light, and the complete system.

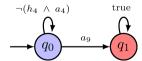
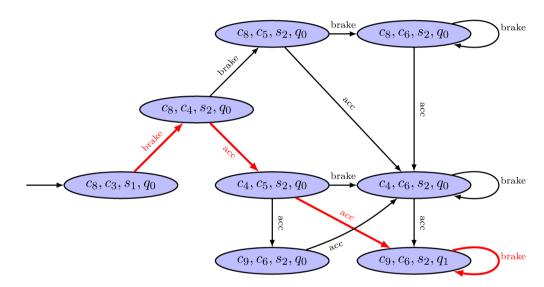


Figure 4.7: The nondeterministic Büchi automaton corresponding to LTL formula  $\neg (h_4 \land a_4) \mathcal{U} a_9$ .



**Figure 4.8:** The product  $TS \otimes \mathcal{A}$  as well as a run on TS that satisfies  $\varphi$  shown in green.

if  $|I| \leq 1$  and, for all  $s \in S$  and  $\alpha \in Act$ ,  $|Post(s, \alpha)| \leq 1$ . The weight of a finite path fragment  $\pi = s_0 s_1 \dots s_n$  is  $W(\pi) = \sum_{i=0}^{n-1} W(s_i, s_{i+1})$ .

We will use finite linear temporal logic (FLTL) to formalize each rule. As opposed to LTL, an FLTL formula  $\varphi$  is interpreted over a finite word  $\sigma = \sigma_0 \sigma_1 \dots \sigma_n \in (2^{\mathcal{AP}})^{n+1}$ . We write  $\sigma \models \varphi$  if and only if  $\sigma$  satisfies  $\varphi$ . For example, for some  $p \in \mathcal{AP}$ ,  $\sigma \models p$  if and only if  $p \in \sigma_0$ . Additionally,  $\sigma \models \Box p$  if and only if, for all  $i \in \{0, \dots, n\}$ , we have  $p \in \sigma_i$ . As a more complicated example, with  $p, p' \in \mathcal{AP}$ , let  $\varphi = \Box(p \Longrightarrow (\bigcirc p \vee \bigcirc p'))$ . Then,  $\sigma \models \varphi$  if and only if, for all  $i \in \{0, \dots, n-1\}$  such that  $p \in \sigma_i$ , we have  $p \in \sigma_{i+1}$  or  $p' \in \sigma_{i+1}$ .

Given an FLTL formula  $\varphi$ , a finite automaton  $\mathcal{A} = (Q, Q_0, \delta, F)$  that accepts all and only finite words that satisfy  $\varphi$  can be automatically constructed (Gunter and Peled, 2002). Here, the definitions of Q,  $Q_0$ ,  $\delta$  and F are similar to those in the definition of a Büchi automaton (see Section 2.2). However, a finite automaton is evaluated over a finite word  $\sigma = \sigma_0 \dots \sigma_n$ . A run of  $\mathcal{A}$  over  $\sigma$  is defined as a finite sequence  $q_0 \dots q_{n+1}$  such that  $q_0 \in Q_0$  and  $(q_i, \sigma_i, q_{i+1}) \in \delta$  for all  $i \in \{0, \dots, n\}$ . We say that  $\mathcal{A}$  accepts  $\sigma$  if and only if there exists a run  $q_0 \dots q_{n+1}$  of  $\mathcal{A}$  over  $\sigma$  such that  $q_{n+1} \in F$ .

We now introduce the concept of prioritized safety specification to formalize a set of rules with unequal importance.

**Definition 4.3.** A prioritized safety specification is a tuple  $\mathcal{P} = (\mathcal{AP}, \Phi, \Psi, \rho)$  where  $\mathcal{AP}$  is a set of atomic propositions,  $\Phi$  is a set of FLTL formulas over  $\mathcal{AP}$ ,  $\Psi = (\Psi_0, \Psi_1, \dots, \Psi_N)$  organizes the formulas in  $\Phi$  into a hierarchy based on their priorities such that  $\Psi_i \subseteq \Phi$ , for all  $i \in \mathbb{N}_{\leq N}$ , and  $\rho : \Phi \to \mathbb{N}$  is a function that assigns the weight to each  $\varphi \in \Phi$ . Throughout the article, we refer to each  $\varphi \in \Phi$  as an atomic safety rule.

**Example 4.5.** Consider an autonomous vehicle that is required to (1) avoid collision, (2) stay on the road, (3) keep sufficient clearance from other vehicles, and (4) stay within a correct lane. To capture these rules, we define the following atomic propositions.

- collision represents a state at which the autonomous vehicle collides with another vehicle or obstacle.
- onroad represents a state at which the autonomous vehicle is fully on the road.
- close represents a state at which the autonomous vehicle overlaps with the clearance zone around another vehicle.
- inlane represents a state at which the autonomous vehicle is fully within a correct lane.

Following Example 2.6, we let  $(x, \theta, v) \in \mathbb{R}^3$  represent the state of the autonomous vehicle,  $\mathsf{FP}(x, \theta) \subset \mathbb{R}^2$  represent its footprint, and  $\mathsf{Obs} \subset \mathbb{R}^2$  represent the union of the footprints of all the vehicles and obstacles in the environment. Additionally, we define the following environment states.

- $\mathsf{RD} \subset \mathbb{R}^2$  is the road, i.e., the area where a vehicle is allowed to drive, and
- $\mathsf{CZ} \subset \mathbb{R}^2$  is the union of the clearance zone around other vehicles constructed from their footprints and the required lateral and longitudinal clearance,
- LN  $\subset \mathbb{R}^2$  is the right lane, i.e., the lane with the correct travel direction for the autonomous vehicle.

The state of the complete system (autonomous vehicle and the environment) is then given by  $(x, \theta, v, \mathsf{Obs}, \mathsf{RD}, \mathsf{CZ}, \mathsf{LN})$ . The labeling function L is defined such that for any  $x, \theta, v, \mathsf{Obs}, \mathsf{RD}, \mathsf{CZ}, \mathsf{LN},$ 

- $\bullet \ \ \operatorname{collision} \in L(x,\theta,v,\operatorname{Obs},\operatorname{RD},\operatorname{CZ},\operatorname{LN}) \ \operatorname{iff} \ \operatorname{FP}(x,\theta) \ \cap \ \operatorname{Obs} \neq \emptyset,$
- onroad  $\in L(x, \theta, v, \mathsf{Obs}, \mathsf{RD}, \mathsf{CZ}, \mathsf{LN})$  iff  $\mathsf{FP}(x, y, \theta) \subseteq \mathsf{RD},$
- close  $\in L(x, \theta, v, \mathsf{Obs}, \mathsf{RD}, \mathsf{CZ}, \mathsf{LN})$  iff  $\mathsf{FP}(x, y, \theta) \, \cap \, \mathsf{CZ} \neq \emptyset$ , and
- inlane  $\in L(x, \theta, v, \mathsf{Obs}, \mathsf{RD}, \mathsf{CZ}, \mathsf{LN})$  iff  $\mathsf{FP}(x, y, \theta) \subseteq \mathsf{LN}.$

We consider the following atomic safety rules, each of which can be expressed by an FLTL formula.

- (i) Avoiding collision:  $\varphi_0 = \Box \neg \text{collision}$ .
- (ii) Staying on the road:  $\varphi_1 = \square$  onroad.
- (iii) Keeping sufficient clearance from other vehicles:  $\varphi_2 = \Box \neg \mathsf{close}$ .
- (iv) Staying within a correct lane:  $\varphi_3 = \Box$ inlane.

Define the prioritized safety specification  $\mathcal{P} = (\mathcal{AP}, \Phi, \Psi, \rho)$  as

$$\begin{split} \mathcal{AP} &= \{ \text{collision}, \text{close}, \text{onroad}, \text{inlane} \}, \\ \Phi &= \{ \varphi_0, \dots, \varphi_3 \}, \\ \Psi &= \{ \Psi_0, \Psi_1, \Psi_2 \}, \text{ and } \\ \rho(\varphi_0) &= 1, \forall i, \end{split}$$

where,  $\Psi_0 = \{\varphi_0\}$ ,  $\Psi_1 = \{\varphi_1\}$ , and  $\Psi_2 = \{\varphi_2, \varphi_3\}$ . In this definition of  $\mathcal{P}$ ,  $\varphi_0$  is the only rule at the top level in the hierarchy. As we will see later, this means that  $\varphi_0$  is the most important rule, and the system should minimize the violation of this rule, even at the cost of infinitely violating the other rules. Next,  $\varphi_1$  is the only rule at the second level. This means that after minimizing  $\varphi_0$ , the system will minimize the violation of  $\varphi_1$ , even at the cost of infinitely violating the other rules at the lower levels. Finally,  $\varphi_2$  and  $\varphi_3$  are at the lowest level. This means that after minimizing the violation of  $\varphi_0$  and  $\varphi_1$ , the system then tries to minimize the weighted violation of  $\varphi_2$  and  $\varphi_3$ , where the weights are given by  $\rho$ . Note that the weights only affect the rules at the same level.

We use the level of unsafety  $\lambda_{\varphi}(\pi)$  to measure the violation of a finite path fragment  $\pi$  with respect to an atomic safety rule  $\varphi$ . The dual concept, including robustness, has been proposed for some variants of temporal logics such as Signal Temporal Logic (Donzé and Maler, 2010; Mehdipour et al., 2019). However, when applying to a system that is subject to multiple requirements, this measure introduces the additional complexity of having to ensure that the system cannot take advantage of robustly satisfying one rule, in order to slightly violate another rule. Additionally, when combining with the concept of rule hierarchy, this measure may cause the system to compromise lower-level rules by robustly satisfying the higher-level rules. For example,

consider the prioritized safety specification in Example 4.5. Using the robustness as the performance measure may incentivize the system to violate staying on the road rule  $\varphi_2$  in order to robustly satisfy the collision avoidance rule  $\varphi_1$  since this behavior allows the system to be as far away from other vehicles as possible, and hence, maximizing the robustness of the top-level rule.

Various definitions of the level of unsafety have been proposed. For example, Tumova et al. (2013) defines  $\lambda_{\varphi}(\pi)$  as the minimum number of states in  $\pi$  that needs to be removed so that the resulting path fragment satisfies  $\varphi$ . Formally, given a finite sequence  $\pi = s_0 s_1 \dots s_n$  and a set  $I \subseteq \{0, \dots, n\}$ , let vanish $(\pi, I)$  represent a subsequence of  $\pi$  obtained by removing all  $s_i$  with  $i \in I$ . Then,  $\lambda_{\varphi}(\pi) = \min_I \{|I| \text{ s.t. vanish}(\pi, I) \models \pi\}$ . In contrast, (Castro et al., 2013) defines the level of unsafety as  $\lambda_{\varphi}(\pi) = \min_I \{\sum_{i \in I \setminus \{n\}} W(s_i, s_{i+1}) \text{ s.t. vanish}(\pi, I) \models \pi\}$ . Following Wongpiromsarn et al. (2021), we restrict atomic safety rules to be expressed using a class of FLTL known as si-FLTL<sub>Gx</sub> defined as follows:

**Definition 4.4.** An si-FLTL<sub>Gx</sub> formula over a set  $\mathcal{AP}$  of atomic propositions is an FLTL formula that is stutter-invariant (see below) and is of the form

$$\varphi = \Box P_{\mathsf{X}},$$

where  $P_X$  belongs to the smallest set defined inductively by the following rules:

- p is a formula for all  $p \in \mathcal{AP} \cup \{True, False\};$
- $\bigcirc p$  is a formula for all  $p \in \mathcal{AP} \cup \{\mathit{True}, \mathit{False}\}$ ; and
- if  $P_X^1$  and  $P_X^2$  are formulas, then so are  $\neg P_X^1$ ,  $P_X^1 \lor P_X^2$ ,  $P_X^1 \land P_X^2$  and  $P_X^1 \Longrightarrow P_X^2$ .

In other words,  $P_X$  is a Boolean combination of propositions from  $\mathcal{AP}$  and expressions of the form  $\bigcirc p$  where  $p \in \mathcal{AP}$ .

Roughly, a specification is stutter-invariant if its satisfaction with respect to any word is not affected by operations that duplicate some letters or remove some duplicate letters in that word. For example, consider  $\sigma =$ 

53

 $\sigma_0\sigma_1\ldots\sigma_n$  and  $\sigma'=\sigma_0\sigma_1\ldots\sigma_{i-1}\sigma_i\sigma_i\sigma_{i+1}\ldots\sigma_n$ , which is constructed from  $\sigma$  by duplicating  $\sigma_i$  for some  $i\in\{0,\ldots,n\}$ . If  $\varphi$  is stutter-invariant, then  $\sigma\models\varphi$  if and only if  $\sigma'\models\varphi$ . We refer the reader to (Peled and Wilke, 1997) for the definition of stutter-invariant properties. See, e.g., (Klein and Baier, 2007; Michaud and Duret-Lutz, 2015) for approaches to check whether a specification is stutter-invariant.

**Example 4.6.** All the rules  $\varphi_0, \ldots, \varphi_3$  defined in Example 4.5 are si-FLTL<sub>GX</sub> formulas. For example, consider  $\varphi_0 = \Box \neg \text{collision}$ . It can be written as  $\varphi_0 = \Box P_X$ , where  $P_X = \neg \text{collision}$ .

Regardless of their simplicity, si-FLTL $_{\mathsf{G}_\mathsf{X}}$  formulas turn out to be sufficiently expressive in many applications; for example, (Wongpiromsarn et~al.,~2011a) shows that all the rules enforced in the DARPA Urban Challenge 2007 can be expressed with si-FLTL $_{\mathsf{G}_\mathsf{X}}$  formulas. Moreover, all of the traffic rules in the examples presented in (Castro et~al.,~2013) can be described using si-FLTL $_{\mathsf{G}_\mathsf{X}}$  formulas.

Wongpiromsarn et al. (2021) show that the violation of a si-FLTL $_{\mathsf{G}_{\mathsf{X}}}$  formula is caused either by visiting an unsafe state or by taking an unsafe transition. As a result, this work defines the level of unsafety as the total time spent in an unsafe state and the total number of unsafe transitions.

Let  $\mathcal{P} = (\mathcal{AP}, \Phi, \Psi, \rho)$  be a prioritized safety specification where  $\Psi = (\Psi_0, \Psi_1, \dots, \Psi_N)$ . We define the level of unsafety of a finite-path fragment  $\pi$  with respect to  $\mathcal{P}$  as

$$\lambda_{\mathcal{P}}(\pi) = (\lambda_{\Psi_0}(\pi), \dots, \lambda_{\Psi_N}(\pi)) \in \mathbb{R}^{N+1},$$

where, for every  $i \in \{0, 1, \dots, N\}$ ,

$$\lambda_{\Psi_i}(\pi) = \sum_{\varphi \in \Psi_i} \rho(\varphi) \lambda_{\varphi}(\pi).$$

Let  $Path_G(TS)$  be the set of finite path fragments that end in a goal state. Given a weighted, finite, and action-deterministic transition system  $TS = (S, Act, \rightarrow, I, \mathcal{AP}, L, W)$  and a prioritized safety specification

 $\mathcal{P} = (\mathcal{AP}, \Phi, \Psi, \rho)$ , the minimum-violation synthesis problem is to compute an optimal path fragment  $\pi^* \in Path_G(TS)$  that minimizes the weight  $W(\pi^*)$  among all the path fragments that minimize the level of unsafety with respect to  $\mathcal{P}$ . Formally, we define the cost function  $J: Path_G(TS) \to \mathbb{R}^{N+2}$  as

$$J(\pi) = (\lambda_{\mathcal{P}}(\pi), W(\pi)). \tag{4.1}$$

Using the cost function J, we now formally define the minimum-violation synthesis problem.

Minimum-Violation Synthesis Problem: Based on the standard lexicographical ordering, compute an optimal finite path fragment  $\pi^*$  such that

$$\pi^* = \arg\min_{\pi \in Path_G(TS)} J(\pi).$$

**Example 4.7.** Consider the prioritized safety specification  $\mathcal{P}$  in Example 4.5. In this case, we have  $\lambda_{\mathcal{P}}(\pi) = (\lambda_{\Psi_0}(\pi), \lambda_{\Psi_1}(\pi), \lambda_{\Psi_2}(\pi))$ , where

$$\begin{split} &\lambda_{\Psi_0}(\pi) = \lambda_{\varphi_0}(\pi), \\ &\lambda_{\Psi_1}(\pi) = \lambda_{\varphi_1}(\pi), \text{ and} \\ &\lambda_{\Psi_2}(\pi) = \lambda_{\varphi_2}(\pi) + \lambda_{\varphi_3}(\pi). \end{split}$$

As a result, the cost function  $J: Path_G(TS) \to \mathbb{R}^4$  is defined as

$$J(\pi) = (\lambda_{\Psi_0}(\pi), \lambda_{\Psi_1}(\pi), \lambda_{\Psi_2}(\pi), W(\pi)).$$

For convenience, we define the following subsets.

- $Path_{G,0}(TS)$  is the subset of  $Path_G(TS)$  that minimize the level of unsafety with respect to  $\varphi_0$ .
- $Path_{G,0,1}$  is the subset of  $Path_{G,0}(TS)$  that minimize the level of unsafety with respect to  $\varphi_1$ .
- $Path_{G,0,1,2}$  is the subset of  $Path_{G,0,1}(TS)$  that minimize the level of unsafety with respect to  $\varphi_2$  and  $\varphi_3$ .

An optimal path  $\pi^*$  is defined as a path in  $Path_{G,0,1,2}$  that minimizes the weight  $W(\pi^*)$ , which typically represents the time or distance on the path fragment.

Tumova et al. (2013) and Castro et al. (2013) solve the minimum-violation synthesis problem by constructing a weighted finite automaton  $\mathcal{A}$  that is the product of a collection of weighted finite automata, each of which corresponds to an atomic safety rule  $\varphi \in \Phi$ . A weighted finite automaton is essentially a finite automaton with weight  $W(q, \sigma, q') \in \mathbb{R}^m_{\geq 0}$  for some  $m \in \mathbb{N}$  assigned to each transition  $(q, \sigma, q')$ . The weights of the transitions of  $\mathcal{A}$  are defined such that the weight of the shortest accepting run over any word  $\sigma$  is the level of unsafety of  $\sigma$ . It can be shown that the minimum-violation synthesis problem is equivalent to finding the shortest path in  $TS \otimes \mathcal{A}$ .

As the size of  $\mathcal{A}$  is exponential in the length of  $\varphi$  (Baier and Katoen, 2008), the approach presented in (Wongpiromsarn et~al., 2021) avoids constructing the product  $TS \otimes \mathcal{A}$  to reduce computational complexity. The main idea is to construct a weighted finite transition system TS' with the same sets S, Act, I, and AP of states, actions, transitions, initial states, and atomic propositions as well as the same transition relation  $\to$  and labeling function L as TS. Instead, the weight W' will be defined such that the minimum-violation synthesis problem can be translated to computing the shortest path on TS'.

To enable such a construction of TS', we first translate a si-FLTL<sub>Gx</sub> formula over  $\mathcal{AP}$  into a si-FLTL<sub>G</sub> formula over  $\mathcal{AP} \times \mathcal{AP}$ .

**Definition 4.5** (si-FLTL<sub>G</sub>). A si-FLTL<sub>G</sub> formula over  $\mathcal{AP} \times \mathcal{AP}$  is a si-FLTL<sub>Gx</sub> formula  $\varphi = \Box P$  where P is a propositional logic formula over  $\mathcal{AP} \times \mathcal{AP}$ .

A propositional logic formula P over  $\mathcal{AP} \times \mathcal{AP}$  is interpreted over a pair  $(l, l') \in 2^{\mathcal{AP}} \times 2^{\mathcal{AP}}$  with the satisfaction relation  $\models$  defined as follows: for  $p, p' \in \mathcal{AP} \cup \{True, False\}$  and  $(l, l') \in 2^{\mathcal{AP}} \times 2^{\mathcal{AP}}$ ,  $(l, l') \models (p, p')$  if and only if  $l \models p$  and  $l' \models p'$ . Here, for any  $l \in 2^{\mathcal{AP}}$ , we have  $l \models True$ ,  $l \not\models False$ , and for any  $p \in \mathcal{AP}$ ,  $l \models p$  if and only if  $p \in l$ . The logic connectives are defined as in the standard propositional logic.

Based on the semantics of FLTL, given a finite word  $\sigma = \sigma_0 \sigma_1 \dots \sigma_n \in (2^{\mathcal{AP}})^{n+1}$  and a si-FLTL<sub>G</sub> formula  $\varphi = \Box P$  over  $\mathcal{AP} \times \mathcal{AP}$ , we say that

 $\sigma$  satisfies  $\varphi$ , written  $\sigma \models_{\mathcal{AP} \times \mathcal{AP}} \varphi$  if and only if  $(\sigma_i, \sigma_{i+1}) \models P$  for all  $i \in \mathbb{N}_{\leq n-1}$  and  $(\sigma_n, \sigma_n) \models P$ . Note that the terminal condition  $(\sigma_n, \sigma_n) \models P$  results from the assumption that  $\varphi$  is stutter-invariant, which ensures that  $\sigma \models \varphi$  if and only if  $\sigma' = \sigma_0 \sigma_1 \dots \sigma_n \sigma_n \models \varphi$ .

Given a si-FLTL<sub>Gx</sub> formula  $\varphi$  over  $\mathcal{AP}$ , we define an operation called denext which constructs a si-FLTL<sub>G</sub> formula over  $\mathcal{AP} \times \mathcal{AP}$  from  $\varphi$  by replacing each instance of p in  $\varphi$  with (p, True) and replacing each instance of p in  $\varphi$  with (True, p), for all  $p \in \mathcal{AP}$ . For example, consider a si-FLTL<sub>Gx</sub> formula  $\varphi = \square(p \implies (\bigcirc p \vee \bigcirc p'))$ . The corresponding si-FLTL<sub>G</sub> formula over  $\mathcal{AP} \times \mathcal{AP}$  is given by

$$\mathsf{denext}(\varphi) = \Box \Big( (p, \mathit{True}) \implies \big( (\mathit{True}, p) \lor (\mathit{True}, p') \big) \Big).$$

**Example 4.8.** Consider Example 4.5. si-FLTL<sub>GX</sub> formulas  $\varphi_0, \ldots, \varphi_4$  over  $\mathcal{AP}$  can be translated to the corresponding si-FLTL<sub>G</sub> formulas over  $\mathcal{AP} \times \mathcal{AP}$  as follows.

- (i) Avoiding collision: denext( $\varphi_0$ ) =  $\Box \neg$ (collision, True).
- (ii) Staying on the road: denext( $\varphi_1$ ) =  $\square$ (onroad, True).
- (iii) Keeping sufficient clearance from other vehicles:  $denext(\varphi_2) = \Box \neg (close, True)$ .
- (iv) Staying within a correct lane: denext( $\varphi_3$ ) =  $\square$ (inlane, True).

Wongpiromsarn et al. (2021) establish the equivalence of the level of unsafety with respect to a si-FLTL<sub>GX</sub> formula over  $\mathcal{AP}$  and the level of unsafety with respect to the corresponding si-FLTL<sub>G</sub> formula over  $\mathcal{AP} \times \mathcal{AP}$ .

The translation of si-FLTL<sub>Gx</sub> formula over  $\mathcal{AP}$  to a si-FLTL<sub>G</sub> formula over  $\mathcal{AP} \times \mathcal{AP}$  allows us to convert the original prioritized safety specification  $\mathcal{P} = (\mathcal{AP}, \Phi, \Psi, \rho)$  to  $\hat{\mathcal{P}} = (\mathcal{AP} \times \mathcal{AP}, \hat{\Phi}, \hat{\Psi}, \hat{\rho})$  with each atomic safety rule obtained from that of  $\mathcal{P}$  by applying denext operation. Formally,  $\hat{\Phi} = \{\text{denext}(\varphi) \mid \varphi \in \Phi\}; \hat{\Psi} = (\hat{\Psi}_0; \hat{\Psi}_1, \dots, \hat{\Psi}_N); \hat{\Psi}_i = \{\text{denext}(\varphi) \mid \varphi \in \Psi_i\}, \text{ for all } i \in \{0, \dots, N\}; \text{ and } \hat{\rho}(\text{denext}(\varphi)) = \rho(\varphi), \text{ for all } \varphi \in \Phi.$ 

For each atomic safety rule  $\varphi \in \Phi$ , define a propositional logic formula  $P_{\varphi}$  such that  $\operatorname{denext}(\varphi) = \Box P_{\varphi}$ . We construct the weighted finite transition system  $TS' = (S, Act, \to, I, \mathcal{AP}, L, W')$  such that the weight W' corresponds to the cost function in (4.1). Formally,  $W' : S \times S \to \mathbb{R}^{N+2}$  is defined by

$$W'(s_1, s_2) = (\lambda_{\hat{\Psi}_0}(s_1, s_2), \dots, \lambda_{\hat{\Psi}_N}(s_1, s_2), W(s_1, s_2)),$$

where  $\lambda_{\hat{\Psi}_i}(s_1, s_2) = \sum_{\varphi \in \hat{\Psi}_i} \hat{\rho}(\varphi) \lambda_{\varphi}(s_1, s_2)$  and  $\lambda_{\varphi}(s_1, s_2)$  is defined based on the satisfaction of the propositional logic formula  $P_{\varphi}$  at  $s_1$  and  $s_2$  as follows:

$$\lambda_{\varphi}(s_1, s_2) = \left\{ \begin{array}{ll} 0 & \text{if } (L(s_1), L(s_2)) \models P_{\varphi} \\ W(s_1, s_2) & \text{if } (L(s_1), l) \not\models P_{\varphi} \text{ for all } l \in 2^{\mathcal{AP}} \\ 1 & \text{otherwise} \end{array} \right\}$$
(4.2)

Note that the first condition of (4.2) corresponds to the case where  $P_{\varphi}$  is satisfied by the transition from  $s_1$  to  $s_2$ , thereby no violation cost being incurred. The second condition corresponds to visiting an unsafe state  $s_1$  and the third condition corresponds to taking an unsafe transition.

**Example 4.9.** Let us revisit Example 4.5 with the corresponding si-FLTL<sub>G</sub> formulas over  $\mathcal{AP} \times \mathcal{AP}$  defined in Example 4.8. For any states  $s_1, s_2 \in S$ , the weight W' is defined as

$$W'(s_1, s_2) = (\lambda_{\hat{\Psi}_0}(s_1, s_2), \lambda_{\hat{\Psi}_1}(s_1, s_2), \lambda_{\hat{\Psi}_2}(s_1, s_2), W(s_1, s_2)),$$

where

- (i)  $\lambda_{\hat{\Psi}_0}(s_1, s_2)$  corresponds to the violation of the avoiding collision rule and is defined as  $\lambda_{\hat{\Psi}_0}(s_1, s_2) = 0$  if collision  $\notin L(s_1)$  (i.e., the autonomous vehicle does not collide with another vehicle or obstacle at state  $s_1$ ) and  $\lambda_{\hat{\Psi}_0}(s_1, s_2) = W(s_1, s_2)$  otherwise,
- (ii)  $\lambda_{\hat{\Psi}_1}(s_1, s_2)$  corresponds to the violation of the staying on the road rule and is defined as  $\lambda_{\hat{\Psi}_1}(s_1, s_2) = 0$  if road  $\in L(s_1)$  (i.e., the vehicle is fully on the road at state  $s_1$ ) and  $\lambda_{\hat{\Psi}_0}(s_1, s_2) = W(s_1, s_2)$  otherwise, and

(iii) 
$$\lambda_{\hat{\Psi}_2}(s_1, s_2) = \lambda_{\varphi_2}(s_1, s_2) + \lambda_{\varphi_3}(s_1, s_2)$$
, where

- $\lambda_{\varphi_2}(s_1, s_2)$  corresponds to the violation of the clearance rule and is defined as  $\lambda_{\varphi_2}(s_1, s_2) = 0$  if close  $\not\in L(s_1)$  (i.e., the autonomous vehicle does not overlap with the clearance zone around another vehicle at state  $s_1$ ) and  $\lambda_{\varphi_2}(s_1, s_2) = W(s_1, s_2)$  otherwise, and
- $\lambda_{\varphi_3}(s_1, s_2)$  corresponds to the violation of the lane rule and is defined as  $\lambda_{\varphi_3}(s_1, s_2) = 0$  if inlane  $\in L(s_1)$  (i.e., the autonomous vehicle is fully within a correct lane at state  $s_1$ ) and  $\lambda_{\varphi_3}(s_1, s_2) = W(s_1, s_2)$  otherwise.

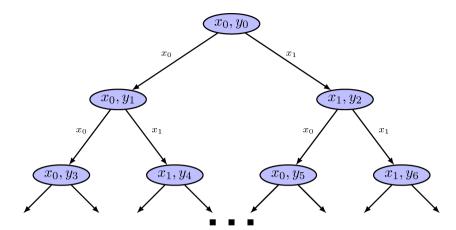
Wongpiromsarn et al. (2021) show that the minimum-violation synthesis problem is equivalent to computing the shortest path (based on the standard lexicographical ordering) on TS', which has the same size as TS. As a result, the construction of TS' allows temporal logic specifications to be handled with the same computational complexity as traditional graph-search algorithms such as Dijkstra and  $A^*$ .

# 5

## Reactive Synthesis

This chapter continues the discussion of correct-by-construction synthesis of control protocols focusing on non-deterministic systems. In particular, we consider open systems whose behaviors can be affected by exogenous inputs. Non-determinism can be used to capture uncertainties in the system, and is particularly useful in capturing uncertainties arising from valid environment behaviors that the system cannot control. Such a system is called reactive as it must react to the environment's behavior. Similarly to the previous chapter, we focus on discrete systems modeled as finite transition systems and we assume that the system's state is observable at all times. This chapter puts emphasis on methods that alleviate some of the difficulties – such as computational complexity and conflicting specifications – that naturally arise when constructing autonomous protocols.

This section incorporates the results from the following publications (Wongpiromsarn, Topcu, and Murray, 2013; Kress-Gazit, Wongpiromsarn, and Topcu, 2011; Dimitrova, Ghasemi, and Topcu, 2018).



**Figure 5.1:** A tree representing a control protocol for non-deterministic systems. The state of the system is a tuple (x, y) where  $x \in \{x_0, x_1\}$  represents the non-deterministic choice of the environment behavior whereas  $y \in \{y_0, y_1, \ldots\}$  represents the state over which the system has control.

#### 5.1 Synthesis of Reactive Control Protocol

Similarly to as for closed systems, we say that a reactive system is correct with respect to specification  $\varphi$  if it satisfies  $\varphi$ . However, for a non-deterministic system, the correctness needs to be interpreted with respect to the non-deterministic choices over which the system does not have control. In this case, we require that the control protocol ensures that specification  $\varphi$  is satisfied for all possible non-deterministic choices, for example, for all of the environment's possible behaviors. As discussed in Pnueli and Rosner (1989) and Piterman et al. (2006), the control protocol synthesis in this case can be treated as a two-player game between the system and the environment, also called the adversary. The system and the environment alternate in picking actions. The environment then attempts to falsify  $\varphi$  while the system attempts to satisfy  $\varphi$ . A provably correct control protocol therefore needs to ensure that  $\varphi$  is satisfied for any possible behavior from the environment. The control protocol may thus be represented by a tree whose branches represent the possible environment actions and capture non-determinism as shown in Figure 5.1.

Solving the above two-player game typically involves computing the winning set, which is defined as the set of initial states from which there exists a strategy for the system to satisfy the specification, for all the possible environment behaviors. Similar to model-checking-based policy synthesis, computing a winning set requires computing the product of the transition system and a finite automata, except for that a non-deterministic Büchi automaton  $\mathcal{A}_{\varphi}$  recognizing  $\varphi$  must be transformed into a deterministic Rabin automaton  $\mathcal{R}_{\varphi}$  (Baier and Katoen, 2008). We call such a transformation determinization. We then compute the product transition system  $TS_p = TS \otimes \mathcal{R}_{\varphi}$  which is defined similarly to as in Definition 3.1. A fixed-point strategy can be applied to  $TS_p$  in order to derive the winning set. Finally, a control protocol can be constructed using the intermediate values in the computation of the winning set.

The size of  $TS_p$  and the runtime of this synthesis algorithm are both at most double exponential in the length of  $\varphi$ . The first exponent results from the construction of the non-deterministic Büchi automaton  $\mathcal{A}_{\varphi}$  from  $\varphi$  and the second exponent results from the determinization of  $\mathcal{A}_{\varphi}$  into a deterministic Rabin automaton  $\mathcal{R}_{\varphi}$ . We refer the reader to (Pnueli and Rosner, 1989; Kloetzer and Belta, 2008b) for more details.

For a class of specifications of the form  $\Box p, \Diamond p, \Box \Diamond p$  and  $\Diamond \Box p$ , where p is a proposition, there exist efficient synthesis algorithms (Asarin  $et\ al.$ , 1998; Alur and La Torre, 2004). The main idea behind these algorithms is that they avoid translating the specification to a non-deterministic Büchi automaton as well as avoiding the determinization of the non-deterministic Büchi automaton into a deterministic Rabin automaton. For example, in a reachability game with specification  $\varphi = \Diamond p$ , we define the set  $W = \{s \in S : s \models p\}$  to be the set of states at which p is satisfied, and we define the predecessor operator  $Pre_{\exists \forall} \mid 2^S \to 2^S$  as follows:  $Pre_{\exists \forall}(R)$  is the set of states whose subsequent successors have at least one successor in R. In particular,  $Pre_{\exists \forall}(R) = \{s \in S \mid \forall s' \in S, s \to s' \text{ implies } \exists s'' \in R \text{ such that } s' \to s''\}$ . The set of states from which the controller can lead the system into W can be computed

efficiently by the iteration sequence

$$R_0 = W,$$

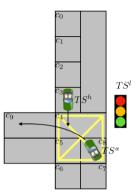
$$R_i = R_{i-1} \cup Pre_{\exists \forall}(R_{i-1}), \forall i > 0.$$

Using the Tarski-Knaster Theorem, it can be shown that there exists a natural number n such that  $R_n = R_{n-1}$ . In addition,  $R_n$  is the minimal solution of the fix-point equation  $R = W \cup Pre_{\exists \forall}(R)$ .

The methodology proposed by Piterman *et al.* (2006) allows for solving a broader class of games efficiently. A summary of the algorithm is as follows: First, a *game structure* is defined as a tuple  $\mathfrak{G} = (V, X, Y, \theta_e, \theta_s, \rho_e, \rho_s, AP, L, \varphi)$  where

- V is a finite set of variables over finite domains,
- $X \subseteq V$  is a set of environment variables,
- $Y = V \setminus X$  is a set of controlled variables,
- $\theta_e(X)$  is a proposition over X characterizing the initial states of the environment,
- $\theta_s(V)$  is a proposition over V characterizing the initial states of the system,
- $\rho_e(V, X')$  is a proposition that relates a state  $s \in dom(V)$  to a possible next input value  $s_X \in dom(X)$  and characterizes the transition relation of the environment,
- $\rho_s(V, X', Y')$  is a proposition that relates a state  $s \in dom(V)$  and an input value  $s_X \in dom(X)$  to an output value  $s_Y \in dom(Y)$  and characterizes the transition relation of the system,
- AP is a set of atomic propositions,
- $L: dom(V) \to 2^{AP}$  is a labeling function, and
- $\varphi$  is the winning condition, characterized by a linear temporal logic (LTL) formula.

We let dom(V), dom(X) and dom(Y) denote the set of all the possible assignments to variables in V, X and Y, respectively. An environment



**Figure 5.2:** An unprotected left turn scenario. As described in Section 4.3, the autonomous vehicle  $TS^a$  needs to make an unprotected left turn with an oncoming vehicle  $TS^h$ .  $TS^l$  denotes the traffic light.

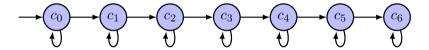


Figure 5.3: The finite transition system that represents the nondeterministic model of the oncoming vehicle.

state  $s_X \in dom(X)$  is a valid input in state  $s \in dom(V)$  if  $(s, s_X) \models \rho_e$ . Analogously, a controlled state  $s_Y \in dom(Y)$  is a valid system output at state  $s \in dom(V)$ , after reading input  $s_X$ , if  $(s, s_X, s_Y) \models \rho_s$ .

**Example 5.1.** Let us revisit the unprotected left turn scenario described in Section 4.3 and illustrated in Figure 5.2. In practice, the behavior of the oncoming vehicle and the traffic light may not be known exactly. For example, in one time step, the oncoming vehicle may stay in the same cell or move to the next cell. Figure 5.3 shows a nondeterministic model corresponding to such behaviors.

We consider the following requirements:

- 1. The autonomous vehicle should eventually go to cell  $c_9$ .
- 2. The autonomous vehicle should not collide with the oncoming vehicle.

The corresponding game structure  $\mathcal{G} = (V, X, Y, \theta_e, \theta_s, \rho_e, \rho_s, AP, L, \varphi)$  is defined as follows:

- $X = \{x_h\}$ , where  $x_h \in \{0, ..., 9\}$  is the label of the cell occupied by the oncoming vehicle.
- $Y = \{x_a\}$ , where  $x_a \in \{0, ..., 9\}$  is the label of the cell occupied by the autonomous vehicle.
- $\theta_e = (x_h = 0 \lor x_h = 1 \lor x_h = 2)$  indicates that the oncoming vehicle can start from cell  $c_0$ ,  $c_1$  or  $c_2$ .
- $\theta_s = (x_a = 4 \lor x_a = 7 \lor x_a = 8 \lor x_a = 9)$  indicates that the autonomous vehicle can start from cell  $c_4$ ,  $c_7$ ,  $c_8$  or  $c_9$ .
- $\rho_e = \bigwedge_{i=0}^5 (x_h = i \implies (\bigcirc x_h = i \lor \bigcirc x_h = i+1)) \land (x_h = 6 \implies (\bigcirc x_h = 6))$  indicates that the oncoming vehicle may stay in the same cell or move to the next cell.
- $\rho_s = \rho_s^{trans} \wedge \rho_s^{safe}$ . Here,  $\rho_s^{tran} = (x_a = 7 \implies (\bigcirc x_a = 7 \lor \bigcirc x_a = 8)) \wedge (x_a = 8 \implies (\bigcirc x_a = 8 \lor \bigcirc x_a = 4)) \wedge (x_a = 4 \implies (\bigcirc x_a = 4 \lor \bigcirc x_a = 9)) \wedge (x_a = 9 \implies (\bigcirc x_a = 9)$  indicates that the autonomous vehicle may stay in the same cell or move to the next cell, and  $\rho_s^{safe} = \neg(x_a = 4 \land x_h = 4)$  indicates that the autonomous vehicle and the oncoming vehicle do not simultaneously occupy cell  $c_4$ .
- $\varphi = \varphi_e \implies \Box \diamondsuit (x_a = 9)$  indicating that the autonomous vehicle is at cell  $c_9$  infinitely often. Here  $\varphi_e$  represents the assumption on the environment, which we will revisit in Example 5.2. Note that from the transition system representing the autonomous vehicle, once the vehicle reaches  $c_9$ , it will stay at  $c_9$  forever. As a result,  $\diamondsuit (x_a = 9)$  and  $\Box \diamondsuit (x_a = 9)$  are equivalent winning conditions.

A game is played as follows: The environment initially chooses an assignment  $s_X \in dom(X)$  such that  $s_X \models \theta_e$ , and the system chooses an assignment  $s_Y \in dom(Y)$  such that  $(s_X, s_Y) \models \theta_e \wedge \theta_s$ . From a state s, the environment chooses an input  $s_X \in dom(X)$  such that  $(s, s_X) \models \rho_e$  and the system chooses an output  $s_Y \in dom(Y)$  such that

 $(s, s_X, s_Y) \models \rho_s$ . Formally, we define a *play* as a maximal sequence of states  $\sigma = s_0 s_1 \dots$  such that  $s_0 \models \theta_e \wedge \theta_s$  and, for every  $j \geq 0$ ,  $(s_j, s_{j+1}) \models \rho_e \wedge \rho_s$ .

A finite memory control protocol for the system can be identified with a partial function  $f: M \times dom(V) \times dom(X) \to M \times dom(Y)$  such that, for every  $s \in dom(V)$ , every  $s_X \in dom(X)$ , and every  $m \in M$ , if  $(s, s_X) \models \rho_e$  and  $f(m, s, s_X) = (m', s_Y)$ , then  $(s, s_X, s_Y) \models \rho_s$ . Here, M is some memory domain with a designated initial value  $m_0 \in M$ .

Protocol f is winning for the system starting from state  $s_0$  if any play  $\sigma = s_0 s_1 \dots$  such that, for all  $i \geq 0$ ,  $f(m_i, s_i, s_{i+1}|_X) = (m_{i+1}, s_{i+1}|_Y)$ , either (i) is infinite and satisfies  $\varphi$ , or (ii) is finite and there is no assignment  $s_X \in dom(X)$  such that  $(s_n, s_X) \models \rho_e$ , where  $s_n$  is the last state in  $\sigma$ . We let  $Win_s$  denote a proposition characterizing the set of states starting from which there exists a winning strategy for the system. A game structure is winning for the system if, for all  $s_X \in dom(X)$  such that  $s_X \models \theta_e$ , there exists  $s_Y \in dom(Y)$  such that  $(s_X, s_Y) \models \theta_s$  and  $(s_X, s_Y) \in Win_s$ .

For certain LTL specifications,  $\mu$ -calculus over game structures can be employed to characterize the set of winning states of the system. The description of  $\mu$ -calculus, however, is beyond the scope of this paper and we refer the reader to (Kozen, 1983; Piterman  $et\ al.$ , 2006). As an example, the  $\mu$ -calculus formula  $\mu R(p\vee \otimes R)$  characterizes the set of states from which the system can force the game to eventually visit p-states, i.e., states that satisfy proposition p. This formula thus provides the solution for the reachability game previously discussed. Here,  $\mu$  is the least fixpoint operator in  $\mu$ -calculus, R is known as a "relational variable" and the operator  $\otimes$  is defined roughly similar to the predecessor operator  $Pre_{\exists \forall}$ .

Piterman et al. (2006) consider a broader class of LTL formula known as generalized reactivity[1] (GR[1]) which covers LTL formulas of the form

$$\varphi = (\Box \Diamond p_1 \wedge \ldots \wedge \Box \Diamond p_m) \implies (\Box \Diamond q_1 \wedge \ldots \wedge \Box \Diamond q_n). \tag{5.1}$$

Roughly, the left hand side of  $\implies$  specifies the assumption on the

environment behavior whereas the right hand side of  $\implies$  specifies the desired property of the system. Piterman et al. (2006) show that there exists a  $\mu$ -calculus formula that characterizes the set of winning states of the system for GR[1] winning conditions. The formulation allows for the synthesis problem to be solved based on fixpoint computation in time proportional to  $nm|dom(V)|^3$ , where |dom(V)| is the size of the state space. The proposed synthesis procedure has been implemented in JTLV (Piterman et al., 2006) and in TuLiP (Wongpiromsarn et al., 2011b). We refer the reader to (Piterman et al., 2006) for more details, including a discussion on the expressiveness of GR[1] and an extension to handle formulas of the form  $\varphi_e \implies \varphi_s$  where  $\varphi_e$  and  $\varphi_s$  are any LTL formulas that can be represented by a deterministic Büchi automaton. A deterministic Büchi automaton is defined as a non-deterministic Büchi automaton with additional constraints that  $|Q_0| \leq 1$  and for any  $q \in Q$  and  $\sigma \in \Sigma$ ,  $(q, \sigma, q') \in \delta$  and  $(q, \sigma, q'') \in \delta$  imply that q'=q''. LTL formulas that can be represented by a deterministic Buchi automaton include those of the form  $\Box(p_1 \implies \Diamond p_2)$  where  $p_1$  and  $p_2$ are propositions.

**Example 5.2.** Consider the unprotected left turn scenario described in Example 5.1. Without any assumption on the oncoming vehicle (i.e.,  $\varphi_e = True$ ), the game structure  $\mathcal{G}$  defined in Example 5.1 is not winning for the system. To see this, consider the case where the oncoming vehicle starts at  $c_2$  while the autonomous vehicle starts at  $c_7$ . In two steps, the oncoming vehicle can reach  $c_4$  and stay there forever, blocking the autonomous vehicle from entering  $c_4$  without violating the safety requirement. A sufficient assumption to ensure that the game structure will be winning is that the oncoming vehicle visits cell  $c_6$  infinitely often, i.e.,  $\varphi_e = \Box \diamondsuit (x_h = 6)$ . With this assumption, a possible control protocol for the autonomous vehicle is to wait until the oncoming vehicle reaches  $c_5$  or  $c_6$  before entering  $c_4$ .

**Remark 5.1.** As demonstrated in Example 5.2, the system (e.g., the autonomous vehicle) has no control over the environment (e.g., the oncoming vehicle). As a result, the synthesis algorithm needs to take into account all the possible values of the environment variables (e.g.,

 $x_h$ ) and ensures that the resulting behavior of the system and the environment satisfies the specification  $\varphi$ . In many cases, in order to obtain a solution, one needs to limit the power of the environment, which is captured by  $\theta_e$ ,  $\rho_e$ , and  $\varphi_e$  in the proposed model.

#### 5.2 Receding Horizon Temporal Logic Planning

The main limitation of the discrete synthesis described in Section 4.2 and Section 5.1 is the state explosion problem. In the worst case, the entire system's state space has to be taken into account. For example, if the system has |V| variables, each can take any of the P possible values. Then, we must consider as many as  $P^{|V|}$  states. This type of computational complexity limits the application of systhesis to relatively small problems.

Similar computational complexity is also encountered in the area of constrained optimal control. In the controls domain, an effective and well-established technique to address this issue is to design and implement control strategies in a receding horizon manner, i.e., optimize over a *shorter* horizon, starting from the currently observed state, implement the initial control action, move the horizon one step ahead, and re-optimize. This approach reduces the computational complexity by essentially solving a sequence of *smaller* optimization problems, each with a specific initial condition (as opposed to optimizing with *any* initial condition in traditional optimal control). Under certain conditions, receding horizon control strategies are known to lead to closed-loop stability (Murray *et al.*, 2003; Mayne *et al.*, 2000; Jadbabaie, 2000). See, e.g., (Goodwin *et al.*, 2004) for a detailed discussion on constrained optimal control, including finite horizon optimal control and receding horizon control.

To partially alleviate the state explosion problem in the synthesis of finite state automata, Wongpiromsarn  $et\ al.\ (2010b)$  and Wongpiromsarn  $et\ al.\ (2012)$  consider reactive module synthesis with GR[1] specifications and show that for systems with a certain structure, the synthesis problem can be solved in a receding horizon fashion, i.e., compute the plan or strategy over a "shorter" horizon, starting from the current state,

$C_{3,1}$	$C_{3,2}$	$C_{3,L}$
$C_{2,1}$	$C_{2,2}$	 $C_{2,L}$
$C_{1,1}$	$C_{1,2}$	$C_{1,L}$

Figure 5.4: The autonomous driving example where the road is partitioned into 3L cells where L is the length of the road.

implement the initial portion of the plan, move the horizon one step ahead, and recompute. This approach essentially reduces the discrete control protocol synthesis problem into a set of smaller problems. The size of these smaller problems depends on the horizon length. For example, consider the autonomous driving problem where an autonomous vehicle needs to navigate the road shown in Figure 5.4 starting from cell  $C_{1,1}$  and with destination  $C_{1,L} \cup C_{2,L} \cup C_{3,L}$ . Suppose the horizon length is l, i.e., the vehicle plans for l cells ahead. Then, the state space for each short-horizon problem contains at most  $3l2^{3l}$  states (whereas the size of the original problem is  $3L2^{3L}$ ). Hence, the horizon length should be made as small as possible, subject to the realizability of the resulting short-horizon specifications; horizons that are too short typically render the specifications unrealizable.

Sufficient conditions that ensure that this receding horizon implementation preserves the desired system-level properties are presented in (Wongpiromsarn et al., 2010b; Wongpiromsarn et al., 2012). For the simplicity of the presentation, in this article, we consider the case where the specification is given by

$$\varphi = (\varphi_{init} \wedge \varphi_{env}) \implies (\varphi_{safety} \wedge \varphi_{goal}),$$
 (5.2)

where  $\varphi_{init}$  is a proposition characterizing the set of initial states,  $\varphi_{env}$  is an LTL formula characterizing the assumption on the environment behavior and can be written as the conjunction of a safety formula and the progress formulas on the left hand side of  $\Longrightarrow$  in (5.1),  $\varphi_{safety}$  is a safety formula and  $\varphi_{goal}$  is of the form  $\varphi_{goal} = \Box \Diamond q$  where q is a proposition characterizing the set of some goal states to be visited infinitely often.

The receding horizon approach works as follows. First, we organize the

discrete state space into a partially ordered set  $(\{W_0, \ldots, W_M\}, \prec_{\varphi})$  such that  $W_0$  only contains the goal states and  $W_0 \prec_{\varphi} W_i$  for all  $i \neq 0$ . The partial order relation  $\prec_{\varphi}$  can be defined based on the notion of "distance" to the goal states.

Next, we define a map  $\mathcal{F}: \{\mathcal{W}_0, \dots, \mathcal{W}_M\} \to \{\mathcal{W}_0, \dots, \mathcal{W}_M\}$  that captures the horizon length and satisfies  $\mathcal{F}(\mathcal{W}_i) \prec_{\varphi} \mathcal{W}_i$  for all  $i \neq 0$ . Finally, we specify a proposition  $\Phi$  that characterizes the receding horizon invariant such that any state that satisfies  $\varphi_{init}$  also satisfies  $\Phi$ . In other words,  $\varphi_{init} \Longrightarrow \Phi$  is a tautology.

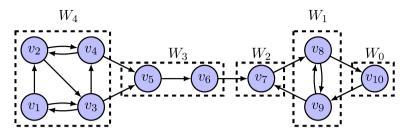
With the partially ordered set  $(\{W_0, \ldots, W_M\}, \prec_{\varphi})$ , the map  $\mathcal{F}$ , and the receding horizon invariant  $\Phi$ , we define a short-horizon specification  $\Psi_i$  associated with each  $W_i$ , for  $i \in \{0, \ldots, M\}$  as

$$\Psi_i = ((\nu \in \mathcal{W}_i) \land \Phi \land \varphi_{env}) \implies (\Box \Phi \land \varphi_{safety} \land \Diamond (\nu \in \mathcal{F}(\mathcal{W}_i))) \quad (5.3)$$

where  $\nu$  is the state of the system. The left-hand side of  $\Longrightarrow$  states that (a) the initial state is assumed to be in  $W_i$  and satisfies  $\Phi$ , and (b) the environment is assumed to satisfy the assumptions  $\varphi_{env}$  stated in the original specification. The right-hand side of  $\Longrightarrow$  then specifies that (a)  $\Phi$  holds throughout an execution, (b) the original safety properties  $\varphi_{safety}$  are satisfied, and (c) the system eventually reaches a state in  $\mathcal{F}(W_i)$ .

According to (5.3),  $\mathcal{F}(W_i)$  essentially defines an intermediate goal for states in  $W_i$ . In addition,  $\Phi$  is introduced to ensure that a provably correct plan exists when the system reaches the end of the current horizon and needs to compute a new plan. We refer the reader to (Wongpiromsarn *et al.*, 2012) for a detailed discussion on this receding horizon framework, including an extension to the case where there are multiple goals that may be visited in an arbitrary order.

Consider a simple example shown in Figure 5.5 where  $\nu_{10}$  is the goal state. The partial order may be defined as  $W_0 \prec_{\varphi} W_1 \prec_{\varphi} \ldots \prec_{\varphi} W_4$  and the map  $\mathcal{F}$  may be defined as  $\mathcal{F}(W_j) = W_{j-2}$ , for all  $j \geq 2$  and  $\mathcal{F}(W_1) = \mathcal{F}(W_0) = W_0$ . The key idea of the receding horizon framework is to synthesize a control protocol for short-horizon specification  $\Psi_4$ , which corresponds to going from  $\nu_1$  only to a state in  $\mathcal{F}(W_4) = W_2$ , rather than



**Figure 5.5:** A graphical description of the receding horizon framework for a special case where there is only one goal  $\nu_{10}$ .  $\nu_{1}, \ldots, \nu_{10}$  are the discrete states.

synthesizing a control protocol for going from the initial state  $\nu_1$  to the goal state  $\nu_{10}$  in one shot, taking into account all the possible behavior of the environment. Once a state in  $\mathcal{W}_3$ , i.e.,  $\nu_5$  or  $\nu_6$  is reached, we then recompute a protocol for the short-horizon specification  $\Psi_3$  for going to a state in  $\mathcal{F}(\mathcal{W}_3) = \mathcal{W}_1$ . This process is then continually repeated. From the finiteness of the set  $\{\mathcal{W}_0, \ldots, \mathcal{W}_M\}$  and its partial order, it can be shown that this receding horizon implementation of the short-horizon strategies ensures the correctness of the global specification, provided that all of the short horizon specifications  $\Psi_i$ , for  $i \in \{0, \ldots M\}$  are realizable (Wongpiromsarn et al., 2012). In this case, the invariant  $\Phi$  is introduced to rule out the states that render the short horizon problems unrealizable.

Given a receding horizon invariant  $\Phi$ , the partial order  $\prec_{\varphi}$  as well as the horizon length defined by the map  $\mathcal{F}$  can be automated by adding an additional component, namely the goal generator to the hierarchical control structure in Figure 5.6. The goal generator works on a graph  $\mathbb{G}$  with  $W_i, i \in \{0, ..., M\}$  being its states. For each  $i \in \{0, ..., M\}$  and  $j \in \{0, ..., M\}$ , a transition from  $W_i$  to  $W_j$  in  $\mathbb{G}$  is added if  $i \neq j$  and the short-horizon specification  $\Psi_i$  is realizable with  $\mathcal{F}(W_i) = W_j$ . After  $\mathbb{G}$  is constructed, the goal generator then performs a graph search to find a path from  $W_i$ , to which the current state of the system belongs, to a goal state in  $W_0$ . This path essentially defines a sequence of intermediate goals for each short-horizon problem. The resulting hierarchical control structure with this implementation of the receding-horizon framework is shown in Figure 5.6. Figure 5.7 shows the similarity

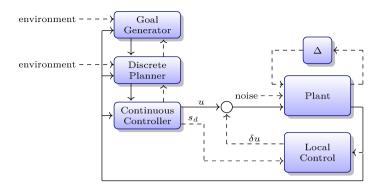


Figure 5.6: The hierarchical control structure with the goal generator.

of this hierarchical control structure with that implemented on Alice, the representative autonomous vehicle in Section 3.4, illustrating that the techniques presented in this article can be utilized to formalize and enable automatic design of the navigation protocol stack of an autonomous system.

Computational Complexity and Completeness: The receding-horizon implementation reduces the computational complexity by restricting the state space considered in each subproblem; however, it is not complete. Even if the original specification is realizable, there may not exist a combination of horizon length, partial order relation, and receding horizon invariant that render all of the short horizon specifications realizable. Nevertheless, its successful applications to autonomous driving problems have been illustrated in Wongpiromsarn et al. (2010a) and Wongpiromsarn et al. (2012). Examples of these applications are provided in Figure 5.9.

Remark 5.2. Computation of the horizon length, partial order relation, and receding horizon invariant requires insights for each problem domain. Automatic construction of these elements is subject to on-going research. Wongpiromsarn et al. (2012) describe automatic construction of certain elements, given other elements, e.g., automatic computation of the horizon length and partial order relation, given a receding horizon invariant, and automatic computation of the receding horizon invariant,

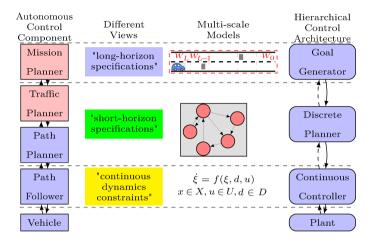


Figure 5.7: The hierarchical control structure with the receding horizon implementation, showing the similarity with the navigation protocol stack implemented on Alice. The goal generator has similar functionality as Mission Planner. It determines a sequence of intermediate goals for the discrete planner such that the original "long-horizon" specification is satisfied. Its computation relies on the graph  $\mathbb{G}$  that encodes the partially ordered set  $(\{W_0,\ldots,W_M\}, \prec_{\varphi})$ . The discrete planner has similar functionality as the composition of Traffic Planner and Path Planner. It computes a discrete plan for the system such that the short-horizon specification in (5.3) with the next intermediate goal computed by the goal generator is satisfied based on a finite, abstract model of the physical system. Finally, the continuous controller deals with the continuous dynamics and constraints to ensure that the physical system follows the plan computed by the discrete planner. This functionality is similar to that of Path Follower in Alice.

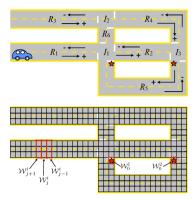


Figure 5.8: The road network and its partition for the autonomous vehicle example. The stars indicate the cells that need to be visited infinitely often.

given a horizon length and partial order relation.

**Example 5.3.** Consider an autonomous driving problem in an urban-like environment. We consider the road network shown in Figure 5.8, which is partitioned into N = 282 cells. Each of these cells may or may not be occupied by an obstacle. The desired properties include:

- Each of the two cells marked by star needs to be visited infinitely often.
- No collision is allowed, i.e., the vehicle cannot occupy the same cell as an obstacle.
- The vehicle stays in the right lane unless there is an obstacle blocking the lane.
- The vehicle can only proceed through an intersection when the intersection is clear.

Wongpiromsarn et al. (2012) show that with some mild assumptions on the environment behavior, there exists a receding-horizon invariant that ensures that all the short-horizon specifications are realizable with horizon length 2, i.e,  $\mathcal{F}(W_i) = \mathcal{F}(W_{i-2})$ . Hence, the size of the state space for each short-horizon problem is at most 4608 whereas the size of

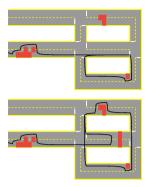


Figure 5.9: Simulation results with (top) no road blockage, (bottom) a road blockage on the middle road. The corresponding movies can be downloaded from http://sourceforge.net/projects/tulip-control/.

the state space of the original problem is in the order of  $10^{87}$ . Roughly, the receding-horizon invariant requires that the vehicle is not surrounded by obstacles and if the vehicle is not in the travel lane, there must be an obstacle blocking the lane. Using JTLV, each short-horizon synthesis problem can be solved in approximately 1.5 seconds on a MacBook with a 2 GHz Intel Core 2 Duo processor and 4 Gb of memory. Simulation results when the receding-horizon approach is applied are shown in Figure 5.9.

# 5.3 Reactive Synthesis with Maximum Realizability

In conventional synthesis, either an implementation is constructed for a given specification, or the specification is identified as unrealizable. Nevertheless, specifications may arise from different design perspectives, especially in large systems, and if they consist of a large number of individual requirements, it is easy to encounter specifications that are unrealizable. In other scenarios, the user may have several alternative requirements in mind, potentially with some preferences, and want to know the best realizable combination of them with respect to some metric. Such cases usually lead to alternating between specification modification and synthesis procedure and hence, defeating the purpose of facilitating the design process.

The possibility of conflict amongst the provided requirements calls for a more comprehensive synthesis procedure that, in the case of unrealizability, can generate an implementation that minimally violates the specifications. In order to define the notion of minimality, one requires a quantitative metric on the satisfaction of LTL formulas. The approach we pursue in this section relies on multiple levels of relaxations of an LTL formula as well as relying on forming a value function that captures the levels of relaxations applied over the specifications. Then, the maximum realizability of a set of LTL formulas turns into seeking an implementation that maximizes the corresponding value.

The backbone of this section's approach toward maximum realizability is bounded synthesis, originally introduced by Schewe and Finkbeiner Schewe and Finkbeiner (2007). Bounded synthesis tackles the computational complexity of reactive synthesis from LTL properties by restricting the size of the search space and incrementing the bound on the size if necessary. More specifically, it searches for a realizable implementation of the size up to a prespecified bound. If no such implementation exists, it increments the bound and repeats the search process. Each instance of bounded search for an implementation can be encoded as a SAT (or QBF, or SMT) problem (Faymonville et al., 2017). The algorithm is complete as a theoretical bound on the maximum size of the implementation exists.

In this section, we formulate maximum realizability as iterative maximum satisfiability (MaxSAT) solving (Biere et al., 2009). In each iteration, we construct a MaxSAT instance that characterizes the existence of an implementation of size within the given bound that not only realizes the hard specification but is also optimal with respect to defined value function for soft specifications. We prove that, for any given finite set of soft specifications, there exists an optimal implementation with a bounded size. Consequently, the proposed algorithm that gradually increases the bound on the implementation size is complete.

#### 5.3.1 Related Work

Maximum realizability and several closely related problems have attracted significant attention in recent years. Tumova et al. (2013) studied the problem of planning over a finite horizon with prioritized safety requirements, where the goal is to synthesize a least-violating control strategy. Kim et al. (2015) studied a similar problem for the case of infinite-horizon temporal logic planning. Lahijanian et al. (2015) describe a method for computing plans for co-safe LTL specifications that minimize the cost of violating each atomic proposition. These approaches are developed for the planning setting without an adversarial environment. Lahijanian and Kwiatkowska (2016) considered the case of probabilistic environments and Lahijanian et al. (2016) studied the problem of partial satisfaction of guarantees in an unknown environment, maximizing the number of soft specifications that are satisfied. Tomita et al. (2017) study a maximum realizability problem in which the specification is a conjunction of a must LTL specification, and a number of weighted desirable LTL specifications, formulated as a mean-payoff optimization.

Two other main research directions related to maximum realizability are quantitative synthesis and specification debugging. Related to quantitative synthesis, the goal in Bloem et al. (2009) is to generate an implementation that maximizes the value of a mean-payoff objective, while possibly satisfying some  $\omega$ -regular specification, while in Almagor et al. (2016) and Tabuada and Neider (2016), the system requirements are formalized in a multi-valued temporal logic. Alur et al. (2008) studied an optimal synthesis problem for an ordered sequence of prioritized  $\omega$ -regular properties, where the classical fixpoint-based game-solving algorithms are extended to a quantitative setting. In specification debugging there is a lot of research dedicated to finding good explanations for the unsatisfiability or unrealizability of temporal logic specifications (Cimatti et al., 2007; Schuppan, 2012; Raman and Kress-Gazit, 2013), and more generally to the analysis of specifications (Cimatti et al., 2008; Ehlers and Raman, 2014).

We start by an overview of the required concepts to formally state the

synthesis problem. Then, we proceed to go over definitions of run graph and annotations to describe the bounded synthesis method and its SAT encoding. Lastly, we provide a brief description of the MaxSAT problem, particularly a class of that called partial weighted MaxSAT.

### **Bounded Synthesis Approach**

The run graph of a universal automaton  $\mathcal{A} = (Q, \Sigma, \delta, Q_0, F)$  on a transition system  $\mathcal{T} = (S, s_0, \tau)$  is the unique graph G = (V, E) with a set of nodes  $V = S \times Q$  and a set of labeled edges  $E \subseteq V \times \Sigma \times V$  such that  $((s, q), \sigma, (s', q')) \in E$  if and only if  $(q, \sigma, q') \in \delta$  and  $\tau(s, \sigma \cap \mathcal{I}) = (s', \sigma \cap \mathcal{O})$ . That is, G is the product of  $\mathcal{A}$  and  $\mathcal{T}$ .

A run graph of a universal Büchi (resp. co-Büchi) automaton is accepting if every infinite path  $(s_0, q_0), (s_1, q_1), \ldots$  contains infinitely (resp. finitely) many occurrences of states  $q_i$  in F. A transition system  $\mathcal{T}$  is accepted by a universal automaton  $\mathcal{A}$  if the unique run graph of  $\mathcal{A}$  on  $\mathcal{T}$  is accepting. We denote with  $\mathcal{L}(\mathcal{A})$  the set of transition systems accepted by  $\mathcal{A}$ .

The bounded synthesis approach is based on the fact that for every LTL formula  $\varphi$  one can construct a universal co-Büchi automaton  $\mathcal{A}_{\varphi}$  with at most  $2^{O(|\varphi|)}$  states such that  $\mathfrak{T} \in \mathcal{L}(\mathcal{A}_{\varphi})$  iff  $\mathfrak{T} \models \varphi$ , for every transition system  $\mathfrak{T}$  (Kupferman and Vardi, 2005).

An annotation of a transition system  $\mathfrak{T}=(S,s_0,\tau)$  with respect to a universal co-Büchi automaton  $\mathcal{A}=(Q,\Sigma,\delta,Q_0,F)$  is a function  $\lambda:S\times Q\to\mathbb{N}\cup\{\bot\}$  that maps nodes of the run graph of  $\mathcal{A}$  on  $\mathcal{T}$  to the set  $\mathbb{N}\cup\{\bot\}$ . Intuitively, such an annotation is valid if every node (s,q) that is reachable from the node  $(s_0,q_0)$  is annotated with a natural number, which is an upper bound on the number of rejecting states visited on any path from  $(s_0,q_0)$  to (s,q). Valid annotations of finite-state transition systems correspond to accepting run graphs. An annotation  $\lambda$  is c-bounded if  $\lambda(s,q)\in\{0,\ldots,c\}\cup\{\bot\}$  for all  $s\in S$  and  $q\in Q$ .

The synthesis method proposed in (Schewe and Finkbeiner, 2007; Finkbeiner and Schewe, 2013) employs the following result in order to reduce the bounded synthesis problem to checking the satisfiability of

propositional formulas: a transition system  $\mathcal{T}$  is accepted by a universal co-Büchi automaton  $\mathcal{A}=(Q,\Sigma,\delta,Q_0,F)$  iff there exists a  $(|\mathcal{T}|\cdot|F|)$ -bounded valid annotation for  $\mathcal{T}$  and  $\mathcal{A}$ . One can estimate a bound on the size of the transition system, which allows to reduce the synthesis problem to its bounded version. Namely, if there exists a transition system that satisfies an LTL formula  $\varphi$ , then there exists a transition system satisfying  $\varphi$  with at most  $(2^{(|\mathsf{subf}(\varphi)|+\log|\varphi|)})!^2$  states, where  $|\varphi|$  denotes the size of the formula  $\varphi$  and  $\mathsf{subf}(\varphi)$  denotes the set of all subformulas of  $\varphi$ .

#### **MaxSAT**

Consider a propositional logic formula in conjunctive normal form (CNF), i.e., a formula that is a conjunction of disjunction of literals, where a literal is a Boolean variable or its negation and a disjunction of literals is called a clause. MaxSAT is the problem of assigning truth values to a set of Boolean variables such that the number of clauses of a propositional logic formula in CNF that are made true, is maximized (Biere et al., 2009). partial weighted MaxSAT is a variant of MaxSAT problem where the clauses are categorized as hard and soft clauses and each of the soft clauses is associated with a positive numerical weight. The objective is to find a truth assignment to the variables that not only makes all the hard clauses true but also maximizes the sum of the weights of the soft clauses that become true.

We exploit the separation of the hard and soft clauses in partial weighted MaxSAT to capture the hard and soft constraints that arise in the encoding of the maximum realizability problem. Furthermore, we design the weights of the soft clauses in a way to promote the quantitative objective associated with the conjunction of the given soft specifications.

The procedure proposed by Finkbeiner and Schewe Finkbeiner and Schewe (2013) provides a SAT encoding of synthesis when the size of the implementation is bounded. Maximum realizability is an optimization variant of synthesis while MaxSAT is an optimization variant of SAT. For a proposed value function, the maximum realizability problem under a bounded implementation size can be reduced to a *partial weighted* 

MaxSAT instance.

### 5.3.2 Maximum Realizability

Let  $\Box \varphi_1, \ldots, \Box \varphi_n$  be a set of LTL specifications, where each  $\varphi_i$  is a safety LTL formula. In order to formalize the maximal satisfaction of  $\Box \varphi_1 \wedge \ldots \wedge \Box \varphi_n$ , we first give a quantitative semantics of formulas of the form  $\Box \varphi$ .

Quantitative semantics of safety specifications. For an LTL formula of the form  $\Box \varphi$  and a transition system  $\mathcal{T}$ , we define the value  $val(\mathcal{T}, \Box \varphi)$  of  $\Box \varphi$  in  $\mathcal{T}$  as

$$val(\mathfrak{I},\square\varphi)\stackrel{\mathrm{def}}{=} \begin{cases} (1,1,1) & \text{if } \mathfrak{I}\models \square\varphi, \\ (1,1,0) & \text{if } \mathfrak{I}\not\models \square\varphi \text{ and } \mathfrak{I}\models \diamondsuit\square\varphi, \\ (1,0,0) & \text{if } \mathfrak{I}\not\models \square\varphi \text{ and } \mathfrak{I}\not\models \diamondsuit\square\varphi \text{ and } \mathfrak{I}\models \square\diamondsuit\varphi, \\ (0,0,0) & \text{if } \mathfrak{I}\not\models \square\varphi \text{ and } \mathfrak{I}\not\models \diamondsuit\square\varphi \text{ and } \mathfrak{I}\not\models \square\diamondsuit\varphi. \end{cases}$$

The value of  $\Box \varphi$  in a transition system  $\mathfrak{T}$  is a vector  $(v_1, v_2, v_3) \in \{0, 1\}^3$ , where the value (1, 1, 1) corresponds to the *True* value in the classical semantics of LTL. When  $\mathfrak{T} \not\models \Box \varphi$ , the values (1, 1, 0), (1, 0, 0) and (0, 0, 0) capture the extent to which  $\varphi$  holds or not along the traces of  $\mathfrak{T}$ . For example, if  $val(\mathfrak{T}, \Box \varphi) = (1, 0, 0)$ , then  $\varphi$  holds infinitely often on each trace of  $\mathfrak{T}$ , but there exists a trace of  $\mathfrak{T}$  on which  $\varphi$  is violated infinitely often. If  $val(\mathfrak{T}, \Box \varphi) = (0, 0, 0)$ , then, on some trace of  $\mathfrak{T}$ ,  $\varphi$  holds for at most finitely many positions. Thus, the lexicographic ordering on  $\{0, 1\}^3$  captures the preference of one transition system over another with respect to the quantitative satisfaction of  $\Box \varphi$ .

**Example 5.4.** Suppose that we want to synthesize a transition system representing a navigation strategy for a robot working at a restaurant. We require that the robot serves the VIP area infinitely often, formalized in LTL as  $\Box \diamondsuit vip\_area$ . We also desire that the robot never enters the staff's office, formalized as  $\Box \neg office$ . Now, suppose that initially the key to the VIP area is in the office. Thus, in order to satisfy  $\Box \diamondsuit vip\_area$ , the robot must violate  $\Box \neg office$ . A strategy in which the office is entered

only once, and satisfies  $\diamondsuit \square \neg office$ , is preferable to one which enters the office over and over again, and only satisfies  $\square \diamondsuit \neg office$ . Thus, we want to synthesize a strategy  $\Im$  maximizing  $val(\Im, \square \neg office)$ .

In order to compare implementations with respect to their satisfaction of a conjunction of several safety specifications  $\Box \varphi_1 \wedge \ldots \wedge \Box \varphi_n$ , we will extend the above definition. We first consider the case where the specifier has not expressed any preference for the individual conjuncts and later on, extend that to the case with a given priority ordering. Consider the following example.

**Example 5.5.** We consider again the restaurant robot, now with two soft specifications. The soft specification  $\Box(req1 \to \bigcirc table1)$  requires that each request by table 1 is served immediately at the next time instance. Similarly,  $\Box(req2 \to \bigcirc table2)$ , requires the same for table number 2. Since the robot cannot be at both tables simultaneously, formalized as the hard specification  $\Box(\neg table1 \lor \neg table2)$ , the conjunction of these requirements is unrealizable. Unless the two tables have priorities, it is preferable to satisfy each of  $req1 \to \bigcirc table1$  and  $req2 \to \bigcirc table2$  infinitely often, rather than serve one and the same table all the time.

Quantitative semantics of conjunctions. To capture the idea illustrated in Example 5.5, we define a value function, which intuitively gives higher values to transition systems in which a fewer number of soft specifications have low values. Formally, let the value of  $\Box \varphi_1 \wedge \ldots \wedge \Box \varphi_n$  in  $\Im$  be

$$val(\mathfrak{T}, \square \varphi_1 \wedge \ldots \wedge \square \varphi_n) \stackrel{\text{\tiny def}}{=} \left(\sum_{i=1}^n v_{i,1}, \sum_{i=1}^n v_{i,2}, \sum_{i=1}^n v_{i,3}\right),$$

where  $val(\mathfrak{T}, \square \varphi_i) = (v_{i,1}, v_{i,2}, v_{i,3})$  for  $i \in \{1, \ldots, n\}$ . To compare transition systems according to these values, we use lexicographic ordering on  $\{0, \ldots, n\}^3$ .

**Example 5.6.** For the specifications in Example 5.5, the defined value function assigns value (2,0,0) to a system satisfying  $\Box \diamondsuit (req1 \rightarrow \bigcirc table1)$  and  $\Box \diamondsuit (req2 \rightarrow \bigcirc table2)$ , but neither of  $\diamondsuit \Box (req1 \rightarrow \bigcirc table1)$ 

and  $\diamondsuit \square (req2 \to \bigcirc table2)$ . It assigns the smaller value (1,1,1) to an implementation that gives priority to table 1 and satisfies  $\square (req1 \to \bigcirc table1)$  but not  $\square \diamondsuit (req2 \to \bigcirc table2)$ .

According to the definition above, a transition system that satisfies all soft requirements to some extent is considered better in the lexicographic ordering than a transition system that satisfies one of them exactly and violates all the others. We could instead inverse the order of the sums in the triple, thus giving preference to satisfying some soft specification exactly, over having some lower level of satisfaction over all of them. The next example illustrates the differences between the two variations.

**Example 5.7.** For the two soft specifications from Example 5.5, reversing the order of the sums in the definition of  $val(\mathfrak{I}, \Box\varphi_1 \wedge \ldots \wedge \Box\varphi_n)$  results in giving the higher value (1,1,1) to a transition system that satisfies  $\Box(req1 \to \bigcirc table1)$  but not  $\Box \diamondsuit(req2 \to \bigcirc table2)$ , and the lower value (0,0,2) to the one that only guarantees  $\Box \diamondsuit(req1 \to \bigcirc table1)$  and  $\Box \diamondsuit(req2 \to \bigcirc table2)$ . The most suitable ordering usually depends on the specific application.

#### 5.3.3 Problem Formulation

Using the definition of quantitative satisfaction of soft safety specifications, we now define the maximum realizability problem, which asks to synthesize a transition system that satisfies a given *hard* LTL specification, and is optimal with respect to the satisfaction of a conjunction of *soft* safety specifications.

Bounded maximum realizability problem: Given an LTL formula  $\varphi$  and formulas  $\Box \varphi_1, \ldots, \Box \varphi_n$ , where each  $\varphi_i$  is a safety LTL formula, and a bound  $b \in \mathbb{N}_{>0}$ , the bounded maximum realizability problem asks to determine if there exists a transition system  $\mathcal{T}$  with  $|\mathcal{T}| \leq b$  such that  $\mathcal{T} \models \varphi$ , and if the answer is positive, to synthesize a transition system  $\mathcal{T}$  such that  $\mathcal{T} \models \varphi$ ,  $|\mathcal{T}| \leq b$  and for every transition system  $\mathcal{T}'$  with  $\mathcal{T}' \models \varphi$  and  $|\mathcal{T}'| \leq b$ , it holds that  $val(\mathcal{T}, \Box \varphi_1 \wedge \ldots \wedge \Box \varphi_n) \geq val(\mathcal{T}', \Box \varphi_1 \wedge \ldots \wedge \Box \varphi_n)$ .

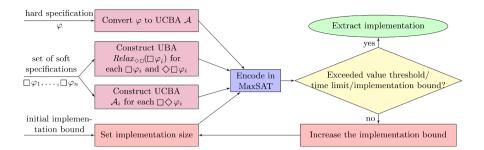


Figure 5.10: Outline of the maximum realizability procedure.

### 5.3.4 Maximum Realizability as Iterative MaxSAT Solving

We now describe the MaxSAT-based approach to maximum realizability proposed by Dimitrova *et al.* (2018). The approach first establishes an upper bound on the minimal size of an implementation that satisfies a given LTL specification  $\varphi$  and maximizes the satisfaction of a conjunction of the soft specifications  $\square \varphi_1, \ldots, \square \varphi_n$ , according to the value function defined in Section 5.3.2. This bound can be used to reduce the maximum realizability problem to its bounded version, which will be encoded as a MaxSAT problem.

For each of the possible values of  $\Box \varphi_1 \land \ldots \land \Box \varphi_n$  there is a corresponding LTL formula that encodes this value in the classical LTL semantics. This property can be utilized to establish an upper bound on the minimal optimal implementation.

**Theorem 5.1.** Given an LTL specification  $\varphi$  and soft safety specifications  $\square \varphi_1, \ldots, \square \varphi_n$ , if there exists a transition system  $\mathfrak{T} \models \varphi$ , then there exists  $\mathfrak{T}^*$  such that

(1) 
$$val(\mathfrak{I}^*, \square \varphi_1 \wedge \ldots \wedge \square \varphi_n) \geq val(\mathfrak{I}, \square \varphi_1 \wedge \ldots \wedge \square \varphi_n)$$
 for all  $\mathfrak{I}$  with  $\mathfrak{I} \models \varphi$ ,

(2) 
$$\mathfrak{T}^* \models \varphi$$
 and  $|\mathfrak{T}^*| \le \left( (2^{(b+\log b)})! \right)^2$ ,

where  $b = \max\{|\mathsf{subf}(\varphi \land \varphi_1' \land \ldots \land \varphi_n')| \mid \forall i : \varphi_i' \in \{\Box \varphi_i, \Diamond \Box \varphi_i, \Box \Diamond \varphi_i\}\}.$ 

The bound above is estimated based on the size of the specifications,

using a worst-case bound on the size of the corresponding automata. Given the automata for all the specifications  $\Box \varphi_i, \Diamond \Box \varphi_i$  and  $\Box \Diamond \varphi_i$ , a potentially better bound can be estimated based on the sizes of these automata.

Figure 5.10 gives an overview of the maximum realizability procedure and the automata constructions it involves. As in the bounded synthesis approach, we construct a universal co-Büchi automaton  $\mathcal{A}$  for the hard specification  $\varphi$ . For each soft specification  $\Box \varphi_j$ , we construct a pair of automata corresponding to the relaxations of  $\Box \varphi_j$ . The relaxation  $\Box \diamondsuit \varphi_j$  is treated as in bounded synthesis. For  $\Box \varphi_i$  and  $\diamondsuit \Box \varphi_i$ , we construct a single universal Büchi automaton and define a corresponding annotation function.

### MaxSAT Encoding of Bounded Maximum Realizability

Let  $\mathcal{A} = (Q, \Sigma, \delta, Q_0, F)$  be a universal co-Büchi automaton for the LTL formula  $\varphi$ . For each syntactically safe formula  $\square \varphi_j$ ,  $j \in \{1, \ldots, n\}$ , we consider two universal automata: the universal automaton  $\mathcal{B}_j = Relax_{\Diamond_{\square}}(\square \varphi_j) = (Q_j, \Sigma, \delta_j, Q_0^j, F_j)$  and a universal co-Büchi automaton  $\mathcal{A}_j = (\widehat{Q}_j, \Sigma, \widehat{\delta}_j, \widehat{Q}_0^j, \widehat{F}_j)$  for the formula  $\square \diamondsuit \varphi_j$ . Given a bound b on the size of the desired transition system, we encode the bounded maximum realizability problem as a MaxSAT problem.

A transition system extracted from an optimal satisfying assignment for the MaxSAT problem is optimal with respect to the value of  $\Box \varphi_1 \wedge \ldots \wedge \Box \varphi_n$ , as stated in the following theorem that establishes the correctness of the encoding.

**Theorem 5.2.** Let  $\mathcal{A}$  be a given co-Büchi automaton for  $\varphi$ , and for each  $j \in \{1, \ldots, n\}$ , let  $\mathcal{B}_j = Relax_{\Diamond \square}(\square \varphi_j)$  be the universal automaton for  $\square \varphi_j$ , and let  $\mathcal{A}_j$  be a universal co-Büchi automaton for  $\square \diamondsuit \varphi_j$ . The constraint system for bound  $b \in \mathbb{N}_{>0}$  is satisfiable if and only if there exists an implementation  $\mathcal{T}$  with  $|\mathcal{T}| \leq b$  such that  $\mathcal{T} \models \varphi$ . Furthermore, from the optimal satisfying assignment to the variables  $\tau_{s,\sigma_I,s'}$  and  $o_{s,\sigma_I}$ , one can extract a transition system  $\mathcal{T}^*$  such that for every transition system  $\mathcal{T}$  with  $|\mathcal{T}| \leq b$  and  $\mathcal{T} \models \varphi$  it holds that

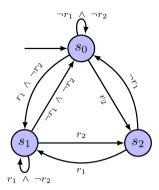


Figure 5.11: An optimal implementation for Example 5.5.

$$val(\mathfrak{T}^*, \square \varphi_1 \wedge \ldots \wedge \square \varphi_n) \geq val(\mathfrak{T}, \square \varphi_1 \wedge \ldots \wedge \square \varphi_n).$$

Figure 5.11 shows a transition system extracted from an optimal satisfying assignment for Example 5.5 with bound 3 on the implementation size. The transitions depicted in the figure are defined by the values of the variables  $\tau_{s,\sigma_I,s'}$ . The outputs of the implementation (omitted from the figure) are defined by the values of  $o_{s,\sigma_I}$ . The output in state  $s_1$  when  $r_1$  is true is  $table_1 \land \neg table_2$ , and the output in  $s_2$  when  $r_2$  is true is  $\neg table_1 \land table_2$ . For all other combinations of state and input, the output is  $\neg table_1 \land \neg table_2$ .

The next proposition establishes the size of the MaxSAT encoding.

**Proposition 5.1.** Let  $\mathcal{A}$  be a given co-Büchi automaton for  $\varphi$ , and for each  $j \in \{1, \ldots, n\}$ , let  $\mathcal{B}_j = Relax_{\Diamond_{\square}}(\square \varphi_j)$  be the universal Büchi automaton for  $\square \varphi_j$ , and let  $\mathcal{A}_j$  be a universal co-Büchi automaton for  $\square \diamondsuit \varphi$ . The constraint system for bound  $b \in \mathbb{N}$  has weights in  $\mathfrak{O}(n^2)$ . It has

$$O\Big((b^{2} + b \cdot |\mathcal{O}|) \cdot 2^{|\mathcal{I}|} + b \cdot |Q| \cdot (1 + \log(b \cdot |Q|)) + \sum_{j=1}^{n} (b \cdot |Q_{j}|(1 + \log(b \cdot |Q_{j}|))) + \sum_{j=1}^{n} (b \cdot |\widehat{Q}_{j}|(1 + \log(b \cdot |\widehat{Q}_{j}|)))\Big)$$

variables, and its size is

$$\mathcal{O}(|Q|^{2} \cdot b^{2} \cdot 2^{|\mathcal{I}|} \cdot (d + \log(b \cdot |Q|)) + \mathcal{O}(\sum_{j=1}^{n} (|Q_{j}|^{2} \cdot b^{2} \cdot 2^{|\mathcal{I}|} \cdot (d_{j} + r_{j} + \log(b \cdot |Q_{j}|))) + \mathcal{O}(\sum_{j=1}^{n} (|\widehat{Q}_{j}|^{2} \cdot b^{2} \cdot 2^{|\mathcal{I}|} \cdot (\widehat{d}_{j} + \log(b \cdot |\widehat{Q}_{j}|)))),$$

where

$$\begin{split} d &= \max_{s,q,\sigma_I,q'} |\delta_{s,q,\sigma_I,q'}|, \qquad \qquad d_j = \max_{s,q,\sigma_I,q'} |\delta_{s,q,\sigma_I,q'}^j|, \\ \widehat{d}_j &= \max_{s,q,\sigma_I,q'} |\widehat{\delta}_{s,q,\sigma_I,q'}^j|, \quad \text{and} \ r_j = \max_{s,q,\sigma_I,q'} |\text{rej}^j(s,q,q',\sigma_I)|. \end{split}$$

#### 5.3.5 A case study

Consider a robotic museum guide in a museum shown in Figure 5.12. The robot has to give a tour of the exhibitions in a specific order, which constitutes the hard specification. The tour starts at the entrance of the museum where the robot picks up newly arrived visitors. The main objective is to take the group through the two exhibitions on that floor and then return to the entrance to pick up a new group of people. Preferably, it also avoids certain locations, such as the library, or the passage when it is occupied. These preferences are encoded in the soft specifications. In particular, on one hand, the robot can only gain access to Exhibition 2 by getting a key from the staff's office. On the other hand, the robot is asked not to disturb the employees in the office. There is a library between Exhibition 1 and Exhibition 2 which can be used to go from one to the other, but it is preferred that visitors do not enter the library. However, it is also desirable that when the other passage between these two exhibitions is occupied, the robot does not go through there.

These specifications cannot be realized in conjunction. Given their priorities, we categorize the requirements into hard and soft specifications, and synthesize a strategy which satisfies the hard specifications and maximizes the satisfaction of the soft specifications. We formalize the problem as follows.

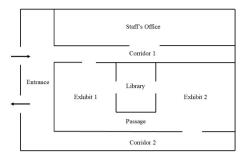


Figure 5.12: Map of the museum.

**Propositions:** The set  $\mathcal{I}$  contains a single Boolean variable occupied that indicates whether the passage between the two exhibitions is occupied. The set of output propositions  $\mathcal{O}$  consists of eight Boolean variables corresponding to the eight locations on the map: entrance, corridor<sub>1</sub>, corridor<sub>2</sub>, exhibition<sub>1</sub>, exhibition<sub>2</sub>, passage, office, library.

The hard specification is the conjunction of the following formulas.

• The robot starts at the entrance:

entrance.

• At each time step, the robot can occupy only one location:

$$\square \bigwedge_{o_1 \in \mathcal{O}} \left( o_1 \to \bigwedge_{o_2 \in \mathcal{O} \setminus \{o_1\}} \neg o_2 \right).$$

• The admissible actions of the robot are to stay in the current location or move to an adjacent one. This leads to eight requirements describing the map. For instance:

$$\square$$
 (corridor<sub>1</sub>  $\rightarrow$   $\bigcirc$  (corridor<sub>1</sub>  $\vee$  office  $\vee$  exhibition<sub>1</sub>)).

Remark: Due to the requirements above, the robot will always be in exactly one valid location, i.e., in a transition system that satisfies the

specifications it is impossible to reach a state where all output variables are false.

•	The robot	must	infinitely	often	visit	both	exhibitions:
---	-----------	------	------------	-------	-------	------	--------------

• The robot has to respect the order of visits, by starting from Exhibition 1, going to Exhibition 2 and finishing at the entrance:

$$\Box(exhibition_1 \to \bigcirc((\neg entrance \land \neg exhibition_1) \ \mathbb{U} \ exhibition_2)),$$
$$\Box(exhibition_2 \to \bigcirc((\neg exhibition_1 \land \neg exhibition_2) \ \mathbb{U} \ entrance)),$$
$$\Box(entrance \to \bigcirc((\neg exhibition_2 \land \neg entrance) \ \mathbb{U} \ exhibition_1)).$$

• The robot does not have access to Exhibition 2 before it visits the office:

 $\neg exhibition_2 \ U \ office.$ 

The set of soft specifications describes the desirable requirements that the robot does not enter the office, the library, or a occupied passage. Formally:

• The robot must not enter the office from corridor 1:

$$\square(corridor_1 \to \bigcirc \neg office)$$
.

• The robot must not enter the library from the exhibitions:

$$\square$$
 (exhibition<sub>1</sub>  $\vee$  exhibition<sub>2</sub>  $\rightarrow \bigcirc \neg library$ ).

• The robot must not enter the passage from the exhibitions when it is occupied:

```
\square ((exhibition<sub>1</sub> \vee exhibition<sub>2</sub>) \wedge \bigcirc occupied \rightarrow \bigcirc \neg passage).
```

The maximum realizability as iterative MaxSAT approach can be applied to synthesize a policy for the robotic museum guide. Table 5.1

summarizes the results. With implementation bound of 8, the hard specification is realizable and a partial satisfaction of soft specifications is achieved. This strategy always selects the passage to transition from Exhibition 1 to Exhibition 2 and hence, avoids the library. It also violates the requirement of not entering the staff's office, to acquire access to Exhibition 2. For implementation bound 10 the solver times out. Notice that strategies with higher values exists, however, they require larger implementation size.

**Table 5.1:** Results of applying synthesis with maximum realizability to the robotic navigation example, with different bounds on implementation size  $|\mathcal{T}|$ . We report on the number of variables and clauses in the encoding, the satisfiability of hard constraints, the value (and bound) of the MaxSAT objective function, the running times of Spot, Open-WBO, and the time of the solver plus the time for generating the encoding.

	Encoding		Solution		Time (s)		
T	# vars	# clauses	sat.	$\Sigma weights$	Spot	Open-WBO	enc.+solve
2	4051	25366	UNSAT	0 (39)	0.93	0.011	0.12
$^4$	19965	125224	UNSAT	0 (39)	0.93	0.079	0.57
6	45897	289798	UNSAT	0(39)	0.93	1.75	2.9
8	95617	596430	SAT	31 (39)	0.93	956	959
10	152949	954532	SAT	- (39)	0.93	time-out	time-out

# **Probabilistic Synthesis and Verification**

Autonomous systems operate in uncertain, dynamic environments and involve many sub-components such as perception, localization, planning, and control. The interaction between all of these components involves uncertainty. The sensors cannot entirely capture the environment around the autonomous system and are inherently noisy. Perception and localization techniques often rely on machine learning, and the outputs of these techniques involve uncertainty. Overall, the autonomous system needs to plan its decisions based on the uncertain output from perception and localization, which leads to uncertain outcomes.

To model decision-making under uncertainty, we start by describing the sequential decision-making model given in Figure 6.1. At any time step t, an agent observes the system's state  $s_t$  and makes a decision  $a_t$  based on this observation. This action yields two results. First, the agent receives an immediate reward  $r_t$  (or a cost.) Second, the system transitions to a new state  $s_{t+1}$  according to a probability distribution at time step t+1, which is determined by the choice of action  $a_t$  in

This section incorporates the results from the following publications (Cubuktepe, Jansen, Junges, Katoen, and Topcu, 2018; Cubuktepe, Jansen, Junges, Katoen, and Topcu, 2020b).

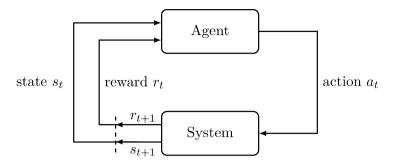


Figure 6.1: A representation of a sequential decision-making problem.

state  $s_t$ . At the subsequent time steps, the agent faces a similar problem and needs to make a decision in a potentially different state in the system and from a potentially different set of actions. We list the key ingredients in this sequential decision model:

- 1. A set of states, which describes all possible configurations of the system.
- 2. A set of available action in each state, which describes the set of decisions in the all states of the system.
- 3. A set of rewards or costs for each state and action, which describes the objective or performance criterion of the system.
- 4. A set of transition probabilities to the states in the system for each state and action, which describes the dynamics of the system.
- 5. A set of (potentially infinite) decision horizon, which describes the planning period.

We focus on a particular sequential decision-making model, which is called Markov decision processes (MDPs) (Puterman, 2014). In MDPs, the set of available actions, the rewards or costs, and the transition probabilities depend only on the current state and implemented action on the system.

**Definition 6.1** (Markov Decision Process (MDP)). A *Markov decision* process MDP is a tuple  $\mathcal{M} = (S, s_{init}, Act, \mathcal{P})$  with a finite set S of states,

an initial state  $s_{init} \in S$ , a finite set Act of actions, and a transition function  $\mathcal{P} \colon S \times Act \times S \to [0,1]$  such that  $\sum_{s' \in S} \mathcal{P}(s,Act,s') = 1$  for all  $s \in S$  and  $Act \in Act(s)$  where Act(s), denotes the set of available actions in the state s. A cost function  $c \colon S \times Act \to \mathbb{R}_{\geq 0}$  associates cost to state-action pairs.

MDPs have numerous applications in various domains thanks to their generality. These applications include but not limited to reinforcement learning (Jansen et al., 2020; Mason et al., 2017; Lecarpentier and Rachelson, 2019), robotics (Liu et al., 2018; Omidshafiei et al., 2017; Ghasemi and Topcu, 2019a), human-robot interaction (Chen et al., 2020; Akash et al., 2020; Cubuktepe et al., 2020a), aircraft collision avoidance (Julian et al., 2019), healthcare (Hosseini et al., 2014), disease management (Radoszycki et al., 2015), finance (Borkar and Jain, 2014), and digital marketing (Thomas et al., 2017).

The central problem for MDPs is to find a control policy, which determines what action to take with the current knowledge of the system at a given time. In other words, the policy is a mapping from the states to the actions. The typical aim is to optimize a given objective, such as minimizing expected cost, for example, minimizing the fuel usage of the system, or maximizing the probability of the successful operation of a system for a (potentially infinite) horizon. Given an MDP, the problem of finding an optimal policy can be cast as a dynamic programming problem, and numerous methods based on value or policy iteration and reinforcement learning exist to find such a policy.

A related problem of finding a policy in an MDP is called *model checking*. Given a model that represents the behavior of the system and a specification that determines the objective, model checking refers to a set of techniques that systematically check whether the system satisfies the given specification. For MDPs, a related technique is *probabilistic model checking*. Given an MDP and a specification, typically expressed as a formula in temporal logic, probabilistic model checking refers to determining whether there exists a policy that satisfies the specification. Many related problems such as minimizing the expected cost of satisfying the specification can also be solved using the probabilistic

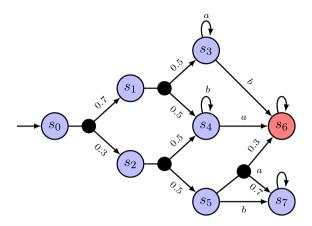
model checking framework.

In Section 6.1, we first consider the probabilistic model checking problem for MDPs and explain different approaches to solve this problem. We illustrate this problem with a case study on human-autonomy interactions in Section 6.2. Then, in Section 6.3, we define an extension of MDPs called parametric MDPs, where functions now give the transition and reward function of an MDP over parameters. We discuss a convex-optimization-based technique by Cubuktepe et al. (2018) to solve the so-called parameter synthesis problem that scales to thousands of parameters as opposed to a handful of parameters for the existing methods. Next, in Section 6.4, we consider a setting where the parameters of the transition probabilities and rewards belong to an uncertainty set parameterized by a collection of random variables. The problem is to compute the satisfaction probability to satisfy a temporal logic specification within any MDP that corresponds to a sample from these unknown distributions. We give a technique by Cubuktepe et al. (2020b) from so-called scenario optimization to compute high confidence bounds on the satisfaction probabilities.

# 6.1 Probabilistic Model Checking in Markov Decision Processes

Probabilistic model checking is a rigorous technique that can precisely solve the aforementioned control problems, and this rigor provides guarantees on appropriate behavior for all possible events in the system. Probabilistic model checking has been extensively studied for MDPs (Baier and Katoen, 2008), and mature tools exist for efficient model checking (Kwiatkowska *et al.*, 2011; Dehnert *et al.*, 2017; Hahn *et al.*, 2014).

In this section, we first give the formal definitions of concepts related to MDPs. Then, we give methods for verification and synthesis in MDPs subject to temporal logic specifications.



**Figure 6.2:** An MDP with the state space  $S = \{s_0, s_1, s_2, s_3, s_4, s_5, s_6, s_7\}$ , action space  $Act = \{a, b\}$  and the initial state  $s_{init} = s_0$ . The transition function  $\mathcal{P}$  is given by (possibly branching edges) in the graph. To avoid clutter, we omit the transitions with probability 1, and we merge action selections that induce the same transition probabilities, but differ only in action names.

#### 6.1.1 Preliminaries for Markov Decision Processes

A probability distribution over a finite or countably infinite set X is a function  $\mu \colon X \to [0, 1] \subseteq \mathbb{R}$  with  $\sum_{x \in X} \mu(x) = 1$ . The set of all distributions on X is denoted by Distr(X).

**Definition 6.2** (Policy). To define measures on MDPs, nondeterministic action choices are resolved by a so-called *policy*  $\sigma: S \to Act$  with  $\sigma(s) \in Act(s)$ . The set of all policies over  $\mathcal{M}$  is  $Pol^{\mathcal{M}}$ .

We note that the above definition of a policy is a deterministic mapping from states to actions. Such policies are called memoryless deterministic and suffice for MDPs for the performance criteria considered in this review (Baier and Katoen, 2008).

Applying a policy to an MDP yields an *induced Markov chain (MC)* where all nondeterminism is resolved.

**Definition 6.3** (Induced Markov Chain (MC)). For MDP  $\mathcal{M} = (S, s_{init}, Act, \mathcal{P})$  and policy  $\sigma \in Pol^{\mathcal{M}}$ , the MC induced by  $\mathcal{M}$  and  $\sigma$  is  $\mathcal{M}^{\sigma} = (S, s_{init}, Act, \mathcal{P}^{\sigma})$ 

with

$$\mathcal{P}^{\sigma}(s,s') = \mathcal{P}(s,\sigma(s),s') \text{ for all } s,s' \in S.$$

Intuitively, the transition probabilities in  $\mathcal{M}^{\sigma}$  are obtained with respect to the action choices of the policy.

**Definition 6.4** (Occupancy Measure). The occupancy measure  $x_{\sigma}$  of a policy  $\sigma$  for an MDP M is defined as

$$x_{\sigma}(s, Act) = \sum_{t=0}^{\infty} \Pr^{\sigma}(s_t = s, Act_t = Act | s_0 = s_{init}), \tag{6.1}$$

where  $\Pr^{\sigma}$  denotes the probability measure induced by  $\sigma$ , and  $s_t$  and  $Act_t$  denote the state and action in  $\mathcal{M}$  at time t.

The occupancy measure  $x_{\sigma}(s, Act)$  is the expected number of times to take action Act at state s under the policy  $\sigma$ .

For an MC  $\mathcal{D}$ , the reachability specification  $\varphi = \mathbb{P}_{\leq \lambda}(\lozenge T)$  asserts that a set  $T \subseteq S$  of target states is reached with probability at most  $\lambda \in [0,1]$ . If  $\varphi$  holds for  $\mathcal{D}$ , we write  $\mathcal{D} \models \varphi$ . Accordingly, for an expected cost specification,  $\psi = \mathrm{ER}_{\leq \kappa}(\lozenge G)$ ,  $\mathcal{D} \models \psi$  holds if and only if the expected cost of reaching a set  $G \subseteq S$  is bounded by  $\kappa \in \mathbb{R}$ . We use standard measures and definitions as in (Baier and Katoen, 2008, Ch. 10).

**Specifications.** We consider specifications that are combinations of reachability specifications and expected cost specifications. A reachability property  $\varphi = \mathbb{P}_{\leq \lambda}(\lozenge T)$  with upper probability bound  $\lambda \in [0, 1] \subseteq \mathbb{Q}$  and target set  $T \subseteq S$  constrains the probability to finally reach T from  $s_{init}$  in  $\mathbb{M}$  to be at most  $\lambda$ . Analogously, expected cost specifications  $\psi = \mathrm{ER}_{\leq \kappa}(\lozenge G)$  impose an upper bound  $\kappa \in \mathbb{R}$  on the expected cost to reach goal states  $G \subseteq S$  with respect to the cost function c. Combining both types of specifications, the intuition is that a set of bad states T shall only be reached with a certain probability  $\lambda$  (safety specification) while the expected cost for reaching a set of goal states G has to be below  $\kappa$  (performance specification). We overload the notation  $\lozenge T$  to denote both a reachability specifications and the set of all paths that

finally reach T from the initial state  $s_{init}$  of an MC. The probability and the expected cost for reaching T from  $s_{init}$  are denoted by  $\Pr(\lozenge T)$  and  $\operatorname{ER}(\lozenge T)$ , respectively. Hence,  $\Pr(\lozenge T) \leq \lambda$  and  $\operatorname{ER}(\lozenge G) \leq \kappa$  express that the properties  $\mathbb{P}_{\leq \lambda}(\lozenge T)$  and  $\operatorname{ER}_{\leq \kappa}(\lozenge G)$  respectively are satisfied by MC  $\mathcal{D}$ . We note that linear temporal logic (LTL) specifications can be reduced to reachability and expected cost specifications, and we refer the reader to (Baier and Katoen, 2008) for a detailed introduction.

An MDP  $\mathcal{M}$  satisfies both reachability specification  $\varphi$  and expected cost specification  $\psi$ , if and only if for all policies  $\sigma$  it holds that the induced MC  $\mathcal{M}^{\sigma}$  satisfies the properties  $\varphi$  and  $\psi$ , i.e.,  $\mathcal{M}^{\sigma} \models \varphi$  and  $\mathcal{M}^{\sigma} \models \psi$ . In our setting, we are also interested in the so-called synthesis problem, where the aim is to find a policy  $\sigma$  such that both specifications are satisfied (while this does not necessarily hold for all specifications). If  $\mathcal{M}^{\sigma} \models \varphi$ , policy  $\sigma$  is said to admit the property  $\varphi$ ; this is denoted by  $\sigma \models \varphi$ .

**Example 6.1.** Consider the MDP in Figure 6.2 with the target set  $T = \{s_6\}$ . Given a reachability specification  $\varphi = \mathbb{P}_{\leq \lambda}(\lozenge T)$  with  $\lambda = 0.85$ , an example policy  $\sigma_1$  that satisfies the specification is given by  $\sigma_1(s_3) = b, \sigma_1(s_4) = a$ , and  $\sigma_1(s_5) = a$ , which induces a reachability probability of 0.895, and is greater than the threshold  $\lambda$ . In this example, the policy  $\sigma_2$ , which is defined by selecting the same actions as  $\sigma_1$  except  $\sigma_2(s_5) = b$ , induces a MC that satisfies the specification with probability 0.85, and also satisfies the specification, even though it may not maximize the reachability probability.

# 6.1.2 Verification and Synthesis in Markov Decision Processes Subject to Temporal Logic Specifications

In this section, we describe the primal and dual linear programming (LP) formulations for verification and synthesis problems in MDPs subject to temporal logic specifications (Puterman, 2014; Forejt et al., 2011). We first start with the primal LP formulation, which is mainly used in verification problems. The variables in the primal LP formulation specify the probability of satisfying the reachability specification  $\varphi$  induced by a maximizing policy, hence it is used in verification problems. We

then describe the dual LP formulation, which is mainly used in policy synthesis problems subject to (multiple) temporal logic specifications. The variables in the dual LP are the expected number of times of taking action in the states of the MDP, and an optimal policy can also be obtained from an optimal solution of the dual LP.

Primal LP. The primal LP formulation computes the maximum probability  $p_{s_{init}}$  of reaching the target set T from the initial state  $s_{init}$ using Bellman's principle of optimality. The result of the LP serves as a certificate for the verification problem (Baier and Katoen, 2008, Ch. 10). The LP reads as follows:

minimize 
$$p_{s_{init}}$$
 (6.2)

subject to

$$p_s = 1, \quad \forall s \in T, \tag{6.3}$$

$$\mathcal{P}(s, Act, s') \ge 0, \quad \forall s, s' \in S \setminus T, \, \forall Act \in Act(s),$$
 (6.4)

$$\mathcal{P}(s, Act, s') \ge 0, \quad \forall s, s' \in S \setminus T, \, \forall Act \in Act(s), \qquad (6.4)$$

$$\sum_{s' \in S} \mathcal{P}(s, Act, s') = 1, \quad \forall s \in S \setminus T, \, \forall Act \in Act(s), \qquad (6.5)$$

$$p_s \ge \sum_{s' \in S} \mathcal{P}(s, Act, s') \cdot p_{s'} \quad \forall s \in S \setminus T, \, \forall Act \in Act(s), \tag{6.6}$$

$$\lambda \ge p_{s_{init}},\tag{6.7}$$

where  $p_s$  are the variables. For  $s \in S$ , the probability variable  $p_s \in [0,1]$ represents an upper bound of the probability of reaching target set  $T \subseteq S$ . We minimize  $p_{s_{init}}$  to compute the maximum probability of reaching set T. The constraint (6.7) ensures that the probability of reaching T is above the threshold  $\lambda$ . This constraint is optional for stating the verification problem, but with that constraint, the LP is feasible if all policies in the MDP satisfy the reachability specification  $\varphi$ .

The constraints of the primal LP have the following meanings. The probability of reaching a state in T from T is set to one (6.3). The constraints (6.4) and (6.7) ensure the validity of the transition probabilities and trivially hold for a valid MDP. For each state  $s \in S \setminus T$  and action  $Act \in Act(s)$ , the probability induced by the maximizing policy is a

lower bound to the probability variables  $p_s$  (6.6). The constraint (6.7) ensures that the probability of reaching T is below the threshold  $\lambda$ .

**Expected cost specifications.** The LP in (6.2) - (6.6) considers reachability probabilities. If we have instead an expected cost specification, we can similarly define the LP as follows:

$$\begin{array}{ll} \text{minimize} & p_{s_{init}} & (6.8) \\ \text{subject to} & \\ p_s = 0, & \forall s \in G, & (6.9) \\ \mathcal{P}(s, Act, s') \geq 0, & \forall s, s' \in S \setminus G, \, \forall Act \in Act(s), \\ & \\ \sum_{s' \in S} \mathcal{P}(s, Act, s') = 1, & \forall s \in S \setminus G, \, \forall Act \in Act(s), \\ & \\ \end{array} \tag{6.11}$$

$$p_s \ge c(s, Act) + \sum_{s' \in S} \mathcal{P}(s, Act, s') \cdot p_{s'}, \quad \forall s \in S \setminus G, \, \forall Act \in Act(s),$$

$$(6.12)$$

$$\kappa \ge p_{s_{init}}.
\tag{6.13}$$

We have  $p_s \in \mathbb{R}$ , as these variables represent the expected cost to reach G. At G, the expected cost is set to zero (6.9); The constraints (6.10) and (6.11) are analogous to (6.4) and (6.5). The actual expected cost for other states is a lower bound to  $p_s$  (6.12). Finally,  $p_{s_{init}}$  is bounded by the threshold  $\kappa$ .

**Dual LP.** In this section, we recall the dual LP formulation to compute a policy that maximizes the probability of satisfying a reachability specification  $\varphi$  in an MDP (Puterman, 2014; Forejt *et al.*, 2011).

The variables of the dual LP formulation are following:

•  $x_{\sigma}(s, Act) \in [0, \infty)$  for each state  $s \in S \setminus T$  and action  $Act \in Act$  defines the occupancy measure of a state-action pair for the policy  $\sigma$ , i.e., the expected number of times for taking action Act in state s.

•  $x_{\sigma}(s) \in [0,1]$  for each state  $s \in T$  defines the probability of reaching a state  $s \in T$ .

$$\text{maximize} \quad \sum_{s \in T} x_{\sigma}(s) \tag{6.14}$$

subject to

$$\sum_{Act \in Act} x_{\sigma}(s, Act) = \sum_{s' \in S \setminus T} \sum_{Act \in Act} \mathcal{P}(s', Act, s) x_{\sigma}(s', Act) + \alpha_s, \quad \forall s \in S \setminus T,$$
(6.15)

$$x_{\sigma}(s) = \sum_{s' \in S \setminus T} \sum_{Act \in Act} \mathcal{P}(s', Act, s) x_{\sigma}(s', Act) + \alpha_s, \quad \forall s \in T,$$
 (6.16)

$$\sum_{s \in T} x_{\sigma}(s) \ge \beta \tag{6.17}$$

where  $\alpha_s = 1$  if  $s = s_{init}$  and  $\alpha_s = 0$  if  $s \neq s_{init}$ . The constraints in (6.15) and (6.16) ensure that the expected number of times transitioning to a state  $s \in S$  is equal to the expected number of times to take action Act that transitions to a different state  $s' \in S$ . The constraint in (6.17) ensures that the specification  $\varphi$  is satisfied with a probability of at least  $\beta$ . We determine the states with probability 0 to reach T by a preprocessing step on the underlying graph of the MDP. To ensure that the variables  $x_{\sigma}(s)$  encode the actual probability of reaching a state  $s \in T$ , we then set the variables of the states with probability 0 to reach T to zero.

For any optimal solution  $x_{\sigma}$  to the LP in (6.14)–(6.17),

$$\sigma(s, Act) = \frac{x_{\sigma}(s, Act)}{\sum_{Act' \in Act} x_{\sigma}(s, Act')}$$
(6.18)

is an optimal policy, and  $x_{\sigma}$  is the occupancy measure of  $\sigma$ , see (Puterman, 2014) and (Forejt *et al.*, 2011) for details.

# 6.2 A Case Study

We now visit an example from (Feng  $et\ al.$ , 2015) as an illustration of synthesis in MDPs for autonomous systems that interact with human

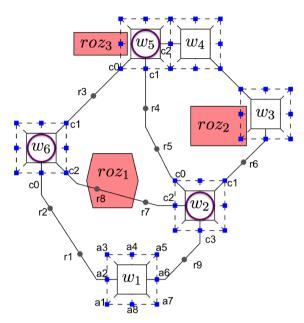


Figure 6.3: A road network for UAV ISR missions (adapted from (Humphrey *et al.*, 2014)).

operators, where a remotely controlled unmanned air vehicle (UAV) is used to perform intelligence, surveillance, and reconnaissance (ISR) missions over a road network. Figure 6.3 shows a map of the road network, which has six surveillance waypoints labeled  $w_1, w_2, \ldots, w_6$ . Approaching a waypoint from certain angles may be better than others, e.g. in order to obtain desired look angles on a waypoint target using an ellipsoidal loiter pattern. Angles of approach are thus discretized in increments of  $45^{\circ}$  around each waypoint, resulting in eight angle points  $a_1, a_2, \ldots, a_8$  around each waypoint. Roads connecting waypoints are discretized into road points  $r_1, r_2, \ldots, r_9$ . Red polygons represent "restricted operating zones" (ROZs), areas in which flying the UAV may be dangerous or lead to a higher chance of being detected by an adversary.

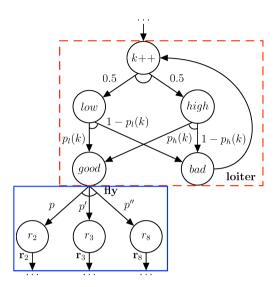
In current practice (Cooke and Pedersen, 2009), at least two human operators are required for a UAV ISR mission: one to pilot the UAV, and the other to steer the onboard sensor and interpret the sensor imagery. Here, we assume the UAV has a certain degree of autonomy that is used

to fulfill most of the piloting functions, e.g., maintaining loiter patterns around waypoints, selecting most of the points that comprise the route, and flying the route. The human operator primarily performs sensor tasks, e.g. steering the onboard sensor to capture imagery of targets at waypoints. However, the operator also retains the ability to affect some of the piloting functions of the UAV. The operator decides how many loiters to perform at each waypoint, since more loiters may be needed if the operator is not satisfied with the sensor imagery obtained on previous loiters. Additionally, waypoints  $w_2$ ,  $w_5$ , and  $w_6$  in Figure 6.3 will be designated as checkpoints. At checkpoints, the operator can directly impact the choices made by the protocol we synthesize by selecting different roads to be taken between waypoints.

The optimal piloting plan for the UAV varies depending on mission objectives. Specification patterns for a variety of UAV missions are presented in (Humphrey et al., 2014), including safety, reachability, coverage, sequencing of waypoints, etc. Through a few concrete examples, e.g. surveillance of the road network with minimum fuel consumption, or flying to certain waypoints while trying to avoid ROZs, the goal is to synthesize the optimal UAV piloting plan for a specific mission objective, which would be implemented by the UAV's onboard automation interface to control the route. In particular, the results shed light on how the uncertainties and imperfections of a human operator's behavior affect the optimal UAV piloting plan. Specifically, what is the influence of an operator's proficiency, workload, and fatigue level on UAV mission performance? Can we synthesize individualized optimal UAV piloting plans for different operators? Can the automation provide informative feedback to operators to assist them in decision-making?

We now introduce the models for the operator, the UAV and the interactions between the two.

The operator model. We build abstractions of the operator's possible behavior as a probabilistic model  $M_{\text{OP}}$ . Figure 6.4 shows a fragment of the model, representing the possible behavior at waypoint  $w_6$ . There is a non-negative integer variable k counting the number of sensor tasks performed by the operator since the beginning of the mission. The



**Figure 6.4:** A fragment of the operator model  $M_{\mathsf{OP}}$  representing the operator's behavior. The red dashed square highlights the common behavior that is repeated at all waypoints, while the blue solid square indicates specific choices of roads at waypoint  $w_6$ . The operator's behavior at other waypoints is omitted from the figure, indicated by  $\cdots$ 

updates "k++" represent increasing the value of k by one. The purpose of using k in the model is to measure the operator's fatigue level. To obtain a finite state model, let the value of k stop increasing once it reaches a certain threshold T (a constant that will be used later in modeling fatigue).

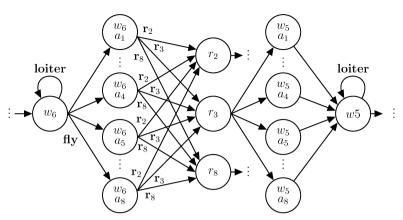
In general, operators' workload levels are driven by a number of factors including mission characteristics, e.g.how many UAVs the operator supervises simultaneously and the phase of the mission. For simplicity and to reduce the complexity of the models (so that the results discussed later are easier to interpret), we model the operator's workload as a uniform distribution over two levels: low and high. Operators' accuracy on vigilance tasks tends to decline with lower levels of proficiency and higher levels of workload (Boff and Lincoln, 1988). We model an oper-

ator's accuracy in steering sensors to capture high resolution imagery of targets as probability distributions that correlate with proficiency, workload, and fatigue. Specifically, when the operator's workload level is low, the probabilities of capturing good and bad quality imagery are  $p_l(k)$  and  $1 - p_l(k)$ , respectively. Here  $p_l(k)$  is a function over the variable k such that  $p_l(k) = p_l(0)$  if k < T and  $p_l(k) = f \cdot p_l(0)$  if k > T, where  $p_l(0)$  is the initial parameter value of the accuracy function, T is the fatigue threshold mentioned earlier, and f is a fatigue discount factor. We define the accuracy function  $p_h(k)$  for high workload analogously. Note that  $p_l(k) \geq p_h(k)$  for any k, modeling the fact that an operator tends to make more errors under higher levels of workload and stress. Furthermore, more proficient operators have higher values for the accuracy parameters  $p_l(0)$  and  $p_h(0)$ . If the quality of the captured imagery is bad, the operator would ask the UAV to continue to loiter at the current waypoint in order to collect more sensor imagery; otherwise, the operator allows the UAV to fly to another waypoint. At each waypoint, the operator repeats the aforementioned behavior (the red dashed square in Figure 6.4).

The operator selects the next road for the UAV at waypoints that are checkpoints ( $w_6$  for example), while the UAV controller chooses the road at any non-checkpoint waypoint. As illustrated in the blue solid square in Figure 6.4, we model the operator's choices at  $w_6$  as following a certain probability distribution, i.e., picking the roads that connect to neighboring road points  $r_2$ ,  $r_3$ , and  $r_8$  with probabilities p, p', and p'', respectively (note that p + p' + p'' = 1).

Suppose the operator chooses  $r_3$ . According to the map shown in Figure 6.3, the next waypoint is  $w_5$ . For simplicity, the operator's behavior at  $w_5$  is omitted from Figure 6.4.

**The UAV model.** We model the UAV's piloting behavior as an MDP  $M_{\text{UAV}}$ , which contains 63 states (6 waypoints,  $6 \times 8$  angle points, and 9 road points). At any waypoint or road point, the UAV can nondeterministically fly to a neighboring angle point or road point. These nondeterministic choices need to be resolved by a strategy. Figure 6.5



**Figure 6.5:** A fragment of the UAV model  $M_{\text{UAV}}$ , representing the UAV loitering and flying over the waypoints  $w_6$  and  $w_5$ .

shows a fragment of the UAV model<sup>1</sup>, illustrating how the UAV loiters and flies over waypoints  $w_6$  and  $w_5$ . If the UAV receives a loiter instruction from the operator, it loiters at the current waypoint, allowing the operator to capture more sensor imagery; otherwise, the UAV randomly picks one of the eight angle points  $a_1, \dots, a_8$  to exit  $w_6$ . Then, a nondeterministic choice between three roads  $r_2, r_3$ , and  $r_8$  needs to be resolved. Suppose  $r_3$  is chosen by the operator; then the UAV can fly to the waypoint  $w_5$  and approach it via one of the eight angles, or the UAV can also fly back to the waypoint  $w_6$  (for clarity, this choice is not drawn in Figure 6.5).

The operator-UAV interactions. We model interactions between the operator and the UAV by composing  $M_{\rm OP}$  and  $M_{\rm UAV}$ , which synchronize over common actions loiter and fly, and obtain a product MDP  $M_{\rm OP} \| M_{\rm UAV}$ , see (Baier and Katoen, 2008) for a formal product model definition. Synchronization between actions in our models abstracts the concrete process of interchanging information between the operator and the UAV, which is assumed via a reliable communication protocol. The model does not distinguish between "sender" and "receiver", or outputs and inputs, of a protocol. We can think of the operator initiating the

<sup>&</sup>lt;sup>1</sup>Our models are shown with several distributions associated with an action name but after composition this can easily be resolved through renaming, and we obtain MDPs.

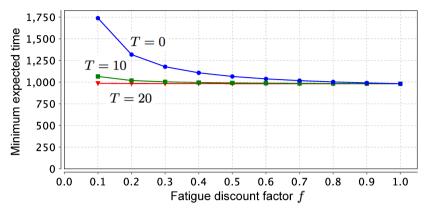
loiter action, which the UAV can receive at any waypoint (see the self-loops at  $w_6$  and  $w_5$  in Figure 6.5). The fly action is initiated as well by the operator, but we can think of the UAV deciding which flight direction to take (see the fly transitions at  $w_6$  in Figure 6.5). Note that synchronization between actions in the model also assumes that the operator and UAV synchronize their behaviors temporally.

Since  $M_{\text{OP}}$  is a discrete time MC and has no nondeterminism, synthesizing a strategy  $\sigma$  for the MDP  $M_{\text{OP}} \| M_{\text{UAV}}$  yields a strategy  $\sigma'$  for  $M_{\text{UAV}}$  such that  $(M_{\text{OP}} \| M_{\text{UAV}})^{\sigma} = M_{\text{OP}} \| M_{\text{UAV}}^{\sigma'}$ . The strategy  $\sigma'$  operates on the MDP model  $M_{\text{UAV}}$  of the UAV. For  $\sigma'$  to be implemented as a flight controller for the UAV, it must know which state of  $M_{\text{UAV}}$  the UAV is in, and hence, from the point of view of the strategy, the model transitions are triggered by the UAV's behavior.

**Representative analysis results.** We present representative results obtained using the above MDP model.

Consider a UAV surveillance mission that requires covering all six waypoints in Figure 6.3, where the objective is to complete the mission as fast as possible. Assume that each loiter takes 10 time units and flying between any neighboring waypoint and/or road point takes 60 time units. Figure 6.6 illustrates the influence of the operator's fatigue threshold T and discount factor f on the minimum expected time to complete the mission. The general trend is that the UAV completes the mission faster if the operator has a higher fatigue threshold T (i.e., less likely to get tired) or a larger value of f (i.e., the accuracy is less discounted). The best UAV performance (i.e., the smallest expected mission completion time) is achieved when f=1, that is, there is no accuracy discount due to fatigue.

The operator's accuracy in steering sensors and capturing good quality imagery are affected by proficiency and workload. Figure 6.7 illustrates the influence of accuracy parameters  $p_l(0)$  and  $p_h(0)$  on the minimum expected time of finishing the mission (i.e., covering all six waypoints). The trends show that a more proficient operator who has higher values for  $p_l(0)$  and  $p_h(0)$  can complete the mission faster. In addition, the



**Figure 6.6:** The effect of operator fatigue on minimum expected mission completion time, for different values of the fatigue threshold T and discount factor f (with fixed parameters  $p_l(0) = 0.9$  and  $p_h(0) = 0.8$ ).

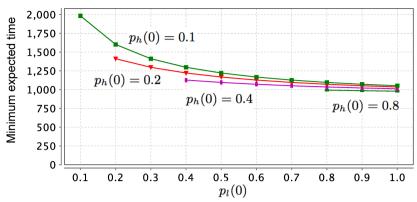
more accuracy declines due to high workload, i.e., the larger the gap between  $p_l(0)$  and  $p_h(0)$ , the longer the time needed to complete the mission.

## 6.3 Parameter Synthesis for Markov Decision Processes

In this section, we consider parametric Markov decision processes (parametric MDPs) whose transitions probabilities are affine functions of a finite set of parameters. Every fixed set of parameters induce an MDP, and the goal is to find a good set of parameters such that the induced MDP satisfies the specifications. We develop an approach that utilizes a sequential convex programming (SCP) method. The techniques improve the runtime and scalability by multiple orders of magnitude compared to black-box SCP by merging ideas from convex optimization and probabilistic model checking. We demonstrate the approaches on a satellite collision avoidance problem with hundreds of thousands of states and tens of thousands of parameters and their scalability on a wide range of commonly used benchmarks.

#### 6.3.1 Parametric Markov Decision Processes

A more general description of expressing the transition and reward function of an MDP is to define them as functions over parameters



**Figure 6.7:** The effect of operator proficiency and workload on minimum expected mission completion time, for different initial values of the accuracy functions  $p_l(0)$  and  $p_h(0)$  (with fixed parameters T = 10 and f = 0.7).

whose values are left unspecified (Daws, 2004; Lanotte et al., 2007; Hahn et al., 2010). Such parametric MDPs describe uncountable sets of MDPs. A well-defined instantiation of the parameters yields an instantiated, parameter-free MDP. For a given finite-state parametric MDP, the parameter synthesis problem is to compute a parameter instantiation such that the instantiated MDP satisfies a specification. We can also think the parameter synthesis problem as a design problem, where the objective is to compute an optimal design as a function of the parameters of the MDP. This design problem can also include computing an optimal policy for the parametric MDP.

Traditionally, approaches for solving the parameter synthesis problems have been built around the notion of abstracting the parametric model into a solution function. The solution function is the probability of satisfying the temporal logic specification as a function of the model parameters (Daws, 2004; Hahn et al., 2010; Dehnert et al., 2015; Filieri et al., 2016). The solution function can be exploited by the probabilistic model checking tools PARAM (Hahn et al., 2010), PRISM (Kwiatkowska et al., 2011) and Storm (Dehnert et al., 2017) to solve the parameter synthesis problem. However, this function is exponentially large in the parameters, and solving the problem is again exponential in the number of parameters, making the whole approach doubly exponential (Baier et al., 2020). Consequently, these approaches

typically can handle millions of states but only a *handful of parameters*. Moreover, these approaches require a fixed policy or has to introduce a parameter for every state/action-pair in the MDP.

Orthogonally, Quatmann et al. (2016) address an alternative parameter synthesis problem which focuses on proving the absence of parameter instantiations. The method iteratively solves simple stochastic games. Spel et al. (2019) consider proving that the parameters behave monotonically, allowing for faster sampling-based approaches. However, this method is limited to a few parameters. A recent survey on parameter synthesis in Markov models can be found in (Junges et al., 2019).

Further variations of parameter synthesis, e.g., consider statistical guarantees for parameter synthesis, often with some prior on the parameter values (Bortolussi and Silvetti, 2018; Calinescu *et al.*, 2016). These approaches cannot provide the absolute guarantees on an answer that the methods in this dissertation provide.

Parametric MDPs generalize interval models (Sen et al., 2006; Chen et al., 2013a). Such interval models have also been considered with convex uncertainties (Puggelli et al., 2013; Hahn et al., 2017; Wu and Koutsoukos, 2008; Chen et al., 2013a). However, the resulting problems with interval models are easier to solve due to the lack of dependencies (or couplings) between parameters in different states.

#### 6.3.2 Preliminaries for Parametric Markov Decision Processes

Let  $V = \{x_1, \ldots, x_n\}$  be a finite set of *variables* over the real numbers  $\mathbb{R}$ . The set of multivariate polynomials over V is  $\mathbb{Q}[V]$ . An *instantiation* for V is a function  $\mathbf{v} \colon V \to \mathbb{R}$ .

**Definition 6.5** ((Affine) parametric Markov decision process (pMDP)). A pMDP is a tuple  $\mathcal{M}_V = (S, s_{init}, Act, V, \mathcal{P})$  with a finite set S of states, an  $initial\ state\ s_{init} \in S$ , a finite set Act of actions, a finite set V of real-valued variables (parameters) and a  $transition\ function$  :  $S \times Act \times S \to \mathbb{Q}[V]$ . A pMDP is affine if (s, Act, s') is an affine function of V for every  $s, s' \in S$  and  $Act \in Act$ .

For  $s \in S$ ,  $Act(s) = \{Act \in Act \mid \exists s' \in S. (s, Act, s') \neq 0\}$  is the set of enabled actions at s. Without loss of generality, we require  $Act(s) \neq \emptyset$  for  $s \in S$ . If |Act(s)| = 1 for all  $s \in S$ , M is a parametric discrete-time Markov chain (pMC). MDPs can be equipped with a state-action cost function  $c: S \times Act \to \mathbb{R}_{>0}$ .

A pMDP  $\mathcal{M}$  is a Markov decision process (MDP) if the transition function yields well-defined probability distributions, i.e., :  $S \times Act \times S \rightarrow [0,1]$  and  $\sum_{s' \in S} (s, Act, s') = 1$  for all  $s \in S$  and  $Act \in Act(s)$ . Applying an instantiation  $\mathbf{v} \colon V \to \mathbb{R}$  to a pMDP  $\mathcal{M}$  yields an instantiated MDP  $\mathcal{M}[\mathbf{v}]$  by replacing each  $f \in \mathbb{Q}[V]$  in  $\mathcal{M}$  by  $f[\mathbf{v}]$ . An instantiation  $\mathbf{v}$  is well-defined for  $\mathcal{M}$  if the resulting model  $\mathcal{M}[\mathbf{v}]$  is an MDP.

## 6.3.3 Formal Problem Statement for Parameter Synthesis

In this section, we state the parameter synthesis problem, which is to compute a parameter instantiation such that the instantiated MDP satisfies the given temporal logic specification. We then discuss the nonlinear program formulation of the parameter synthesis problem, which forms the basis of the considered solution method.

**Problem 6.1** (Parameter Synthesis Problem). Given a parametric MDP  $\mathcal{M}_V = (S, s_{init}, Act, V, \mathcal{P})$ , and a reachability specification  $\varphi = \mathbb{P}_{\leq \lambda}(\lozenge T)$ , compute a well-defined instantiation  $\mathbf{v} \colon V \to \mathbb{R}$  for  $\mathcal{M}_V$  such that  $\mathcal{M}_V[\mathbf{v}] \models \varphi$ .

In words, we seek for an instantiation of the parameters that satisfies  $\varphi$  for all possible strategies in the instantiated MDP. We show necessary adaptions for an expected cost specification  $\psi = \text{ER}_{<\kappa}(\lozenge T)$  later.

For a given well-defined instantiation  $\mathbf{v}$ , Problem 6.1 can be solved by verifying whether  $\mathcal{M}_V[\mathbf{v}] \models \varphi$ . The standard formulation uses a linear program (LP) to minimize the probability  $p_{s_{init}}$  of reaching the target set T from the initial state  $s_{init}$  while ensuring that this probability is realizable under any strategy (Baier and Katoen, 2008, Ch. 10). The straightforward extension of this approach to pMDPs to *compute* a satisfiable instantiation  $\mathbf{v}$  yields the following nonlinear program (NLP) (Cubuktepe et al., 2017; Cubuktepe et al., 2018): with the variables  $p_s$  for  $s \in S$ , and the parameter variables in V in the transition function (s, Act, s') for  $s, s' \in S$  and  $Act \in Act(s)$ :

minimize 
$$p_{s_{init}}$$
 (6.19)

subject to

$$\forall s \in T, \quad p_s = 1, \tag{6.20}$$

$$\forall s, s' \in S \setminus T, \, \forall Act \in Act(s), \quad (s, Act, s') \ge 0, \tag{6.21}$$

$$\forall s \in S \setminus T, \, \forall Act \in Act(s), \quad \sum_{s' \in S} (s, Act, s') = 1, \tag{6.22}$$

$$\lambda \ge p_{s_{init}},$$
 (6.23)

$$\forall s \in S \setminus T, \, \forall Act \in Act(s), \quad p_s \ge \sum_{s' \in S} (s, Act, s') \cdot p_{s'}. \tag{6.24}$$

For  $s \in S$ , the probability variable  $p_s \in [0,1]$  represents an upper bound of the probability of reaching target set  $T \subseteq S$ . The parameters in the set V enter the NLP as part of the functions from  $\mathbb{Q}[V]$  in the transition function. The constraint (6.23) ensures that the probability of reaching T is below the threshold  $\lambda$ . This constraint is optional for stating the problem, but we use the constraint for finding a parameter instantiation that satisfies the specification  $\varphi$ . We minimize  $p_{s_{init}}$  to assign probability variables their minimal values with respect to the parameters V.

The probability of reaching a state in T from T is set to one (6.20). The constraints (6.21) and (6.22) ensure well-defined transition probabilities. Recall that (s, Act, s') is an affine function in V. Therefore, the constraints (6.21) and (6.22) only depend on the parameters in V, and they are affine in the parameters. Constraint (6.23) is optional but necessary later, and ensures that the probability of reaching T is below the threshold  $\lambda$ . For each state  $s \in S \setminus T$  and action  $Act \in Act(s)$ , the probability induced by the maximizing scheduler is a lower bound to the probability variables  $p_s$  (6.24). To assign probability variables to their minimal values with respect to the parameters in V,  $p_{s_{init}}$  is minimized in the objective (6.19). We state the correctness of the NLP in Proposition 6.1.

**Proposition 6.1.** The NLP in (6.19) – (6.24) computes the *minimal probability* of reaching T under a *maximizing* strategy, and an instantiation  $\mathbf{v}$  is feasible to the NLP if and only if  $\mathcal{M}_V[\mathbf{v}] \models \varphi$ .

*Proof.* The NLP in (6.19) – (6.24) is an extension of the LP in (Baier and Katoen, 2008, Theorem 10.105). We refer to (Junges, 2020, Theorem 4.20) for a formal proof.

**Remark 6.1** (Graph-Preserving Instantiations). In the LP formulation for MDPs, states with probability 0 to reach T are determined via a preprocessing on the underlying graph, and their probability variables are set to zero to ensure that the variables encode the actual reachability probabilities. This preprocessing requires the underlying graph of the parametric MDP to be preserved under any valuation of the parameters. Thus, as in (Hahn  $et\ al.,\ 2010$ ; Dehnert  $et\ al.,\ 2015$ ), we consider only graph-preserving valuations. Concretely, we exclude valuations  $\mathbf{v}$  with  $f[\mathbf{v}] = 0$  for  $f \in (s, Act, s')$  for all  $s, s' \in S$  and  $Act \in Act$ . We replace the set of constraints (6.21) by

$$\forall s, s' \in S. \, \forall Act \in Act(s), \quad (s, Act, s') \ge \varepsilon_{\text{graph}},$$
 (6.25)

where  $\varepsilon_{\text{graph}} > 0$  is a small constant.

We demonstrate the constraints for the NLP in (6.19) – (6.25) for a parametric MC by Example 6.2.

**Example 6.2.** Consider the parametric Markov chain in Fig. 6.8 with parameter set  $V = \{v\}$ , initial state  $s_0$ , and target set  $T = \{s_3\}$ . Let  $\lambda$  be an arbitrary constant. The NLP in (6.26) - (6.32) minimizes the

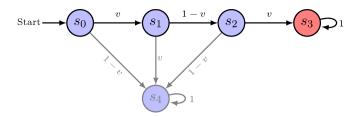


Figure 6.8: A parametric MC with a single parameter v.

probability of reaching  $s_3$  from the initial state:

$$minimize p_{s_0} (6.26)$$

subject to 
$$p_{s_3} = 1$$
, (6.27)

$$\lambda \ge p_{s_0},\tag{6.28}$$

$$p_{s_0} \ge v \cdot p_{s_1},\tag{6.29}$$

$$p_{s_1} \ge (1 - v) \cdot p_{s_2},\tag{6.30}$$

$$p_{s_2} \ge v \cdot p_{s_3},\tag{6.31}$$

$$1 - \varepsilon_{\text{graph}} \ge v \ge \varepsilon_{\text{graph}}.$$
 (6.32)

**Expected Cost Specifications.** The NLP in (6.19) – (6.25) considers reachability probabilities. If we have instead an expected cost specification  $\psi = \text{ER}_{\leq \kappa}(\lozenge G)$ , we replace (6.20), (6.23), and (6.24) in the NLP by the following constraints:

$$\forall s \in G, \, p_s = 0, \tag{6.33}$$

$$\forall s \in S \setminus G, \forall Act \in Act(s), p_s \ge c(s, Act) + \sum_{s' \in S} (s, Act, s') \cdot p_{s'},$$

(6.34)

$$\kappa \ge p_{s_{init}}.$$
(6.35)

We have  $p_s \in \mathbb{R}$ , as these variables represent the expected cost to reach G. At G, the expected cost is set to zero (6.33), and the actual expected cost for other states is a lower bound to  $p_s$  (6.34). Finally,  $p_{s_{init}}$  is bounded by the threshold  $\kappa$ .

For an affine parametric MDP  $\mathcal{M}_V$ , the functions in the resulting NLP (6.19) - (6.23) for parametric MDP synthesis are affine in V. However,

the functions in the constraints (6.24) are *quadratic*, as a result of multiplying affine functions occurring in with the probability variables  $p_{s'}$ . The problem in (6.19) - (6.24) is a quadratically constrained quadratic program (QCQP) (Boyd and Vandenberghe, 2004) and is generally nonconvex (Cubuktepe *et al.*, 2018), and computationally hard to solve. In the rest of the chapter, we discuss two methods to obtain a *locally optimal solution* to the problem in (6.19) - (6.24).

### 6.3.4 Sequential Convex Programming

In this section, we discuss a method by Cubuktepe et al. (2021a) for solving the parameter synthesis problem, which is a SCP approach with trust region constraints (Yuan, 2015; Mao et al., 2018; Chen et al., 2013b). The SCP method computes a locally optimal solution by iteratively approximating a nonconvex optimization problem. The resulting approximate convex problem is an LP, and the convexified functions are no longer upper bounds of the original functions. Therefore, approximation may generate optimal solutions in the convexified problem that are infeasible in the original problem. Therefore, we include trust regions and an additional model checking step to ensure that the new solution improves the objective. The trust regions ensure that the resulting LP accurately approximates the nonconvex QCQP. If the new solution indeed improves the objective, we accept and update the assignment of the variables and enlarge the trust region. Otherwise, we contract the trust region, and do not update the assignment of the variables.

## **Constructing the Affine Approximation**

We now explain in detail how we linearize the bilinear functions in the constraints in (6.24). Recall that this constraint appears as

$$\forall s \in S \setminus T, \, \forall Act \in Act(s), \quad p_s \ge \sum_{s' \in S} (s, Act, s') \cdot p_{s'}.$$

Similar to the previous seciton, consider the bilinear function in the above constraint

$$h(s, Act, s') = (s, Act, s') \cdot p_{s'}$$

$$(6.36)$$

and let

$$(s, Act, s') = 2d \cdot y + c, \text{ and } p_{s'} = z,$$
 (6.37)

where y is the parameter variable, z is the probability, and c, d are constants, similar to the previous chapter. We then convexify h(s, Act, s') as

$$h_{a}(s, Act, s') := 2d \cdot ((\hat{y} + \hat{z}) + \hat{y} \cdot (z - \hat{z}) + \hat{z} \cdot (y - \hat{y})) + c \cdot z, \quad (6.38)$$

where  $\langle \hat{y}, \hat{z} \rangle$  are any assignments to y and z. Note that the function  $h_{\mathbf{a}}(s, Act, s')$  is affine in the parameter variable y and the probability variable z.

After the linearization, the set of constraints (6.24) is replaced by the convex constraints

$$\forall s \in S \setminus T. \forall Act \in Act(s), \quad p_s \ge \sum_{s' \in S} h_a(s, Act, s').$$

**Remark 6.2.** If the parametric MDP is not affine, i.e., (s, Act, s') is not affine in V for every  $s, s' \in S$  and  $Act \in Act(s)$ , then h(s, Act, s') will not be a quadratic function in V and probability variables  $p'_s$ . In this case, we can compute  $h_a(s, Act, s')$  by computing a first order approximation with respect to V and  $p_{s'}$  around the previous assignment.

We use penalty variables  $k_s$  for all  $s \in S \setminus T$  to all linearized constraints, ensuring that they are always feasible. We minimize the sum of penalty variables to minimize the violation of the constraints in (6.45) However, the functions in these constraints do not over-approximate the functions in the original constraints. Therefore, a feasible solution to the linearized problem is potentially infeasible to the parameter synthesis problem. To make sure that the linearized problem accurately approximates the parameter synthesis problem, we use a trust region contraint around

the previous parameter instantiations. The resulting LP is:

minimize 
$$p_{s_{init}} + \tau \sum_{\forall s \in S \setminus T} k_s$$
 (6.39)

subject to

$$\forall s \in T, \quad p_s = 1, \tag{6.40}$$

$$\forall s, s' \in S \setminus T, \forall Act \in Act(s), \quad (s, Act, s') \ge \varepsilon_{\text{graph}},$$
 (6.41)

$$\forall s \in S, \forall Act \in Act(s), \quad \sum_{s' \in S} (s, Act, s') = 1, \tag{6.42}$$

$$\lambda \ge p_{s_{init}},$$
 (6.43)

$$\forall s \in S \setminus T, \forall Act \in Act(s), \quad k_s + p_s \ge \sum_{s' \in \backslash T} h_{\mathbf{a}}(s, Act, s'), \quad (6.44)$$

$$\forall s \in S \setminus T, \quad k_s \ge 0, \tag{6.45}$$

$$\forall s \in S \setminus T, \quad \hat{p}_s / \delta' \le p_s \le \hat{p}_s \cdot \delta', \tag{6.46}$$

 $\forall s, s' \in S \setminus T, \forall Act \in Act(s),$ 

$$(s, Act, s')/\delta' \le (s, Act, s') \le (s, Act, s') \cdot \delta',$$
 (6.47)

where  $\tau > 0$  is a constant, and (s, Act, s') and  $\hat{p}_s$  denotes the previous assignment for the parameter and probability variables. The constraints (6.46)–(6.47) are the trust region constraints.  $\delta > 0$  is the size of the trust region, and  $\delta' = \delta + 1$ . We demonstrate the linearization in Example 6.3.

**Example 6.3.** Recall the parametric Markov Chain (MC) in Fig. 6.8 and the QCQP from Example 6.2. After linearizing around an assignment for  $\hat{v}, \hat{p}_{s_0}, \hat{p}_{s_1}$ , and  $\hat{p}_{s_2}$ , the resulting LP with a trust region radius  $\delta > 0$ 

is given by

minimize 
$$p_{s_0} + \tau \sum_{i=0}^{2} k_{s_i}$$
  
subject to  $p_{s_3} = 1, \ \lambda \ge p_{s_0},$   
 $k_{s_0} + p_{s_0} \ge \hat{v} \cdot \hat{p}_{s_1} + \hat{p}_{s_1} \cdot (v - \hat{v}) + \hat{v} \cdot (p_{s_1} - \hat{p}_{s_1}),$   
 $k_{s_1} + p_{s_1} \ge p_{s_2} - \hat{v} \cdot \hat{p}_{s_2} - \hat{p}_{s_2} \cdot (v - \hat{v}) - \hat{v} \cdot (p_{s_2} - \hat{p}_{s_2}),$   
 $k_{s_2} + p_{s_2} \ge \hat{v} \cdot \hat{p}_{s_3} + \hat{p}_{s_3} \cdot (v - \hat{v}) + \hat{v} \cdot (p_{s_3} - \hat{p}_{s_3}),$   
 $k_{s_0} \ge 0, k_{s_1} \ge 0, k_{s_2} \ge 0,$   
 $\hat{p}_{s_0}/\delta' \ge p_{s_0} \ge \hat{p}_{s_0} \cdot \delta', \ \hat{p}_{s_1}/\delta' \ge p_{s_1} \ge \hat{p}_{s_1} \cdot \delta',$   
 $\hat{p}_{s_2}/\delta' > p_{s_2} > \hat{p}_{s_2} \cdot \delta', \ \hat{v}/\delta' > v > \hat{v} \cdot \delta'.$ 

We detail our SCP method in Fig. 6.9. We initialize the method with a guess for the parameters  $\hat{\mathbf{v}}$ , for the probability variables  $\hat{\mathbf{p}}$ , and the trust region  $\delta > 0$ . Then, we solve the LP (6.39)–(6.46) that is linearized around  $\hat{\mathbf{v}}$  and probability variables  $\hat{\mathbf{o}}$ .

After obtaining an instantiation to the parameters v, we model check the instantiated MDP  $\mathcal{M}_V[\mathbf{v}]$  to obtain the values of probability variables  $res(\mathbf{p})$  for the instantiation  $\mathbf{v}$ . If the instantiated MDP indeed satisfies the specification, we return the instantiation v. Otherwise, we check whether the probability of reaching the target set  $\beta$  is larger than the previous best value  $\hat{\beta}$ . If  $\beta$  is larger than  $\hat{\beta}$ , we update the values for the probability and the parameter variables, and enlarge the trust region. Else, we reduce the size of the trust region, and resolve the problem that is linearized around  $\mathbf{v}$  and  $\mathbf{p}$ . This procedure is repeated until a parameter instantiation that satisfies the specification is found, or the value of  $\delta$  is smaller than  $\omega > 0$ . The intuition behind enlarging the trust region is as follows: If the instantiation to the parameters  $\mathbf{v}$  increases the probability of reaching the target set  $\beta$  over the previous solution, then we conclude that the linearization is accurate. Consequently, the SCP method may take a larger step in the next iteration for faster convergence in practice.

For expected cost specifications, the resulting algorithm is similar,

except, we accept the parameter instantiation if the expected cost is reduced compared to the previous iteration, and initialize  $\hat{\beta}$  with a large constant. For various numerical results about the SCP Method, we refer the reader to (Cubuktepe *et al.*, 2018).

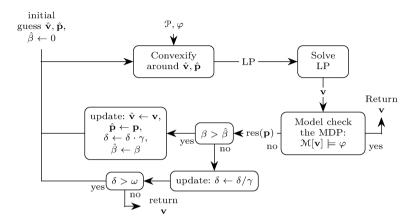


Figure 6.9: SCP with model checking in the loop. The NLP (6.19)–(6.24) is linearized around  $\hat{\mathbf{v}}$ ,  $\hat{\mathbf{p}}$ . Then, we solve the LP (6.39)–(6.46) an optimal solution to the parameter values, denoted by  $\hat{\mathbf{v}}$ . After each iteration, we model check the instantiated MDP  $M[\mathbf{v}]$  to determine whether the specification is satisfied. If the instantiated MDP satisfies the specification, we return the parameter instantiation. Otherwise, we check whether the reachability probability, denoted by  $\beta$ , is improved compared to the previous iteration, denoted by  $\hat{\beta}$ . If the probability is improved, we accept this step, update the assignment for the parameters and the probability variables, and increase the size of the trust region  $\delta$  by  $\gamma$ . Otherwise, we do not update the assignment, and decrease the size of the trust region.

## 6.4 Scenario-Based Verification in Uncertain Markov Decision Processes

In this chapter, we consider Markov decision processes (MDPs) in which the transition probabilities and rewards belong to an uncertainty set parametrized by a collection of random variables. The probability distributions for these random parameters are unknown. The problem is to compute the probability to satisfy a temporal logic specification within any MDP that corresponds to a sample from these unknown distributions. In general, this problem is undecidable, and we resort to techniques from so-called scenario optimization. Based on a finite

number of samples of the uncertain parameters, each of which induces an MDP, the proposed method estimates the probability of satisfying the specification by solving a finite-dimensional convex optimization problem. The number of samples required to obtain a high confidence on this estimate is independent from the number of states and the number of random parameters. Experiments on a large set of benchmarks show that a few thousand samples suffice to obtain high-quality confidence bounds with a high probability.

There are several approaches for verification of uncertain MDPs. Bacci et al. (2019) consider the analysis of Markov models in the presence of uncertain rewards, utilizing statistical methods to reason about the probability of a parametric MDP satisfying an expected cost specification. This approach is restricted to reward parameters and does not explicitly compute confidence bounds. Llerena et al. (2018) compute bounds on the long-run probability of satisfying a specification with probabilistic uncertainty for Markov chains. Other related techniques include multi-objective model checking to maximize the average performance with probabilistic uncertainty sets (Scheftelowitsch et al., 2017), sampling-based methods which minimize the regret with uncertainty sets (Ahmed et al., 2017), and Bayesian reasoning to compute parameter values that satisfy a metric temporal logic specification on a continuous-time Markov chain (Bortolussi and Silvetti, 2018). Arming et al. (2018) consider a variant of the problem in this dissertation where the probability distribution of the uncertainty sets is assumed to be known. This work formulates the policy synthesis problem as a partially observable Markov decision process (POMDP) synthesis problem and use off-the-shelf point-based POMDP methods (Pineau et al., 2003; Cassandra et al., 1997). The work in (Puggelli et al., 2013; Wolff et al., 2012) consider the verification of MDPs with convex uncertainties. However, the uncertainty sets for different states in an MDP are restricted to be independent, which does not hold in the considered problem setting where we have parameter dependencies.

Uncertainties in MDPs have received quite some attention in the artificial intelligence and planning literatures. Interval MDPs (Puggelli et al., 2013; Givan et al., 2000) use probability intervals in the transi-

tion probabilities. Dynamic programming, robust value iteration and robust policy iteration have been developed for MDPs with uncertain transition probabilities whose parameters are statistically independent, also referred to as rectangular, to find a policy ensuring the highest expected total reward at a given confidence level (Nilim and El Ghaoui, 2005; Wolff et al., 2012). The work in (Wiesemann et al., 2013) relaxes this independence assumption a bit and determines a policy that satisfies a given performance with a pre-defined confidence provided an observation history of the MDP is given by using conic programming. State-of-the art exact methods can handle models of up to a few hundred of states (Ho and Petrik, 2018). Multi-model MDPs (Steimle et al., 2018) treat distributions over probability and cost parameters and aim at finding a single policy maximizing a weighted value function. For deterministic policies this problem is NP-hard, and it is PSPACE-hard for history-dependent policies.

#### 6.4.1 Uncertain Markov Decision Processes

We now introduce the setting that we study in this chapter. Specifically, we use parameters to define the uncertainty in the transition probabilities and cost functions of an MDP. Each random parameter follows an unknown probability distribution from which we can sample the parameter values.

**Definition 6.6** (Uncertain Markov decision process (uMDP)). An uMDP  $\mathcal{M}_{\mathbb{P}}$  (uMDP) is a tuple  $\mathcal{M}_{\mathbb{P}} = (\mathcal{P}, \mathbb{P})$  where  $\mathcal{P}$  is a parametric MDP (pMDP), and  $\mathbb{P}$  is a probability distribution over the parameter space  $V_{\mathcal{P}}$ . If  $\mathcal{P}$  is a pMC, then we call  $\mathcal{M}_{\mathbb{P}}$  a uMC.

Intuitively, a uMDP is a pMDP with an associated distribution over possible (graph-preserving) parameter instantiations. That is, a realization of  $\mathbb{P}$  yields a concrete MDP  $\mathcal{P}[u]$  with the respective instantiation  $u \in V_{\mathbb{P}}$  (and  $\mathbb{P}(u) > 0$ ).

**Remark 6.3.** In a uMDP, we distinguish *controllable* and *uncontrollable* parameters. The uncontrollable parameters follow the probability

distribution  $\mathbb{P}$ . In contrast, we can actively *instantiate* the controllable parameters. In the following, we specifically allow cost parameters to be controllable.

**Definition 6.7** (Satisfaction Probability). Let  $\mathcal{M}_{\mathbb{P}} = (\mathcal{P}, \mathbb{P})$  be a uMDP and  $\varphi$  a specification. The (weighted) satisfaction probability of  $\varphi$  is

$$F(\mathfrak{M}_{\mathbb{P}},\varphi) = \int_{V_{\mathbb{P}}} I_{\varphi}(u) \ d\,\mathbb{P}(u)$$

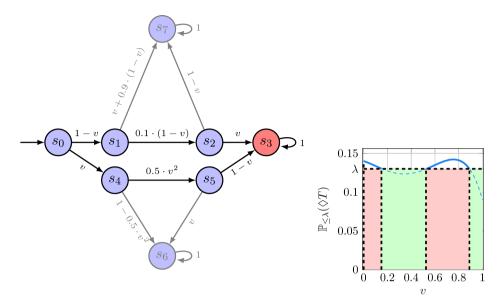
with  $u \in V_{\mathcal{P}}$  and  $I_{\varphi} \colon V_{\mathcal{P}} \to \{0,1\}$  is the indicator for  $\varphi$ , i.e.  $I_{\varphi}(u) = 1$  iff  $\mathcal{P}[u] \models \varphi$ .

Note that  $I_{\varphi}$  is measurable, as  $V_{\mathbb{P}}$  is the finite union of semi-algebraic sets (Basu *et al.*, 2010). Moreover, we have that  $F(\mathfrak{M}_{\mathbb{P}}, \varphi) \in [0, 1]$  and  $F(\mathfrak{M}_{\mathbb{P}}, \varphi) + F(\mathfrak{M}_{\mathbb{P}}, \neg \varphi) = 1$ .

**Example 6.4.** Consider the uMC in the left figure of Figure 6.10 with the uncontrollable parameter set  $V = \{v\}$ , initial state  $s_0$ , target set  $T = \{s_3\}$  and an uniform distribution for the parameter v over the interval [0,1]. We plot the probability of satisfying the specification  $\varphi = \mathbb{P}_{\leq \lambda}(\lozenge T)$  as a function of v in the right figure of Figure 6.10. We also show the satisfying region and its complementary as green and red regions. The satisfying region is given by the union of the intervals [0.13, 0.525] and [0.89, 1.0], and the satisfaction probability  $F(\mathfrak{M}_{\mathbb{P}}, \varphi_r)$  is 0.395 + 0.11 = 0.505.

#### 6.4.2 Formal Problem Statement

In this section, we state the problem that we study in this chapter. We seek to compute the satisfaction probability of the parameter space for a reachability or an expected cost specification  $\varphi$  on an uMDP. Intuitively, we seek the probability that a randomly sampled instantiation from the parameter space induces an MDP which satisfies  $\varphi$ . Formally: Given an uncertain MDP  $\mathcal{M}_{\mathbb{P}} = (\mathcal{P}, \mathbb{P})$ , and a specification  $\varphi$ , compute the satisfaction probability  $F(\mathcal{M}_{\mathbb{P}}, \varphi)$ . However, as mentioned, the problem is in general undecidable (Arming *et al.*, 2018). Therefore, we consider an approximation of computing the satisfaction probability:



**Figure 6.10:** Left: A uMC with parameter v. Right: The probability of satisfying the reachability specification  $\varphi = \mathbb{P}_{\leq \lambda}(\lozenge T)$  versus the value of the parameter v. Intervals that satisfy  $\varphi$  are green, intervals that violate  $\varphi$  are red.

**Problem 6.2.** Given an uncertain MDP  $\mathfrak{M}_{\mathbb{P}} = (\mathfrak{P}, \mathbb{P})$ , a reachability specification  $\varphi_r = \mathbb{P}_{\leq \lambda}(\Diamond T)$ , and a tolerance probability  $\nu$ , compute a confidence probability  $\alpha_{\nu}$  such that  $F(\mathfrak{M}_{\mathbb{P}}, \varphi_r) \geq 1 - \nu$  holds with a probability of at least  $1 - \alpha_{\nu}$ .

We illustrate the problem statement with the following example.

**Example 6.5.** For a UAV motion planning example, consider the question "What is the probability on a given day such that there exists a policy for the UAV to successfully finish the mission." A possible result is, e.g., 0.78 (confidence probability: 0.99) and 0.81 (confidence probability: 0.95). Then, with a confidence probability of 0.99, the actual satisfaction probability is indeed greater than 0.78, and with a (slightly lower) confidence probability of 0.95 it is greater than 0.81. Such a result shows that it is quite likely that the UAV will finish the mission successfully with a probability that is at least 81%.

#### 6.4.3 Scenario-Based Verification

In this section, we present an approach by Cubuktepe et al. (2020b) to solve Problem 6.2, that is, to approximate the satisfaction probability with respect to a specification. We first consider the robust verification problem that accounts for all possible values in the uncertainty set, potentially leading to a very pessimistic result. This problem can be formulated as a semi-infinite convex optimization problem, which is NP-hard (Wiesemann et al., 2013). Here, we exploit the structure of this problem, which includes finitely many variables but infinitely many constraints. The presented approach is based on scenario optimization (Calafore and Campi, 2006; Campi and Garatti, 2008): We sample a finite number of parameter values and restrict the semi-infinite problem to these samples. The resulting finite-dimensional convex optimization problem can be solved efficiently (Boyd and Vandenberghe, 2004). Based on the solution of the optimization problem, we compute high confidence in the estimate of the satisfaction probability. The estimate also generalizes to the samples from the probability distribution that are not in the sample set.

**Remark 6.4.** For ease of presentation, we focus on uncertain Markov chains (uncertain MCs). The results and methods generalize to uncertain MDPs (uncertain MDPs).

We first develop the main results for the simple setting where all sampled instantiated MCs from the parameter space  $V_{\mathcal{D}}$  satisfy the reachability specification  $\varphi_r$ . This assumption does not imply that all instantiated MCs satisfy  $\varphi_r$ : The sample set does not contain an MC that violates  $\varphi_r$  even though there exists such an MC in the parameter space. In Chapter 6.4.5, we drop this assumption and allow sampled points in  $V_{\mathcal{D}}$  to violate  $\varphi_r$ . This completes our treatment of Problem 6.2.

## 6.4.4 Restriction to Satisfying Samples

In this section, we assume that all instantiated MCs satisfy  $\varphi_r$ . We then generalize the presented method to any values of  $\nu$ . We want to check if an uncertain MC  $\mathcal{D}$  satisfies a reachability specification  $\varphi = \mathbb{P}_{\leq \lambda}(\lozenge T)$ 

for all instantiations in the sample set  $\mathcal{U}$ . For each instantiation, we can formulate a linear program (LP) that is feasible if and only if  $\varphi_r$  is satisfied (Puterman, 2014). For a subset  $\mathcal{U} \subseteq V_{\mathcal{D}}$  of the parameter space  $V_{\mathcal{D}}$  of the uncertain MC  $\mathcal{D}$ , we can then write the conjunction of these LPs. We assume that  $|\mathcal{U}|$  is finite and sampled from the probability distribution  $\mathbb{P}$  over the parameter space  $V_{\mathcal{D}}$ .

For each instantiation  $u \in \mathcal{U}$ , we introduce a set of linear constraints that are parametrized by u. We assume that each sample has a unique index. We use the following variables. For  $s \in S$  and  $u \in \mathcal{U}$ , the variable  $p_s^u \in [0,1]$  represents the probability of reaching the target set  $T \subseteq S$  from state s. The variable  $\tau$  represents an upper bound on the probability of satisfying  $\varphi_r$  for all instantiations in  $\mathcal{U}$ . Note that  $\tau$  is a variable in our formulation, whereas  $\lambda$  is the threshold of the reachability specification, and thus constant. The set  $\neg \exists \Diamond T$  represents the set of states which cannot reach the target set T. The probability of reaching T from these states is zero, and the set  $\neg \exists \Diamond T$  does not change for different graph-preserving instantiations (Hahn  $et \ al.$ , 2010). The set  $\neg \exists \Diamond T$  can be found in polynomial time in the size of an uncertain MC by using standard graph-based search algorithms (Baier and Katoen, 2008). We solve the following LP  $\mathcal{L}_r(\mathcal{U})$ , which is parametrized by each instantiation u in  $\mathcal{U}$ ,

minimize 
$$\tau$$
 (6.48)

subject to  $\forall u \in \mathcal{U}$ ,

$$p_{s_{init}}^u \le \tau, \tag{6.49}$$

$$p_{s_{init}}^u \le \lambda, \tag{6.50}$$

$$\forall s \in T, \quad p_s^u = 1, \tag{6.51}$$

$$\forall s \in \neg \exists \lozenge T, \quad p_s^u = 0, \tag{6.52}$$

$$\forall s \in S \setminus (T \cup \neg \exists \Diamond T), \quad p_s^u = \sum_{s' \in S} \mathcal{P}(s, s')[u] \cdot p_{s'}^u. \tag{6.53}$$

The objective (6.48) minimizes the maximal probability that can be achieved by all MCs induced by  $\mathcal{U}$ . The constraint (6.49) represents an upper bound on the reachability probability for all instantiations.

We minimize the upper bound to compute the maximal probability of satisfying  $\varphi_r$  for all instantiated MCs. The constraint (6.50) ensures that the probability of reaching T from the initial state  $s_{init}$  is below the threshold  $\lambda$ . The constraint (6.51) sets the probability to reach a state in T from T to 1. The constraint (6.52) sets the reachability probabilities from the states in  $\neg \exists \Diamond T$  to zero. The constraint (6.53) computes the probability of satisfying the specification for each non-target state  $s \in S$  in the standard way.

There are infinitely many constraints in the semi-infinite LP  $\mathcal{L}_r(V_{\mathcal{D}})$  as the cardinality of  $(V_{\mathcal{D}})$  is infinite. To obtain a LP with finitely many constraints, we instantiate the parameters  $u \in V_{\mathcal{D}}$  by sampling the probability distribution  $\mathbb{P}$ . Then, for a given violation probability  $\nu \in (0,1)$ , we compute a solution that violates the constraints in the LP  $\mathcal{L}_r(V_{\mathcal{D}})$  with a tolerance probability that is not larger than  $\nu$ . We first give some properties of the LP  $\mathcal{L}_r(\mathcal{U})$ . For proofs in this section, we refer the reader to (Cubuktepe et al., 2020b).

**Theorem 1.** Let uncertain MC  $\mathcal{D}$  and the sample sets  $\mathcal{U} \subseteq V_{\mathcal{D}}$  with  $K = |\mathcal{U}| \geq 2$ . Assume for all  $u \in \mathcal{U}$ ,  $\mathcal{D}[u] \models \varphi$ . For a given tolerance probability  $\nu \in [0,1)$ , let the associated confidence probability

$$\alpha_{\nu} = \sum_{i=0}^{1} {K \choose i} (1 - \nu)^{K-i} \nu^{i}. \tag{6.54}$$

Then, with a probability of at least  $1 - \alpha_{\nu}$ , we have

$$F(\mathcal{D}_{\mathbb{P}}, \varphi_r) \ge 1 - \nu. \tag{6.55}$$

Remark 6.5 (Independence to Model Size). The confidence probability in Theorem 1 is in fact independent from the number of states, transitions, or random parameters of the uncertain MC. From a practical perspective, the number of samples that are needed for a certain confidence does not depend on the model size.

Finally, Theorem 1 asserts that with a probability of at least  $1 - \alpha_{\nu}$ , the next sampled point from  $V_{\mathbb{D}}$  will satisfy the specification with

a probability of at least  $1 - \nu$ . Note that  $\alpha_{\nu}$  is the tail probability of a binomial distribution. It converges exponentially rapidly to 0 in  $|\mathcal{U}|$  (Campi and Garatti, 2008).

### 6.4.5 Satisfaction Probability by Treating Violating Samples

Theorem 1 assumes that all sampled points, that is, the induced MCs, satisfy the specification  $\varphi_r$ . This is a severe assumption in general. To lift this assumption, we consider the discarding approach from (Campi and Garatti, 2011). Specifically, after sampling a set of instantiations  $\mathcal{U}$  from  $V_{\mathcal{D}}$  according to the probability distribution  $\mathbb{P}$ , we remove the constraints for the MCs that violate the specification  $\varphi_r$  from the LP. We construct the set  $\mathcal{R} = \mathcal{U} \setminus \mathcal{Q}$ , where  $\mathcal{Q}$  denotes the set of samples that induce MCs violating the specification  $\varphi_r$ . Therefore, the set  $\mathcal{R}$  denotes the set of sampled MCs that satisfy the specification  $\varphi_r$ . We then solve the LP  $\mathcal{L}_r(\mathcal{R})$ 

minimize 
$$\tau$$
  
subject to  $\forall u \in \mathbb{R}$ , (6.56)  
 $(6.49) - (6.53)$ ,

where for  $u \in \mathcal{R}$  and  $s \in S$ ,  $p_s^u$  gives the probability of satisfying the reachability specification of the instantiated MC  $\mathcal{D}[u]$  at state s. The other constraints in the optimization problem in LP  $\mathcal{L}_r(\mathcal{R})$  are identical to the LP  $\mathcal{L}_r(\mathcal{U})$ . We give the main result of this section.

**Theorem 2.** Let uncertain MC  $\mathcal{D}$  and the sample sets  $\mathcal{U}, \mathcal{Q} \subseteq V_{\mathcal{D}}$ , with  $K = |\mathcal{U}| \geq 2$  and  $L = |\mathcal{Q}|$ . For a given tolerance probability  $\nu \in [0, 1)$ , the associated confidence probability is

$$\alpha_{\nu} = {L+1 \choose L} \sum_{i=0}^{L+1} {K \choose i} (1-\nu)^{K-i} \nu^{i}.$$
 (6.57)

Then, with a probability of at least  $1 - \alpha_{\nu}$ , we have

$$F(\mathcal{D}_{\mathbb{P}}, \varphi_r) \ge 1 - \nu. \tag{6.58}$$

## 6.4.6 Building Scenario-Based Algorithms

The question remains how we leverage the theoretical results to compute an estimate on the satisfaction probability to solve the problem in this section. For instance, let  $\nu$  be a violation probability and  $\mathcal{U}$  the sample set. Then, we can use Theorem 2 to compute the confidence probability  $\alpha_{\nu}$  by using the discarding approach from (Campi and Garatti, 2011). Similarly, for a the sample set  $\mathcal{U}$  and a threshold on the confidence probability  $\alpha_{\nu}$  we do a bisection on  $\nu$ . Specifically, we repeatedly apply Theorem 2 for different values of  $\nu \in (0,1)$ , to see if the corresponding confidence probability  $\alpha_{\nu}$  is below the threshold. We then approximate the lower and upper bounds on  $\nu$ .

The correctness of the approach is based on scenario-based optimization. However, it also applies to an obtained solution by any procedure (Campi et al., 2018). For instance, for any obtained value for the controlled parameters, we can construct a scenario program by sampling from random parameters. We can then apply Theorem 2 to compute the confidence probability  $\alpha_{\nu}$  or the violation probability  $\nu$ .

**Generalization to Uncertain MDPs.** Recall that we want to compute the satisfaction probability for an uncertain MDP, which is the probability that for any sampled MDP, we are able to synthesize a policy that satisfies the specification  $\varphi_r$ . To generalize the presented results to uncertain MDPs, we can modify the constraint (6.53) in the LP  $\mathcal{L}_r(\mathcal{U})$  as

$$\forall s \in S \setminus (T \cup \neg \exists \lozenge T), \ \forall Act \in Act(s), \quad p_s^u \leq \sum_{s' \in S} \mathcal{P}(s, Act, s')[u] \cdot p_{s'}^u. \tag{6.59}$$

The constraints (6.49)–(6.52) and (6.59) assert that, for each non-target state  $s \in S$  and action  $Act \in Act(s)$ , the probability induced by the minimizing policy is an upper bound to the probability variables  $p_s^u$ . Recall that, the reachability specification  $\varphi_r$  is satisfied if and only if the reachability probability at the initial state induced by the minimizing policy is less than  $\lambda$ . Then, our theoretical results apply to the uncertain MDPs.

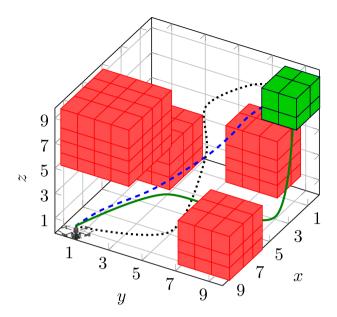


Figure 6.11: An example of a 3D UAV benchmark with obstacles and a target area.

## 6.4.7 Case Study: UAV Motion Planning

We implemented the approach from Chapter 6.4.3 using the model checker Storm (Dehnert et al., 2017) to construct and analyze samples of MDPs. To solve the scenario optimization problems with cost parameters, we used the SCS solver (O'Donoghue et al., 2016). All computations ran on a computer with 8 2.2 GHz cores, and 32 GB of RAM. We note that further benchmarks evaluating scenario-based algorithms with a varying number of samples in (Cubuktepe et al., 2020c; Badings et al., 2022). Specifically, the obtained confidence probabilities decrease exponentially rapidly with an increasing number of samples.

In the benchmark, we consider a UAV motion planning example to model a realistic problem with a high number of random parameters. We model the problem as an uncertain MDP, where the parameters represent how the weather conditions affect the movement of the UAV, and how the weather may change. In particular, different wind conditions induce specific satisfaction probabilities. We assume that the planning

area is a certain valley where we have historic weather data which provide distributions over parameter values. The mission of the UAV is to transport a payload to a specific location and return safely to its initial position. The problem is to compute the satisfaction probability, that is, the probability that for any sampled MDP for this scenario we are able to synthesize a UAV policy that satisfies the specification.

We model the problem as follows: States encode the position of the UAV, the current weather situation, and the general wind direction in the valley. Parameters describe how the weather affects the position of the UAV for different zones in the valley, and how the weather/wind may change during the day. Fig. 6.11 shows an example environment with zones to avoid (red) and a target zone (green). We define four different weather conditions that each induce certain probability distributions over the eight different wind directions. The parameters of the model determine the probabilities of transitioning between different weather and wind conditions at each time step. The specification is to reach the target zone safely with a probability of at least 0.9. The number of states in our example is 266 880, and the number of parameters is 2500.

For the distributions over parameter values, that is, over weather conditions, we consider the following cases. First, we assume a uniform distribution over the different weather conditions in each zone. Second, the probability for a weather condition inducing a wind direction that pushes the UAV into the positive y-direction is five times more likely than others. Similarly, in the third case, it is five times more likely to push the UAV into the negative x-direction. We depict some example trajectories of the UAV for three different conditions in Fig. 6.11. The trajectory given by the blue dashed line represents the expected trajectory for the first case, taking a direct route to reach the target area. Similarly, the trajectories given by the black dotted and solid green lines represent the expected trajectories for the second and third cases. For the second case, we observe that the UAV tries to avoid to get closer to the obstacles in x direction as the wind may push the UAV to the obstacles. For the third case, the UAV avoids the obstacle at the bottom and then reaches the target area.

We sample 1 000 parameters for each case and approximate the maximal satisfaction probability with a confidence probability of at least  $1 - \alpha_{\nu}$ , with  $\alpha_{\nu} = 10^{-6}$ . The highest satisfaction probability is given by the first weather condition with 0.86, and the other conditions have a satisfaction probability of 0.78 and 0.75, showing that it may be harder to navigate around the obstacles with non-uniform probability distributions. The average time to compute the satisfaction probabilities is 1 341 seconds.

Finally, we introduce costs to a 2-dimensional example, where hitting an obstacle causes (1) a cost of 100 and (2) the UAV to return to the initial position. Specifically, we introduce cost parameters for transitions that steer the UAV towards x or y-directions. We minimize the maximal possible expected cost (under all parameter values) to reach the target location. The specification asserts that the resulting expected cost should be less than 20.

We uniformly sample 1000 parameter values for weather conditions and note that the UAV policies favor on average transitioning to y-direction more compared to the x-direction to minimize the cost while ensuring that the probability of hitting an obstacle is minimized. The average expected cost of the induced MDPs is 7.41 and the satisfaction probability is 0.71. The solving time for this example is 2274 seconds.

## 7

## **Dealing with Information Limitations**

Methods for the synthesis and verification of policies in Markov decision processes (MDPs), parametric Markov decision processes (pMDPs), and uncertain Markov decision processes (uMDPs) assume that the agent is able to observe the underlying state of the system. However, many sensor or communication limitations may lead to imperfect or limited observations of the system's state in practice. In this section, we study the problem of sequential decision-making under uncertainty when the decision-making agent has only limited observational capabilities.

There exist several formalisms for modeling decision-making under imperfect perception. Bai et al. (2014) model an integrated perception and planning problem using partially observable Markov decision process (POMDP)s. In Ghasemi and Topcu (2019b), the authors propose a perception-aware point-based POMDP solver. In Benenson et al. (2006), the authors integrate simultaneous localization and mapping (SLAM) with a partial motion planner for autonomous navigation. Fu et al. (2016) consider a robot with a temporal logic task in a probabilistic map obtained from a semantic SLAM.

We begin this section by presenting the POMDP, a commonly used

model for decision-making under partial observability. We then present uncertain partially observable Markov decision processs (uPOMDPs)—in addition to having imperfect information about the current state, the transition and observation functions belong to uncertainty sets. Following the exposition of Cubuktepe et al. (2021b), we present an algorithm for the synthesis of policies that robustly satisfy specifications in uncertain POMDPs. We then present an application of uPOMDPs through a spacecraft motion planning case study. Finally, as an example of policy synthesis under an alternate model of partial observability we present an algorithm, and corresponding case studies, for the setting in which an MDP model of the system is available, but the semantic labeling of the environment is only partially known.

### 7.1 Partially Observable Markov Decision Processes

POMDPs generalize MDPs by introducing an *observation function*, which defines a probability distribution over a set of possible observations given the current state of the system.

**Definition 7.1** (POMDP). A POMDP is a tuple  $\mathfrak{M}_Z = (\mathfrak{M}, Z, O)$ , with  $\mathfrak{M}$  the underlying MDP of  $\mathfrak{M}_Z$ , a finite set of observations Z, and observation function  $O \colon S \to Z$ .

For brevity, we use so-called deterministic observation functions which may be derived from the more standard stochastic observation functions  $O: S \to Distr(Z)$  via a (polynomial) reduction (Chatterjee *et al.*, 2016). For POMDPs, observation-action sequences are based on a finite path  $\pi \in \mathsf{Paths}^{\mathcal{M}}_{fin}$  of the underlying MDP  $\mathcal{M}$  and have the form:  $\pi_o = O(\pi) = O(s_0) \xrightarrow{Act_0} O(s_1) \xrightarrow{Act_1} \cdots O(s_n)$ . The set of all finite observation-action sequences for a POMDP  $\mathcal{M}_Z$  is  $\mathsf{ObsSeq}^{\mathcal{M}_Z}_{fin}$ .

While the agent acts within the environment, it encounters certain observations, according to which it can infer the probability of the system being in a certain state. Technically, this *belief* b is a distribution  $b \in Distr(S)$ , such that b(s) describes the probability of being in state  $s \in S$ .

Recall that a policy in an MDP is a mapping from states to actions. A policy in a POMDP is a mapping from the observations (or a history of observations) to the actions.

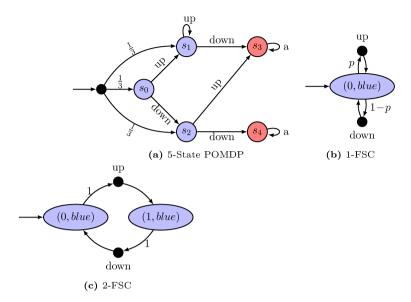
The problem of computing an optimal policy for POMDPs is undecidable (Madani *et al.*, 1999). To achieve computational tractability, policies are often restricted to finite memory by computing a finite state controller (FSC) (Meuleau *et al.*, 1999; Amato *et al.*, 2010), which is an NP-hard problem (Vlassis *et al.*, 2012).

**Definition 7.2** (Finite-Memory Policies.). An observation-based policy  $\sigma: (Z \times Act)^* \times Z \to Distr(Act)$  for a POMDP maps a trace, i.e., a sequence of observations and actions, to a distribution over actions. An FSC consists of a finite set of memory states and two functions. The action mapping  $\gamma(n,0)$  takes an FSC memory state n and an observation 0, and returns a distribution over POMDP actions. To change a memory state, the memory update  $\eta(n,0,Act)$  returns a distribution over memory states and depends on the action Act selected by  $\gamma$ . An FSC induces an observation-based policy by following a joint execution of these functions upon a trace of the POMDP. An FSC is memoryless if there is a single memory state; memoryless FSCs encode policies  $\sigma: Z \to Distr(Act)$ .

We start with a simple 5-state reachability example to highlight the utility of FSCs as finite-memory POMDP policies.

**Example 7.1.** Consider the POMDP in Figure 7.1. The POMDP has three observations (*blue*,  $s_3$  and  $s_4$ ), where observation *blue* is received upon visiting  $s_0$ ,  $s_1$ , and  $s_2$ . That is, the agent is unable to distinguish between these states. The specification is  $\varphi = \Pr_{\geq 0.9}(\diamondsuit s_3)$ ; the agent should reach state  $s_3$  with at least probability 0.9. We define the 1-FSC  $\mathcal{A}_1$ , illustrated in Figure 7.1b, with one memory node 0:

$$\alpha(0,blue) = \begin{cases} up & \text{with probability} \quad p,\\ down & \text{with probability} \quad 1-p,\\ \delta(0,0,Act) = 0 \quad \forall 0 \in Z, Act \in Act. \end{cases}$$



**Figure 7.1:** (a) POMDP for Example 7.1 with three observations  $\{blue, s_3, s_4\}$  with (b) 1-FSC and (c) 2-FSC. Both FSCs are defined for observing "blue" and subsequent action choices that may result in a change of memory node for the 2-FSC.

A 2-FSC with two memory nodes (0 and 1), see Figure 7.1c, allows for greater expressivity, i.e. the policy can base its decision on larger observation sequences.

With this memory structure, we can create an FSC  $A_2$  that ensures the satisfaction of  $\varphi$ :

$$\alpha(0,blue) = \begin{cases} up & \text{with probability} & 1, \\ down & \text{with probability} & 0, \end{cases}$$
 
$$\alpha(1,blue) = \begin{cases} up & \text{with probability} & 0, \\ down & \text{with probability} & 1, \end{cases}$$
 
$$\delta(0,blue,up) = 1,$$
 
$$\delta(1,blue,down) = 0.$$

While POMDPs provide an expressive modeling framework for partially observable settings, in many cases the transition and observation

functions of the model may not be known exactly. In these cases, it is necessary to also model our uncertainty in the POMDP model. We discuss such uPOMDPs models in the following section.

### 7.2 Uncertain Partially Observable Markov Decision Processes

By combining uMDPs and POMDPs, we obtain uPOMDPs—in addition to imperfect information pertaining to the current state, the transition and observation functions belong to uncertainty sets.

**Definition 7.3** (UPOMDP). An uPOMDP is a tuple  $\mathfrak{M}_{Z,\mathfrak{P}}=(S,\,s_{init},\,Act,\,\mathbb{I},\,\mathfrak{P},\,Z,\,O,\,r)$  with a finite set S of states, an initial state  $s_{init}\in S$ , a finite set Act of actions, a set  $\mathbb{I}$  of probability intervals, an uncertain  $transition function <math>\mathfrak{P}\colon S\times Act\times S\to \mathbb{I}$ , a finite set Z of observations, an  $uncertain\ observation\ function\ O\colon S\times Z\to \mathbb{I}$ , and a reward function  $R\colon S\times Act\to \mathbb{R}_{>0}$ .

Nominal probabilities are point intervals where the upper and lower bounds coincide. If the probabilities of all transitions and observations are nominal, the model is a (nominal) POMDP. Without a loss of generality, we may express any uPOMDP as a set of nominal POMDPs that vary only in their transition functions. For a transition function  $P: S \times Act \times S \to \mathbb{R}$ , we write  $P \in \mathcal{P}$  if for all  $s, s' \in S$  and  $\alpha \in Act$  we have  $P(s, Act, s') \in \mathcal{P}(s, Act, s')$  and  $P(s, Act, \cdot)$  is a probability distribution over S. Finally, we note that a fully observable uPOMDP where each state has unique observations is an uMDP.

Existing approaches for policy synthesis in uPOMDPs rely on dynamic programming (Wolff et al., 2012), convex optimization (Wolff et al., 2012), or value iteration (Wolff et al., 2012). While the complexity of solving a standard MDP is polynomial in the number of states and actions, solving an uMDP is NP-hard in general (Wiesemann et al., 2013). The existing approaches for uPOMDPs rely on sampling (Burns and Brock, 2007) or robust value iteration (Osogami, 2015) on the belief space of the uPOMDP. The policy synthesis algorithm presented by Suilen et al. (2020) is based on convex optimization and searches over

memoryless policies. However, their resulting optimization problems are exponentially larger than ours, and they only consider memoryless policies. Meanwhile, Burns and Brock, 2007 utilize sampling-based methods and Osogami, 2015 employ a robust value iteration on the belief space of the uPOMDP. Ahmadi et al., 2018 use sum-of-squares optimization to verify uncertain POMDPs against temporal logic specifications. Itoh and Nakamura, 2007 assume distributions over the probability values of the uncertainty set. Finally, Chamie and Mostafa, 2018 consider a convexified belief space and computes a policy that is robust over this space.

# 7.2.1 Synthesizing Robust Finite-State Controllers for Uncertain Partially Observable Markov Decision Processes

Following the presentation of Cubuktepe *et al.*, 2021b, we now present an algorithm for the synthesis of finite-state controllers that robustly satisfy specifications in uPOMDPs.

#### **Problem Formulation**

We begin by introducing observation-based policies, which are similar to memoryless FSCs for POMDPs. We then introduce specifications for POMDPs, followed by the notion of robustly satisfying a specification in a uPOMDP.

**Definition 7.4** (Observation-based policy). An observation-based policy  $\sigma \colon Z \to Distr(Act)$  for an uPOMDP maps observations to distributions over actions. Note that such a policy is referred to as memoryless and randomized. More general types of policies take an (in)finite sequence of observations and actions into account. We use  $\Sigma^{\mathcal{M}_{Z,\mathcal{P}}}$  to denote the set of observation-based strategies for  $\mathcal{M}_{Z,\mathcal{P}}$ . Applying  $\sigma \in \Sigma^{\mathcal{M}_{Z,\mathcal{P}}}$  to  $\mathcal{M}_{Z,\mathcal{P}}$  resolves all choices and partial observability and results in an induced (uncertain) Markov chain  $\mathcal{M}_{Z,\mathcal{P}}^{\sigma}$ .

**Specifications in POMDPs.** We constrain the undiscounted expected reward (the value) of a policy for a POMDP using *specifications*: For a POMDP  $\mathcal{M}_Z$  and a set of goal states G the specification  $\mathrm{ER}_{<\kappa}(\lozenge G)$ 

states that the expected reward before reaching G is at least  $\kappa$ . For brevity, we require that the POMDP has no dead-ends, i.e., that under every policy, we eventually reach G. Reachability specifications to a subset of G and discounted rewards are special cases (Puterman, 2014).

Robustly Satisfying Specifications in uPOMDPs. A policy  $\sigma$  satisfies a specification  $\varphi = \operatorname{ER}_{\leq \kappa}(\lozenge G)$ , if the expected reward to reach G induced by  $\sigma$  is at least  $\kappa$ . POMDP  $\mathcal{M}_{Z,\mathcal{T}}[P]$  denotes the *instantiated* uPOMDP  $\mathcal{M}_{Z,\mathcal{T}}$  with a fixed transition function  $P \in \mathcal{T}$ . A policy robustly satisfies  $\varphi$  for the uPOMDP  $\mathcal{M}_{Z,\mathcal{T}}$ , if it does so for all  $\mathcal{M}_{Z,\mathcal{T}}[P]$  with  $P \in \mathcal{T}$ . Thus, a (robust) policy for uPOMDPs accounts for all possible instantiations  $P \in \mathcal{T}$ .

Intuitively, to robustly satisfy a specification in a uPOMDP, we require a policy that satisfies the specification for all POMDP instantiations from  $\mathcal{M}_{Z,\mathcal{P}}[P]$ . If we have several (expected cost or reachability) specifications  $\varphi_1, \ldots, \varphi_m$ , we write  $\sigma \models \varphi_1 \land \ldots \land \varphi_n$  where  $\sigma$  robustly satisfies all specifications. Note that general temporal logic constraints can be reduced to reachability specifications (Baier and Katoen, 2008; Bouton et al., 2020), therefore we omit a detailed introduction to the underlying logic.

**Assumption 7.1.** For the correctness of the presented method, we require the lower bounds of the intervals to be strictly larger than zero, that is, an instantiation cannot "eliminate" transitions. Put differently, either a transition exists in *all* instantiations of the uMDP, or in none. That assumption is standard and, for instance, also employed in (Wiesemann *et al.*, 2013). Moreover, the problem statement would be different and theoretically harder to solve, see (Winkler *et al.*, 2019). We allow the upper and lower bound of an interval to be the same, resulting in *nominal* transition probabilities.

**Problem 7.1** (Robust Synthesis for Uncertain POMDPs). Given an uncertain uPOMDP  $\mathcal{M}_{Z,\mathcal{P}}$  and an expected reward specification  $\varphi = \mathrm{ER}_{\leq \kappa}(\Diamond G)$ , compute an FSC that yields an observation-based policy  $\sigma$  which robustly satisfies  $\varphi$  for  $\mathcal{M}_{Z,\mathcal{P}}$ .

# Optimization Problem for Uncertain Partially Observable Markov Decision Processes

We now reformulate the above problem statement as a semi-infinite nonconvex optimization problem, with finitely many variables but infinitely many constraints.

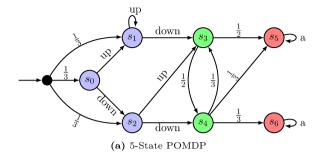
To do so, we begin by adopting a small extension of (PO)MDPs in which only a subset of the actions are available in any given state, i.e., the transition function should be interpreted as a partial function. We denote the set of actions at state s by Act(s). We ensure that any states that share an observation also share the set of available actions. Moreover, we translate the observation function to be deterministic without uncertainty, i.e., of the form  $O \colon S \to Z$ , by expanding the state space (Chatterjee  $et\ al.$ , 2016).

**Definition 7.5** (Binary/Simple uncertain POMDP). An uncertain POMDP is *binary*, if  $|Act(s)| \leq 2$  for all  $s \in S$ . A binary uncertain POMDP is *simple* if for all  $s \in S$ , the following is true:

$$|Act(s)| = 2 \text{ implies } \forall Act \in Act(s). \exists s' \in S, \ \mathcal{P}(s, Act, s') = 1,$$

Simple uPOMDPs differentiate the states with *action* choices and *uncertain* outcomes. All uPOMDPs can be transformed into simple uPOMDPs. We refer to (Junges *et al.*, 2018) for a transformation. This transformation preserves the optimal expected reward of an uPOMDP. We denote  $S_{\rm a}$  by the states with action choices, and  $S_{\rm u}$  the states with uncertain outcomes. We now give an example to a simple and binary POMDP in Figure 7.2.

We now introduce the optimization problem with the nonnegative reward variables  $\{r_s \geq 0 \mid s \in S\}$  denoting the expected reward before reaching goal set G from state s, and positive variables  $\{\sigma_{s,Act} > 0 \mid s \in S, Act \in Act(s)\}$  denoting the probability of taking an action Act in a state s for the policy. Note that we only consider policies where for all states s and actions Act it holds that  $\sigma_{s,Act} > 0$ , such that applying the policy



**Figure 7.2:** A simple and binary POMDP with four observations  $\{blue, green, s_3, s_4\}$ . This POMDP is simple as the states  $\{s_0, s_1, s_2\}$  only have the action choices in  $\{up, down\}$ , and the actions have deterministic outcomes. On the other hand, the states  $\{s_3, s_4\}$  only have *uncertain* outcomes with probabilistic transitions belonging to intervals.

to the uncertain POMDP does not change the underlying graph.

maximize 
$$r_{s_I}$$
 (7.1)

subject to 
$$r_{s_I} \ge \kappa$$
, (7.2)

$$\forall s \in G, \quad r_s = 0, \tag{7.3}$$

$$\forall s \in S, \quad \sum_{Act \in Act(s)} \sigma_{s,Act} = 1,$$
 (7.4)

$$\forall s, s' \in S, \forall Act \in Act(s), \quad O(s) = O(s') \implies \sigma_{s,Act} = \sigma_{s',Act}, \quad (7.5)$$

$$\forall s \in S_{\mathrm{u}}. \forall P \in \mathcal{P}, \quad r_s \leq R(s) + \sum_{s' \in S} P(s, s') \cdot r_{s'}. \tag{7.6}$$

$$\forall s \in S_{a}, \quad r_{s} \leq \sum_{Act \in Act(s)} \sigma_{s,Act} \cdot \left( R(s,Act) + \sum_{s' \in S} \mathcal{P}(s,Act,s') \cdot r_{s'} \right),$$

$$(7.7)$$

The objective is to maximize the expected reward  $r_{s_I}$  at the initial state. The constraint (7.3) encodes the specification requirement and assigns the expected reward to 0 in the states of goal set G. We ensure that the policy is a valid probability distribution in each state by (7.4). Next, (7.5) ensures that the policy is observation-based. We encode the computation of expected rewards for states with uncertain outcomes by (7.6) and with action choices by (7.7). We omit denoting the unique actions in the transition function P(s, s') and reward function R(s)

in (7.6) for states with uncertain outcomes.

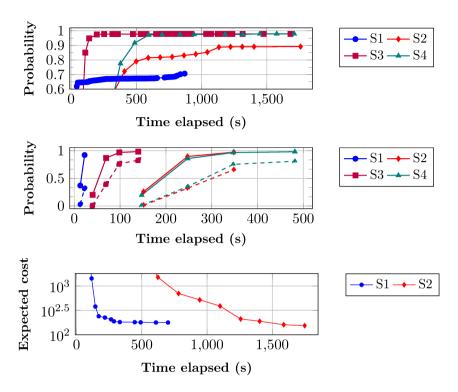
Let us consider some properties of the optimization problem. First, the functions in (7.7) are quadratic. Essentially, the policy variables  $\sigma_{s,Act}$  are multiplied with the reward variables  $r_s$ . In general, these constraints are nonconvex, and we later linearize them. Second, the values of the transition probabilities P(s,s') for  $s,s' \in S_u$  in (7.6) belong to continuous intervals. Therefore, there are infinitely many constraints over a finite set of reward variables. These constraints are similar to the LP formulation for MDPs (Puterman, 2014), and are affine; there are no policy variables. We refer the reader to (Cubuktepe  $et\ al.$ , 2021b) for a solution approach of the optimization problem.

## 7.2.2 Spacecraft Motion Planning Case Study

We evaluate the sequential convex programming (SCP)-based approach from (Cubuktepe *et al.*, 2021b) for solving the uPOMDP problem on a case study based on a satellite motion planning problem.

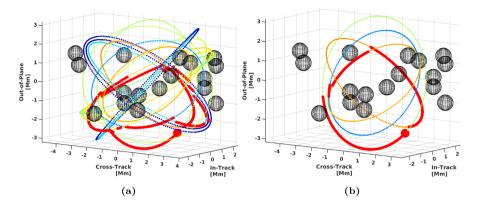
This case study considers the robust spacecraft motion planning system presented by Frey et al. (2017) and Hobbs and Feron (2020). The spacecraft orbits the earth along a set of predefined natural motion trajectories (NMTs) (Kim et al., 2007). While the spacecraft follows its current NMT, it does not consume fuel. Upon an imminent close encounter with other objects in space, the spacecraft may be directed to switch into a nearby NMT at the cost of a certain fuel usage. We consider two objectives: (1) To maximize the probability of avoiding a close encounter with other objects and (2) to minimize the fuel consumption, both within successfully finishing an orbiting cycle. Uncertainty enters the problem in the form of potential sensing and actuating errors. In particular, there is uncertainty about the spacecraft position, the location of other objects in the current orbit, and the failure rate of switching to a nearby NMT.

Modeling Spacecraft Motion Planning Using uPOMDPs. We encode the problem as an uPOMDP with two-dimensional state variables for the NMT  $n \in \{1, ..., 36\}$ , and the (discretized) time index  $i \in \{1, ..., I\}$ 



**Figure 7.3:** Computational effort versus the performance of the different policies for the spacecraft motion planning case study. Top: Obtained probability of avoiding close encounters between the spacecraft and other objects in the orbit. Middle: The performance of policies synthesizes using a nominal POMDP model, applied to the nominal model (solid lines) and to the uPOMDP model (dashed lines). Bottom: The obtained expected cost of successfully finishing an orbit.

for a fixed NMT. We use different values of resolution I in the examples. Every combination of  $\langle n,i\rangle$  defines an associated point in the 3-dimensional space. The transition model is built as follows. In each state, there are two types of actions, (1) staying in the current NMT, which increments the time index by 1, and (2) switching to a different NMT if two locations are close to each other. More precisely, we allow a transfer between two points in space defined by  $\langle n,i\rangle$  and  $\langle n',i'\rangle$  if the distance between the two associated points in space is less than 250km. A switching action may fail. In particular, the spacecraft may transfer to an unintended nearby orbit. The observation model contains 1440



**Figure 7.4:** We show the obtained trajectory from a policy in red that finishes an orbit around the origin. (a) Trajectory from a memoryless policy. (b) Trajectory from policy with 5 memory states. We highlight the initial location by a big red circle. We depict the other objects by black spheres, and all NMTs that were used as a part of the trajectory in blue, green, or yellow.

different observations of the current locations of the orbit, which give partial information about the current NMT and time index in orbit. Specifically, for each NMT, we can only observe the location up to an accuracy of 40 different locations in each orbit. The points that are close to each other have the same observation.

**Experimental Setup and Algorithm Variants.** We consider three benchmarks. S1 is our standard model with a discretized trajectory of I=200 time indices. S2 denotes an extension of S1 where the considered policy is an FSC with 5 memory states. S3 uses a higher resolution (I=600). Finally, S4 is an extension of S3, where the policy is an FSC with 2 memory states. In all models, we use an uncertain probability of switching into an intended orbit and correctly locating the objects, given by the intervals [0.50, 0.95]. The four benchmarks have approximately 3.6e4, 3.5e5, 1.1e5 and 3.4e5 states as well as 6.5e4, 7.0e5, 2.0e5 and 6.7e5 transitions, respectively. In this example, the objective is to maximize the probability of avoiding a close encounter with objects in the orbit while successfully finishing a cycle.

Memory Yields the Best Policies. Fig. 7.3 (top) shows the convergence of the reachability probabilities for each model, specifically the probability versus the runtime in seconds. First, we observe that after 20 minutes of computation, using larger POMDPs that have a higher resolution or memory in the policy yields superior policies. Second, the policy with memory is superior to the policy without memory. Finally, we observe that larger models indeed require more time per iteration, which means that on the smaller uncertain POMDP S1, the algorithm converges faster to its optimum.

Comparing the Policy Performances. Fig. 7.4 shows a comparison of policies and depicts the resulting spacecraft trajectories. In particular, we show the trajectories from two different policies, a memoryless policy (one memory state) in Fig. 7.4a—computed on S1—and from a policy with finite memory (five memory states) in Fig. 7.4b— computed on S2. The trajectory from the memoryless policy switches its orbit 17 times, whereas the trajectory from the finite-memory policy switches its orbit only 4 times. Additionally, the average length of the trajectory with the finite-memory policy is 188, compared to 375 for the memoryless one. These results demonstrate the utility of finite-memory to improve the reachability probability and to minimize the number of switches.

Robust Policies are Indeed More Robust. We demonstrate the utility of computing robust policies against uncertainty in Fig. 7.3 (middle). Intuitively, we compute policies on nominal models and use them on uncertain models. We give results on the nominal transition probabilities of the four considered models. The performance of the policies on the nominal models has solid lines, and the performance of the policies on the uncertain models has dashed lines. Note that when we apply the policies synthesized using nominal models on the uPOMDP, they perform significantly worse and fail to satisfy the objective in each case. The results clearly demonstrate the trade-offs between computing policies for nominal and uncertain problems. In each case, the computation time for the problem with uncertainty is roughly an order of magnitude larger. Yet, the resulting policies are better: we observe that the probability

of a close encounter with another object increases up to 60 percentage points, if we do not consider the uncertainty in the model.

**Expected Energy Consumption.** Finally, we consider an example where there is a cost for switching orbits. In this example we only consider models S1 and S2, which both have discretized trajectories with I = 200 time indices. The objective is to minimize the cost of successfully finishing a cycle in orbit. We obtain the cost of each switching action according to the parameters in (Frey et al., 2017). Additionally, we define the cost of a close encounter to be  $10000N \cdot s$ , a much higher cost than all possible switching actions. We reduce the uncertainty in the model by setting the previously mentioned intervals for these models to [0.90, 0.95]. In particular, the worst-case probability to correctly detect objects is now higher than before, reducing the probability of close encounters with those objects. Fig. 7.3 shows the convergence of the costs for each model. The costs of the resulting policies for models S1 and S2 are  $178N \cdot s$  and  $153N \cdot s$ , respectively. Similar to the previous example, the results demonstrate the utility of finite-memory policies in reducing fuel costs in spacecraft motion planning problems.

#### 7.3 Task-Oriented Active Perception and Planning

In the previous subsections we presented FSCs for POMDPs, as well as an approach to compute robust policies for uPOMDPs. In general, POMDPs provide a useful modeling framework for policy synthesis when perception and planning are closely coupled. However, in some cases the uncertainty stemming from perception-based information limitations must be decoupled from the the uncertainty that arises from the stochasticity of the underlying system dynamics. For instance, in the case of probabilistic knowledge over atomic propositions, a POMDPs model would require an exponential expansion in its state space, resulting in a computationally intractable problem.

Following the presentation of Ghasemi *et al.* (2020), we now present a model and an accompanying algorithm that is able to separately reason about these two sources of uncertainty. More specifically, we consider

an agent that is assigned a temporal logic task in an environment that is only partially known. We represent the environment using semantic labeling, i.e., with a set of state properties (labels) captured by atomic propositions over which the agent holds a probabilistic belief that is updated as new sensory measurements arrive. The goal is to design a joint perception and planning strategy for the agent that realizes the task with high probability. We present a planning strategy that takes the semantic uncertainties into account and by doing so provides probabilistic guarantees on the task success. Furthermore, as new data arrives, the belief over the atomic propositions evolves and the planning strategy adapts accordingly.

Temporal logic planning under imperfect perception has been studied from different perspectives. Jones et al. (2013) proposed a new type of logic, called distribution temporal logic, to enable expressions over belief-based predicates. In Ding et al. (2011), the authors considered an agent moving over a graph where the truth values of the predicates over the nodes depend on known probabilities. There is also a family of solutions that rely on sampling methods (Vasile and Belta, 2013). In another related work, Silva et al. (2019) propose a synthesis algorithm for probabilistic temporal logic over reals specifications in the belief space. Montana et al. (2017) propose a sampling-based solution to temporal logic planning under imperfect state information, relying on constructing a transition system by sampling from a feedback-based information roadmap. The work of Kress-Gazit et al. (2009) considers uncertainty in the environment propositions and proposed to design a reactive controller offline such that it can satisfy the task for all admissible environments.

Many solutions resort to replanning techniques, such as the iterative receding-horizon planning algorithm by Wongpiromsarn *et al.* (2012) mentioned in Section 5.2. For a subclass of temporal logic formulas, Livingston *et al.* (2012) introduce a method to locally patch a nominal controller once a change in the environment is detected. Lahijanian *et al.* (2016) propose an iterative replanning strategy that can relax the constraints imposed by the task if the discovery of a new obstacle deems the task unrealizable. Fu and Topcu (2015) and Fu and Topcu (2016)

design an alternating active sensing and planning strategy for temporal logic tasks. In the work of Guo *et al.* (2013) the agent generates a plan according to its initial knowledge over a deterministic model and after new real-time information is gathered, it revises the plan.

#### 7.3.1 Problem Formulation

Returning to the presentation of Ghasemi *et al.* (2020), we now introduce the agent model, the environment model, the observation model and the task specification language that we use in the formal problem statement.

#### **Agent and Environment Models**

We model the interaction between the agent and the environment by a MDP. In particular, we use  $\mathcal{M} = (S, s_{init}, Act, \mathcal{P})$  to represent the MDP modeling the agent's decision-making.

The agent perceives its environment through atomic propositions. Different from the approaches discussed above, we use a time-varying labeling function to capture the belief of the agent about the truth values of the atomic propositions.

**Definition 7.6.** An environment model is a tuple  $\mathcal{E} = (S, \mathcal{AP}, \bar{\mathcal{L}})$  where S is a finite, discrete state space,  $\mathcal{AP}$  is a set of atomic propositions, and  $\bar{\mathcal{L}}: S \to 2^{\mathcal{AP}}$  is a deterministic labeling function that captures the true state of the environment. The agent's belief at time t is a probabilistic labeling function  $\mathcal{L}_t: S \times 2^{\mathcal{AP}} \to [0,1]$ . For a state s and a subset of atomic propositions  $P \subseteq \mathcal{AP}, \mathcal{L}_t(s,P)$  assigns the probability of the event that P holds true at s, i.e.,  $\Pr(\bigcap_{p \in P} s \models p)$ .

We denote by  $\mathcal{L}_0$  the agent's prior belief. This prior belief may, for example, be an uninformative prior distribution. We assume that the truth values of the propositions are mutually independent in each state, facilitating the update of the labeling function over time. Nevertheless, so long as the joint distribution model is known, the updates can be computed.

#### **Observation Model**

At each time step, the agent's perception module processes a set of sensory measurements regarding the atomic propositions. While the measurements may be from multiple sensing units, for ease of notation, we consider their joint model by a general observation function.

**Definition 7.7.** Let  $\mathcal{Z}(s_1, s_2, p) \in \{ \mathit{True}, \mathit{False} \}$  denote the perception output of the agent at state  $s_1$  for the atomic proposition p at state  $s_2$ . The joint observation model of the agent is

$$O: S \times S \times \mathcal{AP} \times \{\mathit{True}, \mathit{False}\} \rightarrow [0, 1],$$

where  $\mathcal{O}(s_1, s_2, p, b)$  represents the probability that  $\mathcal{Z}(s_1, s_2, p) = True$  if the truth value of p is given by the Boolean variable b.

More specifically,  $\mathcal{Z}(s_1, s_2, p)$  represents a measurement of the property captured by p at state  $s_2$ , i.e., whether p holds at  $s_2$  or not, when the agent takes this measurement from state  $s_1$ , i.e., the current state of the agent is  $s_1$ . Hence,  $\mathcal{Z}(s_1, s_2, p)$  is a Bernoulli random variable and its distribution is dictated by the observation model  $\mathcal{O}(s_1, s_2, p, b)$ , which depends on the true value of p at state  $s_2$  captured by p, i.e., p indicates whether p in reality holds at p or not. An accurate observation model is one for which the output probability of  $\mathcal{O}(s_1, s_2, p, b)$  is one for p and zero for p and zero

In the Bayesian framework, the observation model is used to update the agent's belief. Nevertheless, in the absence of such an observation model, one can perform the update in a frequentist way.

#### Task Specification Language

We use syntactically co-safe linear temporal logic (scLTL) (Kupferman and Vardi, 2001) to specify finite-horizon tasks for the agent. Notice that scLTL is a variant of linear temporal logic (LTL), introduced in Section 2.2, that deals only with finite-horizon tasks. Hence, one would need to investigate finite traces of a system for verification purposes. The language defined by an scLTL formula can equivalently be represented

through a deterministic finite automaton (DFA)  $\mathcal{D} = (\mathcal{Q}, q_{init}, \Sigma, \delta, \mathcal{F})$  (Kupferman and Vardi, 2001).

#### **Problem Statement**

We consider an agent whose interaction with the environment is captured by an MDP. The agent is tasked with an scLTL specification that can be successfully completed within finite time. The agent is unaware of the environment state. However, it starts with a prior knowledge and over time, gathers observations that can be used to further revise its belief. The formal definition of the problem is stated next.

**Problem 7.2.** Given an MDP  $\mathcal{M} = (S, s_{init}, Act, \mathcal{P})$ , an environment model with unknown labeling function  $\mathcal{E} = (S, \mathcal{AP}, \_)$ , an observation model  $\mathcal{O}$ , and an scLTL task  $\varphi$ , find a policy  $\pi$  that maximizes the probability of satisfying the task conditioned on the true labeling function, i.e.,

$$\pi^* = \operatorname*{arg\,max}_{\pi} \Pr(\mathcal{M}^{\pi} \models \varphi \,|\, \bar{\mathcal{L}}).$$

#### 7.3.2 Joint Active Perception and Planning

Before introducing the algorithm, it is necessary to first describe the challenges of having a probabilistic view of the environment and, in particular, atomic propositions. In a setting where the agent is uncertain about the valuations of all atomic propositions over the whole environment, there may be up to  $2^{|S||\mathcal{AP}|}$  possibilities for how the environment is configured. In this case, computing policies that can account for all possible configurations, as offline reactive synthesis does (Baier and Katoen, 2008), is indeed computationally intractable. Additionally, if the environment is not dynamically changing, then such a comprehensive policy that accounts for all possible configurations, is not necessary. Another fundamental challenge is the fact that the nature of the probabilistic perception differs from the stochasticity of the agent model. Therefore, as seen in Example 7.2, one cannot combine the belief probabilities on the perception side with the transition probabilities of the MDP.

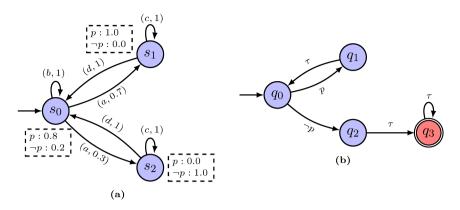


Figure 7.5: The MDP transitions behave in a stochastic way while the uncertainty in knowledge does not. (a) An MDP. The edge labels represent an action and a transition probability, respectively, while the node labels capture agent's knowledge about the value of property p at each state. (b) A DFA with knowledge uncertainty. The edge labels represent properties' valuations that lead to a transition where  $\tau$  denotes any valuation (tautology). Node  $q_3$  is the accepting state.

**Example 7.2** (Stochastic transition versus knowledge uncertainty). Consider the simple MDP and DFA in Figure 7.5. The DFA accepts any path on the MDP whose induced word is in the form of  $(p\tau)^*\neg p\tau$ . For instance, if the true labels are  $s_0 \models p$ ,  $s_1 \models p$ , and  $s_2 \models \neg p$ , then the path

$$s_0 \xrightarrow{a} s_2 \xrightarrow{c} s_2 \xrightarrow{d} s_0$$

on the MDP generates the run

$$q_0 \xrightarrow{p} q_1 \xrightarrow{\neg p} q_0 \xrightarrow{\neg p} q_2 \xrightarrow{p} q_3$$

on the DFA which is accepting and satisfies the task.

A key observation is that the nature of the probabilities of the MDP transitions are distinct from that of the agent's belief. A stochastic transition means that if the agent takes the same action at the same state multiple times, every time the next state is determined by the given probabilities. Therefore, a path like

$$s_0 \xrightarrow{a} s_1 \xrightarrow{d} s_0 \xrightarrow{a} s_2$$

is possible on the MDP. On the other hand, the true labels of the states are fixed and the distribution of the agent's belief does not translate

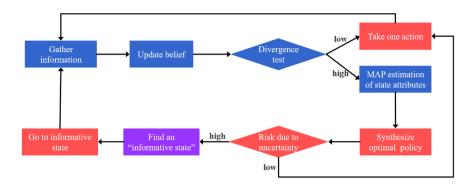


Figure 7.6: The schematic of the perception-planning loop. The blue blocks refer to pure perception modules and the red blocks refer to pure planning modules, and the purple block indicates a combined perception and planning module of the algorithm.

into similar behavior. For instance, the path

$$s_0 \xrightarrow{b} s_0 \xrightarrow{b} s_0 \xrightarrow{b} s_0$$

on the MDP cannot generate the run

$$q_0 \stackrel{p}{\longrightarrow} q_1 \stackrel{p \vee \neg p}{\longrightarrow} q_0 \stackrel{\neg p}{\longrightarrow} q_2 \stackrel{p \vee \neg p}{\longrightarrow} q_3$$

on the DFA. Since in this run, the truth value assignment of property p at state  $s_0$  is inconsistent.

The outline of the proposed algorithm is illustrated in Figure 7.6. At each state, the agent gathers some perception outputs, e.g., sensing measurements, and uses them to update its belief about the environment. If the updated belief, called the posterior belief, is significantly different from the previous one, called the prior belief, the agent must replan. Otherwise, it will continue with its previous policy to take a step. To check the significance of the added information from the new perception data, we compute the difference between the prior and posterior beliefs using statistical distance measured by a divergence metric. If the threshold (a hyperparameter given to the algorithm) on this difference is exceeded, the agent first estimates the most probable configuration of the environment. Then, the agent applies a synthesis algorithm with the estimated environment model. This synthesis algorithm outputs a

strategy that maximizes the probability of satisfying the given temporal logic task. The given strategy induces a Markov chain that is used to calculate the risk due to perception uncertainty. If the risk is lower than a threshold (another hyperparameter given to the algorithm), the agent uses the computed policy to take a step. Otherwise, it will find an active perception strategy to reduce its perception uncertainty. We now proceed to explain different stages of the proposed algorithm.

#### **Information Processing**

Consider the agent's state to be  $s_t$  at time t. The agent will receive new perception outputs according to the observation model  $\mathcal{O}(s_t,.,.,.)$  for all states and atomic propositions.

The agent employs the observations to update its learned model of the environment in a Bayesian approach. For ease of notation, let  $\mathcal{L}_t \equiv \mathcal{L}_t(s,p)$  and  $\mathcal{O}(b) \equiv \mathcal{O}(s_t,s,p,b)$ . Given the prior belief of the agent  $\mathcal{L}_{t-1}$  and the received observations  $\mathcal{Z}$ , the posterior (updated) belief follows

$$\Pr(s \models p | \mathcal{Z}(s_t, s, p) = \mathit{True}) = \frac{\mathcal{L}_{t-1} \mathcal{O}(\mathit{True})}{\mathcal{L}_{t-1} \mathcal{O}(\mathit{True}) + (1 - \mathcal{L}_{t-1}) \mathcal{O}(\mathit{False})},$$

$$\Pr(s \models p | \mathcal{Z}(s_t, s, p) = False) = \frac{\mathcal{L}_{t-1}(1 - \mathcal{O}(True))}{\mathcal{L}_{t-1}(1 - \mathcal{O}(True)) + (1 - \mathcal{L}_{t-1})(1 - \mathcal{O}(False))},$$

for all  $s \in S$  and  $p \in \mathcal{AP}$ . Depending on the truth value observed for p,  $\mathcal{L}_t(s,p)$  will be updated according to one of the above expressions. Besides, for any  $P \subseteq \mathcal{AP}$ ,

$$\mathcal{L}_t(s, P) = \prod_{p \in \mathcal{AP}} \mathcal{L}_t(s, p).$$

#### Divergence Test on the Belief

If the agent's knowledge about the environment configuration has not significantly changed from its knowledge in the previous state, then the agent will continue with its previous policy. Nevertheless, if its knowledge has significantly changed, the agent will synthesize a new policy. We use the Jensen-Shannon divergence to quantify the change in the belief distribution between two consecutive time steps. The cumulative Jensen-Shannon divergence over the states and the propositions can be expressed as

$$D_{\mathcal{J}\mathcal{S}\mathcal{D}}(\mathcal{L}_{t-1}||\mathcal{L}_t) = \frac{1}{2}D_{\mathcal{K}\mathcal{L}}(\mathcal{L}_{t-1}||\mathcal{L}_m^t) + \frac{1}{2}D_{\mathcal{K}\mathcal{L}}(\mathcal{L}_t||\mathcal{L}_m^t),$$

where  $D_{\mathcal{KL}}(.\|.)$  is the Kullback–Leibler divergence between two probability distributions and  $\mathcal{L}_m^t = 1/2 \left(\mathcal{L}_{t-1} + \mathcal{L}_t\right)$  is the average probability distribution. One of the input parameters to the algorithm is a threshold  $\gamma_d$  on the above divergence. If  $\gamma_d$  is not exceeded, the agent uses its previous policy  $\pi_{t-1}$  to pick an action and transitions according to its outcome. Otherwise, it has to synthesize a new policy.

For the policy synthesis step, the agent first estimates the most probable environment configuration from the distribution dictated by its updated belief. Let  $\hat{\mathcal{L}}_t: S \to 2^{\mathcal{AP}}$  indicate the agent's inference of the environment configuration at time t. The maximum a posteriori estimation is fairly simple as it decomposes into finding the mode of the posterior distribution for each property at each state. For binary-valued atomic propositions that follow Bernoulli distributions, the inference turns into picking the more probable outcome for each property at each state, i.e.,

$$\hat{\mathcal{L}}_t(s) = \{ p \in \mathcal{AP} \mid \mathcal{L}_t(s, p) \ge 0.5 \}.$$

### **Policy Synthesis**

Finding an optimal policy, i.e., a policy that maximizes the probability of realizing a temporal logic task, translates into a reachability criterion on the product MDP. Let  $\mathcal{F}_{\mathcal{D}}^{\mathcal{M}} = S \times \mathcal{F}$  denote the equivalent accepting states on the product MDP. The agent must find a policy that with high probability reaches to  $\mathcal{F}_{\mathcal{D}}^{\mathcal{M}}$ .

Given that there exists an optimal memoryless deterministic policy on the product MDP (Baier and Katoen, 2008), one can restrict the search space to that of memoryless deterministic policies and formulate

$$\pi_{t} = \operatorname*{arg\,max}_{\pi \in \Pi_{nm,d}} \Pr(\mathcal{M}_{\mathcal{D}}^{\pi} \models \Diamond \mathcal{F}_{\mathcal{D}}^{\mathcal{M}} \mid \hat{\mathcal{L}}_{t}), \tag{7.8}$$

where  $\Pi_{nm,d}$  is the set of memoryless and deterministic policies. To find  $\pi_t$  in (7.8), we use a linear programming (LP) approach (Baier and Katoen, 2008). The optimal value of the linear program is the maximum probability of reaching the set of accepting states. In order to find the corresponding optimal policy, it suffices to find the actions for which the corresponding LP constraints are active. If there are more than one action with active constraint for a state, any of those actions can be chosen arbitrarily.

#### Risk Assessment of Imperfect Perception

To assess the risk of the computed policy due to perception uncertainties, we now factor in the probabilistic belief of the agent over the environment properties. In particular, we first generate the induced Markov chain  $\mathcal{M}_{\mathcal{D}}^{\pi_t}$  by applying the policy  $\pi_t$  over the product MDP  $\mathcal{M}_{\mathcal{D}}$ . Next, we verify the induced Markov chain with the uncertain labels against the task specification. We develop an algorithm via a computation graph that yields the exact probability of the task realization. However, due to the complexities explained in Example 7.2, such quantitative analysis has exponential complexity, as formalized in the next theorem.

**Theorem 7.1.** Let  $\mathcal{M}_{\mathcal{D}}^{\pi_t}$  to be a Markov chain with n states and  $\mathcal{L}_t$  to be a fully probabilistic labeling function (i.e., all labels are uncertain) over m atomic propositions. Quantitative verification of  $\mathcal{M}_{\mathcal{D}}^{\pi_t}$  against a reachability specification has a complexity of  $O(n2^{nm})$ .

Since the complexity of an exact quantitative analysis is prohibitive, we instead propose a statistical verification. More specifically, we approximate the expected value of the probability of the task realization over all possible instances of the environment

$$\mathbb{E}_{\mathcal{L} \sim Dist(\mathcal{L})} \left[ \Pr \left( \mathcal{M}_{\mathcal{D}}^{\pi_t} \models \varphi \right) \right]$$

by an empirical expectation with N samples

$$\hat{\mathbb{E}}_{\mathcal{L} \sim Dist(\mathcal{L})} \left[ \Pr \left( \mathcal{M}_{\mathcal{D}}^{\pi_t} \models \varphi \right) \right] = \frac{1}{N} \sum_{i=1}^{N} \Pr \left( \mathcal{M}_{\mathcal{D}}^{\pi_t} \models \varphi | \mathcal{L}_i \right),$$

where  $\mathcal{L}_i$  are samples drawn from  $Dist(\mathcal{L}) = \mathcal{L}_t$ . By the application of Hoeffding's inequality (Hoeffding, 1994), we establish the following concentration result for this approximation.

**Theorem 7.2.** Let  $f(\mathcal{L}_i) = \Pr\left(\mathcal{M}_{\mathcal{D}}^{\pi_t} \models \varphi | \mathcal{L}_i\right)$  denote the output of verification for an environment modeled by  $\mathcal{L}_i$ . The empirical expectation of  $\mathbb{E}_{\mathcal{L} \sim Dist(\mathcal{L})}[f(\mathcal{L})]$  with N samples has the following concentration bound

$$\Pr\left(\left|\mathbb{E}_{\mathcal{L}\sim Dist(\mathcal{L})}\left[f(\mathcal{L})\right] - \frac{1}{N}\sum_{i=1}^{N}f(\mathcal{L}_i)\right| \ge \epsilon\right) \le 2\exp(-2N\epsilon^2).$$

It is worth noting that  $\Pr\left(\mathcal{M}_{\mathcal{D}}^{\pi_t} \models \varphi | \mathcal{L}_i\right)$  itself is the output of a system of linear equations (Baier and Katoen, 2008) that depend on the sampled labeling function  $\mathcal{L}_i$ . Therefore, further characterization of  $f(\mathcal{L}_i)$  such as bounding its higher moments is very difficult.  $\Pr\left(\mathcal{M}_{\mathcal{D}}^{\pi_t} \models \varphi | \mathcal{L}_i\right)$  can also be computed via statistical verification techniques by sampling paths over the Markov chain (Agha and Palmskog, 2018).

For a policy  $\pi_t$ , we define a risk parameter

$$\mathcal{R}(\mathcal{M}_{\mathcal{D}}, \pi_{t}, \mathcal{L}_{t}, \varphi) = \left| \Pr(\mathcal{M}_{\mathcal{D}}^{\pi_{t}} \models \varphi \,|\, \hat{\mathcal{L}}_{t}) - \mathbb{E}_{\mathcal{L} \sim Dist(\mathcal{L})} \left[ \Pr\left(\mathcal{M}_{\mathcal{D}}^{\pi_{t}} \models \varphi\right) \right] \right|,$$
(7.9)

which accounts for the variation in the task realization guarantee of the policy with respect to the perception uncertainty. Another input parameter to the proposed perception and planning algorithm is a threshold  $\gamma_r$  on the risk due to perception uncertainty. If  $\gamma_r$  is not exceeded, the agent acts according to the computed policy  $\pi_t$ . Otherwise, it takes an active perception strategy as explained next.

# **Active Perception Strategy**

We develop an algorithm to compute an active perception strategy in the form of a sequence of actions that the agent follows to reduce its

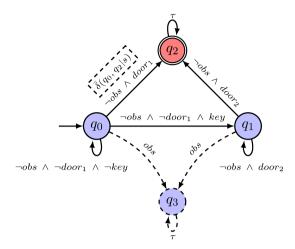


Figure 7.7: Total probabilistic finite automaton for the formula  $\varphi = (\neg obs \ U \ door_1) \lor ((\neg door_2 \ U \ key) \land (\neg obs \ U \ door_2))$ . State  $q_3$ , a sink state, and the transitions to it have been added to make the automaton total. The transitions between the states of the automaton are probabilistic where the probabilities depend on the belief over the states' propositions. For example,  $\bar{\delta}(q_0, q_2|s)$  indicates how  $\mathcal{L}_t(s, .)$  determines the probability of the transition between  $q_0$  and  $q_2$ .

perception uncertainty. We consider three criteria for computing an active perception strategy. First, such a strategy should enable the agent to reduce its uncertainty about the value of the propositions that affect the task progress. These propositions are the ones that enable the transitions from the current stage of the task, i.e., state of the automaton, to the next ones. For example in Figure 7.7, if the agent is at state  $q_1$ , the propositions that matter are obs and  $door_2$ . To measure the uncertainty reduction, we use expected entropy of the said propositions over the whole state space. Second, an active perception strategy must not affect the stage of the task and so, the agent has to remain in the same state of the automaton. Third, after the agent completes the sequence of actions, it should be able to return to the point from which it started the active perception strategy.

Based on these criteria, we propose an algorithm to construct an active perception strategy. The algorithm takes a bound  $C_{Act}$  on the number of actions and uses that to construct a tree of depth  $C_{Act}$ . Each node of

the tree has four parameters: distribution over the states of the MDP, distribution over the states of the DFA, expected entropy reduction, and reachability probability back to the root node. The distribution over the states of the MDP depends on the MDP's transitions while the distribution over the states of the DFA depend on the agent's belief over the propositions, as shown in Figure 7.7. Once the tree is generated, the algorithm picks the best node using a hyperparameter  $\beta$  that weighs safety versus information quality. Safety refers to the ability of the agent to remain in the same state of the automaton as well as its ability to return to the root node while information quality refers to the amount of entropy reduction. The sequence of actions leading to the optimum node with respect to a combination of safety and information quality results in the active perception strategy. After following this sequence of actions, the perception-planning loop starts over.

# 7.3.3 Task-Oriented Active Perception and Planning Case Studies Simulating Planar Navigation with Finite-Horizon Tasks

We consider an agent that navigates in a discretized 2D environment and has a finite-horizon task. For instance, the task encoded as a DFA in Figure 7.7 asks the agent to either go to the state where  $door_1$  is located or find a key and go to the state where  $door_2$  is located, while avoiding the obstacles. We implemented different versions of the task-oriented perception and planning algorithm to evaluate the effect of each module on the performance.

Table 7.1 reports the results for a reach-avoid task in an environment with 64 states and with randomly generated obstacles and target. In the table, No perc. refers to a baseline scenario where the agent estimates the environment configuration with its prior knowledge and plans according to that. Perc. w/ no update is a perception strategy that incorporates only the most recent perception output. Perc. w/ update is perception with a Bayesian update, as described in Section 7.3.2. With the exception of the first algorithm, all algorithms have a replanning module, however, the ones with div. replan only if the divergence threshold over the change in the belief is exceeded. info. means that active perception is

**Table 7.1:** Results of planar navigation under a reach-avoid task using different versions of the proposed algorithm. Success is the percentage of runs that complete the task. #Step is the average length of the runs and #Plan is the number of times that the agent synthesizes a new policy.

Algorithm	Success	#Step	#Plan
No perc.	0%	50	1
Perc. $w/$ no update $+$ replan	0%	38.4	38.4
Perc. $w/update + replan$	84%	21.8	21.8
Perc. $w/update + div.$	80%	22.8	14.8
Perc. $w/update + replan + info.$	92%	19.4	19.4
Perc. w/ update + div. + info.	86%	22.6	14.6

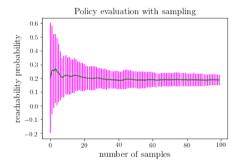


Figure 7.8: Statistical verification of an induced Markov chain with uncertain atomic propositions.

enabled and hence the agent will perform active perception strategies when the risk due to perception uncertainty is high. The results show that adding the divergence test reduces the number of policy synthesis steps. Furthermore, the divergence test reduces the success rate. On the other hand, adding the active perception module increases the success rate. Figure 7.8 depicts the results from a risk assessment step for the MDP with 64 states and 2 atomic propositions. Even though the size of the sampling space is large  $(2^{64})$ , it can be seen that the empirical expectation of the reachability probability quickly converges with about 20 samples.



Figure 7.9: A scene from the created urban environment.

#### **Drone Navigation in Simulated Urban Environment**

In the AirSim Shah et al., 2017 simulator, we designed an urban environment as shown in Figure 7.9 and tasked a drone to fly from an initial state to a specific flagged building while avoiding collision with other entities of the environment. The drone is equipped with 4 cameras and 4 depth sensors. The perception module processes the cameras' readings as well as the depth measurements to map the semantic labels to a discretized model of the environment. We applied the proposed perception and planning scheme. However, in contrast to the previous simulation scenario, an observation model does not exist here. Therefore, we used a frequentist update rule for the agent's belief. Details of the simulation setting as well as recordings of the resulting behavior of the drone are available in Ghasemi et al., 2020.

# **Runtime Assurance via Shielding**

Runtime verification detects violations of a set of specified properties while a system is executing (Leucker and Schallhart, 2009). An extension of this idea is to perform runtime enforcement of specified properties, in which violations are not only detected but also overwritten in a way that specified properties are maintained.

A general approach for runtime enforcement of specified properties is *shield synthesis*, in which a shield monitors the system and instantaneously overwrites incorrect outputs. A shield must ensure both *correctness*, i.e., it corrects system outputs such that all properties are always satisfied, as well as *minimum deviation*, i.e., it deviates from system outputs only if necessary and as rarely as possible. The latter requirement is important because the system may satisfy additional noncritical properties that are not considered by the shield but should be retained as much as possible.

Merging formal methods with shielding we require the notion of k-

This section follows an adapted exposition from Humphrey, Könighofer, Könighofer, and Topcu (2016) and Könighofer, Alshiekh, Bloem, Humphrey, Könighofer, Topcu, and Wang (2017).

stabilizing shields (Bloem et al., 2015). Given a safety specification, we can identify wrong outputs, i.e., outputs after which the specification is violated or, more precisely, after which the environment can force the specification to be violated. A "wrong" trace is then a trace that ends in a wrong output. The idea of shields is that they may modify the outputs so that the specification always holds, but that such deviations last for at most k consecutive steps after a wrong output. If a second violation happens during the k-step recovery phase, the shield enters a mode where it only enforces correctness, but no longer minimizes the deviation.

#### 8.0.1 Overview of Shielding

We apply shielding to reactive systems, which model well with Mealy machines—finite-state machines that depend both on state and current inputs. We can view a Mealy machine formally as the tuple  $\mathcal{M} = (S, q_0, \Sigma_I, \Sigma_O, \delta, \lambda)$  composed of the following data.

- A finite set of states, S.
- A starting state,  $q_0$ .
- An input alphabet,  $\Sigma_I$ .
- An output alphabet,  $\Sigma_O$ .
- A transition function,  $\delta \colon S \times \Sigma_I \to S$ .
- An output function,  $\lambda \colon S \times \Sigma_I \to \Sigma_O$ .

This tuple is the fundamental unit for the formalism of a system model, termed the *design*, and its *shield*.

**Definition 8.1** (Shield, Bloem *et al.* (2015)). Let  $\mathcal{D} = (S, q_0, \Sigma_I, \Sigma_O, \delta, \lambda)$  be a design,  $\varphi$  be a set of properties, and  $\varphi^v \subseteq \varphi$  be a valid subset such that  $\mathcal{D} \models \varphi^v$ . A reactive system  $\mathcal{S} = (S', q'_0, \Sigma, \Sigma_O, \delta', \lambda')$  is a *shield* of  $\mathcal{D}$  with respect to  $(\varphi \setminus \varphi^v)$  if and only if  $(\mathcal{D} \circ \mathcal{S}) \models \varphi$ . We also require that for  $\mathcal{S}$  to be a it must be shield of *any* design  $\mathcal{D}$  such that  $\mathcal{D} \models \varphi^v$ .

#### Correctness property

With correctness we refer to the property that the shield corrects any design's output such that a given safety specification is satisfied.

Since a shield must work for any design, the synthesis procedure does not need to consider the design's implementation. This property is crucial because the design may be unknown or too complex to analyze. On the other hand, the design may satisfy additional (noncritical) specifications that are not specified in  $\varphi$  but should be retained as much as possible.

#### Minimum deviation property

Minimum deviation requires a shield to deviate only if necessary, and as infrequently as possible. To ensure minimum deviation, a shield can only deviate from the design if a property violation becomes unavoidable. Given a safety specification  $\varphi$ , a shield S does not deviate unnecessarily if for any design  $\mathcal{D}$  and any trace that is not wrong and  $\mathcal{D}$  does not violate  $\varphi$ , S keeps the output of  $\mathcal{D}$  intact.

#### **Admissibility property**

To address shortcoming with k-stabilizing shields (1), we guarantee the following: (a) Admissible shields are subgame optimal. That is, for any wrong trace, if there is a finite number k of steps within which the recovery phase can be guaranteed to end, the shield will always achieve this. (b) The shield is admissible, that is, if there is no such number k, it always picks a deviation that is optimal in that it ends the recovery phase as soon as possible for some possible future inputs. As a result, admissible shields work well in settings in which finite recovery can not be guaranteed, because they guarantee correctness and may well end the recovery period if the system does not pick adversarial outputs. To address shortcoming (2), admissible shields allow arbitrary failure frequencies and in particular failures that arrive during recovery, without losing the ability to recover.

To create an admissible shield we must first define what k-stabilizing a trace means. Given an unavoidable violation occurring in design  $\mathcal{D}$ ,

the shield can deviate from the expected outputs of the design for at most k consecutive time steps. This means that violations cannot be consecutive and can only be tolerated after k steps. In the event that the design violates the specification within k steps, then the shield transitions the system to a fail-safe mode

**Definition 8.2** (Admissible Shield). A shield S is admissible if for any trace, whenever there exists a k and a shield S' such that S' adversely k-stabilizes the trace, then S adversely k-stabilizes the trace. If such a k does not exist for trace, then S collaboratively k-stabilizes the trace for a minimal k.

#### 8.1 Shielding a UAV mission

In this section, we apply shields on a scenario in which a UAV must maintain certain properties while performing a surveillance mission in a dynamic environment. We show how a shield can be used to enforce the desired properties, where a human operator in conjunction with a lower-level autonomous planner is considered as the reactive system that sends commands to the UAV's autopilot. We discuss how we would intuitively want a shield to behave in such a situation.

A common UAV control architecture consists of a ground control station that communicates with an autopilot onboard the UAV (Chao et al., 2010). The ground control station receives and displays updates from the autopilot on the UAV's state, including position, heading, airspeed, battery level, and sensor imagery. It can also send commands to the UAV's autopilot, such as waypoints to fly to. A human operator can then use the ground control station to plan waypoint-based routes for the UAV, possibly making modifications during mission execution to respond to events observed through the UAV's sensors. However, mission planning and execution can be very workload intensive, especially when operators are expected to control multiple UAVs simultaneously (Donmez et al., 2010). To address this issue, methods for UAV command and control have been explored in which operators issue high-level commands, and automation carries out low-level execution details.

Several errors can occur in this type of human-automation paradigm (Chen and Barnes, 2012). For instance, in issuing high-level commands to the low-level planner, a human operator might neglect required safety properties due to high workload, fatigue, or an incomplete understanding of exactly how the autonomous planner might execute the command. The planner might also neglect these safety properties either because of software errors or by design. Waypoint commands issued by the operator or planner could also be corrupted by software that translates waypoint messages between ground station and autopilot specific formats or during transmission over the communication link.

As the mission unfolds, waypoint commands will be sent periodically to the autopilot. If a waypoint that violates the properties is received, a shield that monitors the system inputs and can overwrite the waypoint outputs to the autopilot would be able to make corrections to ensure the satisfaction of the desired properties.

Consider a mission map (Figure 8.1) (Feng et al., 2016b), which contains three tall buildings (illustrated as blue blocks), over which a UAV should not attempt to fly. It also includes two unattended ground sensors (UGS) that provide data on possible nearby targets, one at location  $loc_1$  and one at  $loc_x$ , as well as two locations of interest,  $loc_y$  and  $loc_z$ . The UAV can monitor  $loc_x$ ,  $loc_y$ , and  $loc_z$  from several nearby vantage points. The map also contains a restricted operating zone (ROZ), illustrated with a red box, in which flight might be dangerous, and the path of a possible adversary that should be avoided (the pink dashed line). Inside the communication relay region (large green area), communication links are highly reliable. Outside this region, communication relies on relay points with lower reliability. Given this scenario, properties of interest include:

- Connected waypoints. The UAV is only allowed to fly to directly connected waypoints.
- No communication. The UAV is not allowed to stay in a location with reduced communication reliability.
- Restricted operating zones. The UAV has to leave a ROZ

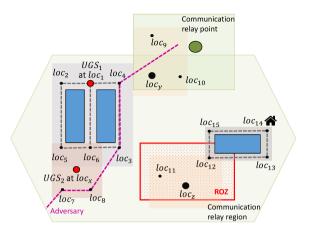


Figure 8.1: A map for UAV mission planning.

within 2 time steps.

- **Detected by an adversary.** Locations on the adversary's path cannot be visited more than once over any window of 3 time steps.
- UGS. If a UGS reports a possible nearby target, the UAV should visit a respective waypoint within 7 steps (for  $UGS_1$  visit  $loc_1$ , for  $UGS_2$  visit  $loc_5$ ,  $loc_6$ ,  $loc_7$ , or  $loc_8$ ).
- Go home. Once the UAV's battery is low, it should return to a designated landing site at  $loc_{14}$  within 10 time steps.

The task of the shield is to ensure these properties during operation. In this setting, the operator in conjunction with a lower-level planner acts as a reactive system that responds to mission-relevant inputs; in this case data from the UGSs and a signal indicating whether the battery is low. In each step, the next waypoint is sent to the autopilot, which is encoded in a bit representation via outputs  $l_4$ ,  $l_3$ ,  $l_2$ , and  $l_1$ . The shield monitors (Figure 8.2) mission inputs and waypoint outputs, correcting outputs immediately if a violation of the safety properties becomes unavoidable.

We represent each of the properties by a safety automaton, the product of which serves as the shield specification. Figure 8.3 models the "connected

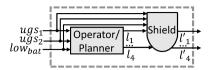


Figure 8.2: The interaction between the operator/planner (acting as a reactive system) and the shield.

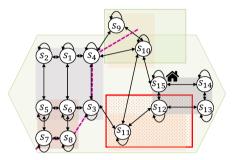


Figure 8.3: Safety automaton of Property 8.1 over the map in Figure 8.1.

waypoints" property, where each state represents a waypoint with the same number. Edges are labeled by the values of the variables  $l_4 \dots l_1$ . For example, the edge leading from state  $s_5$  to state  $s_6$  is labeled by  $\neg l_4 l_3 l_2 \neg l_1$ . For clarity, we drop the labels of edges in Figure 8.3. The automaton also includes an error state, which is not shown. Missing edges lead to this error state, denoting forbidden situations.

How should a shield behave in this scenario? If the human operator wants to monitor a location in a ROZ, he or she would like to simply command the UAV to "monitor the location in the ROZ and stay there", with the planner handling the execution details. If the planner cannot do this while meeting all the safety properties, it is appropriate for the shield to revise its outputs. Yet, the operator would still expect his or her commands to be followed to the maximum extent possible, leaving the ROZ when necessary and returning whenever possible. Thus, the shield should minimize deviations from the operator's directives as executed by the planner.

#### Using a k-stabilizing shield

As a concrete example, assume the UAV is currently at  $loc_3$ , and the operator commands it to monitor  $loc_{12}$ . The planner then sends commands to fly to  $loc_{11}$  then  $loc_{12}$ , which are accepted by the shield. The planner then sends a command to loiter at  $loc_{12}$ , but the shield must overwrite it to maintain Property 8.1, which requires the UAV to leave the ROZ within two time steps. The shield instead commands the UAV to go to  $loc_{15}$ . Suppose the operator then commands the UAV to fly to  $loc_{13}$ , while the planner is still issuing commands as if the UAV is at  $loc_{12}$ . The planner then commands the UAV to fly to  $loc_{13}$ , but since the actual UAV cannot fly from  $loc_{15}$  to  $loc_{13}$  directly, the shield directs the UAV to  $loc_{14}$  on its way to  $loc_{15}$ . The operator might then respond to a change in the mission and command the UAV fly back to  $loc_{12}$ , and the shield again deviates from the route assumed by the planner, and directs the UAV back to  $loc_{15}$ , and so on. Therefore, a single specification violation can lead to an infinitely long deviation between the UAV's actual position and the UAV's assumed position. A k-stabilizing shield is allowed to deviate from the planner's commands for at most k consecutive time steps. Hence, no k-stabilizing shield exists.

#### Using an admissible shield

Recall the situation in which the shield caused the actual position of the UAV to "fall behind" the position assumed by the planner, so that the next waypoint the planner issues is two or more steps away from the UAV's current waypoint position. The shield should then implement a best-effort strategy to "synchronize" the UAV's actual position with that assumed by the planner. Though this cannot be guaranteed, the operator and planner are not adversarial towards the shield, so it will likely be possible to achieve this re-synchronization, for instance when the UAV goes back to a previous waypoint or remains at the current waypoint for several steps. This possibility motivates the concept of an admissible shield. Assume that the actual position of the UAV is  $loc_{14}$  and the its assumed position is  $loc_{13}$ . If the operator commands the UAV to loiter at  $loc_{13}$ , the shield will be able to catch up with the

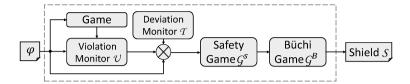


Figure 8.4: Synthesizing admissible shields.

state assumed by the planner and to end the deviation by the next specification violation.

#### 8.2 Synthesizing Admissible Shields

We will illustrate shield synthesis using the admissible shielding definition, but the process is similar for other types of shields. Starting from a safety specification  $\varphi = (Q, q_0, \Sigma, \delta, F)$  with  $\Sigma = \Sigma_I \times \Sigma_O$ , the admissible shield synthesis procedure consists of five steps (Figure 8.4).

#### Step 1. Constructing the Violation Monitor ${\mathfrak U}$

From  $\varphi$  we build the automaton  $\mathcal{U} = (U, u_0, \Sigma, \delta^u)$  to monitor property violations by the design. The goal is to identify the latest point in time from which a specification violation can still be corrected with a deviation by the shield. This constitutes the start of the recovery period, in which the shield is allowed to deviate from the design. In this phase the shield monitors the design from all states that the design could reach under the current input and a correct output. A second violation occurs only if the next design's output is inconsistent with all states that are currently monitored. In case of a second violation, the shield monitors the set of all input-enabled states that are reachable from the current set of monitored states.

The first phase of the construction of the violation monitor  $\mathcal{U}$  considers  $\varphi = (Q, q_0, \Sigma, \delta, F)$  as a *safety game* and computes its winning region  $W \subseteq F$  so that every reactive system  $\mathcal{D} \models \varphi$  must produce outputs such that the next state of  $\varphi$  stays in W. Only in cases in which the next state of  $\varphi$  is outside of W the shield is allowed to interfere.

The second phase expands the state space Q to  $2^Q$  via a subset construction, with the following rationale. If the design makes a mistake (i.e., picks outputs such that  $\varphi$  enters a state  $q \notin W$ ), we have to "guess" what the design actually meant to do and we consider all output letters that would have avoided leaving W and continue monitoring the design from all the corresponding successor states in parallel. Thus,  $\mathcal U$  is essentially a subset construction of  $\varphi$ , where a state  $u \in U$  of  $\mathcal U$  represents a set of states in  $\varphi$ .

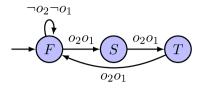
The third phase expands the state space of  $\mathcal{U}$  by adding a counter  $d \in \{0,1,2\}$  and a output variable z. Initially d is 0. Whenever a property is violated d is set to 2. If d > 0, the shield is in the recovery phase and can deviate. If d = 1 and there is no other violation, d is decremented to 0. In order to decide when to decrement d from 2 to 1, we add an output z to the shield. If this output is set to True and d = 2, then d is set to 1.

The final violation monitor is  $\mathcal{U} = (U, u_0, \Sigma^u, \delta^u)$ , with the set of states  $U = (2^Q \times \{0, 1, 2\})$ , the initial state  $u_0 = (\{q_0\}, 0)$ , the input/output alphabet  $\Sigma^u = \Sigma_I \times \Sigma_O^u$  with  $\Sigma_O^u = \Sigma_O \cup z$ , and the next-state function  $\delta^u$ , which obeys the following rules:

- 1.  $\delta^u((u,d),(\sigma_I,\sigma_O)) = (\{q' \in W \mid \exists q \in u, \sigma_O' \in \Sigma_O^u . \delta(q,(\sigma_I,\sigma_O')) = q'\}, 2)$  if  $\forall q \in u . \delta(q,(\sigma_I,\sigma_O)) \notin W$ , and
- 2.  $\delta^u((u,d),\sigma) = (\{q' \in W \mid \exists q \in u \, . \, \delta(q,\sigma) = q'\}, \operatorname{dec}(d))$  if  $\exists q \in u \, . \, \delta(q,\sigma) \in W$ , and  $\operatorname{dec}(0) = \operatorname{dec}(1) = 0$ , and if z is  $\mathit{True}$  then  $\operatorname{dec}(2) = 1$ , else  $\operatorname{dec}(2) = 2$ .

Our construction sets d=2 whenever the design leaves the winning region, and not when it enters an unsafe state. Hence, the shield  $\mathcal{S}$  can take a remedial action as soon as "the crime is committed", before the damage is detected, which would have been too late to correct the erroneous outputs of the design.

**Example 8.1.** We illustrate the construction of  $\mathcal{U}$  using the specification  $\varphi$  from Figure 8.5 over the outputs  $o_1$  and  $o_2$ . (Figure 8.5 represents a safety automaton if we make all missing edges point



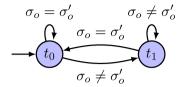


Figure 8.5: Safety automaton of Example 8.1.

Figure 8.6: The deviation monitor  $\mathfrak{I}$ .

to an (additional) unsafe state.) The winning region consists of all safe states, i.e.,  $W = \{F, S, T\}$ . The resulting violation monitor is  $\mathcal{U} = (\{F, S, T, FS, ST, FT, FST\} \times \{0, 1, 2\}, (F, 0), \Sigma^u, \delta^u)$ . The transition relation  $\delta^u$  is illustrated in Table 8.1 and lists the next states for all possible present states and outputs. Lightning bolts denote specification violations. The update of counter d, which is not included in Table 8.1, is as follows: Whenever the design commits a violation d is set to 2. If no violation exists, d is decremented in the following way: if d = 1 or d = 0, d is set to 0. If d = 2 and d is set to 1, else d remains 2. In this example, d is set to d is set d is set to d is set d is d is d is d is d is set d is d in

Let us take a closer look at some entries of Table 8.1. If the current state is  $(\{F\},0)$  and we observe the output  $\neg o_2o_1$ , a specification violation occurs. We assume that  $\mathcal{D}$  meant to give an allowed output, either  $o_2o_1$  or  $\neg o_2\neg o_1$ . The shield continues to monitor both F and S; thus,  $\mathcal{U}$  enters the state  $(\{F,S\},2)$ . If the next observation is  $o_2o_1$ , which is allowed from both possible current states, the possible next states are S and T, therefore  $\mathcal{U}$  traverses to state  $(\{S,T\},2)$ . However, if the next observation is again  $\neg o_2o_1$ , which is neither allowed in F nor in S, we know that a second violation occurs. Therefore, the shield monitors the design from all three states and  $\mathcal{U}$  enters the state  $(\{F,S,T\},2)$ .

#### Step 2. Constructing the Deviation Monitor ${\mathfrak T}$

We build  $\mathfrak{T}=(T,t_0,\Sigma_O\times\Sigma_O,\delta^t)$  to monitor deviations between the shield and design outputs. Here,  $T=\{t_0,t_1\}$  and  $\delta^t(t,(\sigma_O,\sigma_O'))=t_0$  if and only if  $\sigma_O=\sigma_O'$ . That is, if there is a deviation in the current time step, then  $\mathfrak{T}$  will be in  $t_1$  in the next time step. Otherwise, it will

	$\neg o_1 \neg o_2$	$\neg o_1 o_2$ or $o_1 \neg o_2$	$o_1 o_2$
{F}	{F}	$\{F,S\}$	{S}
${S}$	$\{T\}^{f}$	$\{\mathrm{T}\}^{{\color{red} t}}$	$\{T\}$
$\{T\}$	{F}⁴	$\{\mathrm{F}\}^{{\color{red} t}}$	$\{F\}$
$\{F,S\}$	{F}	$\{\mathrm{F,S,T}\}^{rac{f}{2}}$	$\{S,T\}$
$\{S,T\}$	{F,T}⁴	$\{\mathrm{F},\mathrm{T}\}^{m{\ell}}$	$\{F,T\}$
$\{F,T\}$	$\{F\}$	$\{\mathrm{F,S,T}\}^{rac{m{\ell}}{2}}$	$\{F,S\}$
$\{F,S,T\}$	$\{F\}$	$\{\mathrm{F,S,T}\}^{rac{m{\ell}}{2}}$	$\{F,S,T\}$

**Table 8.1:** Transition relation  $\delta^u$  of monitor  $\mathcal{U}$  for example 8.1.

be in  $t_0$ . This deviation monitor is shown in Figure 8.6.

#### Step 3. Constructing and Solving the Safety Game $S^s$

Given the automata  $\mathcal{U}$  and  $\mathcal{T}$  and the safety automaton  $\varphi$ , we construct a safety game  $\mathcal{G}^s = (G^s, g_0^s, \Sigma_I^s, \Sigma_O^s, \delta^s, F^s)$ , which is the synchronous product of  $\mathcal{U}$ ,  $\mathcal{T}$ , and  $\varphi$ , such that  $G^s = U \times T \times Q$  is the state space,  $g_0^s = (u_0, t_0, q_0)$  is the initial state,  $\Sigma_I^s = \Sigma_I \times \Sigma_O$  is the input of the shield,  $\Sigma_O^s = \Sigma_O \cup \{z\}$  is the output of the shield,  $\delta^s$  is the next-state function, and  $F^s$  is the set of safe states such that  $\delta^s((u, t, q), (\sigma_I, \sigma_O), (\sigma_O', z)) =$ 

$$(\delta^u[u,(\sigma_I,(\sigma_O,z))],\delta^t[t,(\sigma_O,\sigma_O')],\delta[q,(\sigma_I,\sigma_O')]),$$

and 
$$F^s = \{(u, t, q) \in G^s \mid q \in F \land u = (w, 0) \to t = t_0\}.$$

We require  $q \in F$ , which ensures that the output of the shield satisfies  $\varphi$ , and that the shield can only deviate in the recovery period (i.e., if d=0, no deviation is allowed). We use standard algorithms for safety games (cf. Faella (2009)) to compute the winning region  $W^s$  and the most permissive non-deterministic winning strategy  $\rho_s: G \times \Sigma_I \to 2^{\Sigma_O}$  that is not only winning for the system, but also contains all deterministic winning strategies.

# Step 4. Constructing the Büchi Game $\mathfrak{G}^b$

Implementing the safety game ensures correctness  $(\mathcal{D} \circ \mathcal{S} \models \varphi)$  and that the shield  $\mathcal{S}$  keeps the output of the design  $\mathcal{D}$  intact, if  $\mathcal{D}$  does

not violate  $\varphi$ . The shield still has to keep the number of deviations per violation to a minimum. Therefore, we would like the recovery period to be over infinitely often. This can be formalized as a Büchi winning condition. We construct the Büchi game  $\mathcal{G}^b$  by applying the non-deterministic safety strategy  $\rho^s$  to the game graph  $\mathcal{G}^s$ .

Given the safety game  $\mathcal{G}^s=(G^s,g^s_0,\Sigma^s_I,\Sigma^s_O,\delta^s,F^s)$  with the non-deterministic winning strategy  $\rho^s$  and the winning region  $W^s$ , we construct a Büchi game  $\mathcal{G}^b=(G^b,g^b_0,\Sigma^s_I,\Sigma^s_O,\delta^b,F^b)$  such that  $G^b=W^s$  is the state space, the initial state  $g^b_0=g^s_0$  and the input/output alphabet  $\Sigma^b_I=\Sigma^s_I$  and  $\Sigma^b_O=\Sigma^s_O$  remain unchanged,  $\delta^b=\delta^s\cap\rho^s$  is the transition function, and  $F^b=\{(u,t,q)\in W^s\mid (u=(w,0)\vee u=(w,1))\}$  is the set of accepting states. A play is winning if  $d\leq 1$  infinitely often.

# Step 5. Solving the Büchi Game $\mathfrak{S}^b$

Most likely, the Büchi game  $\mathcal{G}^b$  contains reachable states, for which  $d \leq 1$  cannot be enforced infinitely often. We implement an admissible strategy that enforces to visit  $d \leq 1$  infinitely often whenever possible. This criterion essentially asks for a strategy that is winning with the help of the design.

The admissible strategy  $\rho^b$  for a Büchi game  $\mathcal{G}^b = (G^b, g_0^b, \Sigma_I^b, \Sigma_O^b, \delta^b, F^b)$  can be computed as follows Faella (2009):

- 1. Compute the winning region  $W^b$  and a winning strategy  $\rho_w^b$  for  $\mathcal{G}^b$  (cf. Mazala (2001)).
- 2. Remove all transitions that start in  $W^b$  and do not belong to  $\rho^b_w$  from  $\mathcal{G}^b$ . This results in a new Büchi game  $\mathcal{G}^b_1 = (G^b, g^b_0, \Sigma^b_I, \Sigma^b_O, \delta^b_1, F^b)$  with

$$(g, (\sigma_I, \sigma_O), g') \in \delta_1^b$$
if  $(g, \sigma_I, \sigma_O) \in \rho_w^b$ 
or if  $\forall \sigma_O' \in \Sigma_O^b . (g, \sigma_I, \sigma_O') \notin \rho_w^b \land (g, (\sigma_I, \sigma_O), g') \in \delta^b$ .

3. In the resulting game  $\mathcal{G}_1^b$ , compute a cooperatively winning strategy  $\rho^b$ . In order to compute  $\rho^b$ , one first has to transform all input

variables to output variables. This results in the Büchi game  $\mathcal{G}_2^b = (G^b, g_0^b, \emptyset, \Sigma_I^b \times \Sigma_O^b, \delta_1^b, F^b)$ . Afterwards,  $\rho^b$  can be computed with the standard algorithm for the winning strategy on  $\mathcal{G}_2^b$ .

The strategy  $\rho^b$  is an admissible strategy of the game  $\mathcal{G}^b$ , since it is winning and cooperatively winning (Faella, 2009). Whenever the game  $\mathcal{G}^b$  starts in a state of the winning region  $W^b$ , any play created by  $\rho^b_w$  is winning. Since  $\rho^b$  coincides with  $\rho^b_w$  in all states of the winning region  $W^b$ ,  $\rho^b$  is winning. We know that  $\rho^b$  is cooperatively winning in the game  $\mathcal{G}^b$ . A proof that  $\rho^b$  is also cooperatively winning in the original game  $\mathcal{G}^b$  can be found in Faella (2009).

**Theorem 8.1.** A shield that implements the admissible strategy  $\rho^b$  in the Büchi game  $\mathcal{G}^b = (G^b, g_0^b, \Sigma_I^b, \Sigma_O^b, \delta^b, F^b)$  in a new reactive system  $\mathcal{S} = (G^b, g_0^b, \Sigma_I^b, \Sigma_O^b, \delta', \rho^b)$  with  $\delta'(g, \sigma_I) = \delta^b(g, \sigma_I, \rho^b(g, \sigma_I))$  is an admissible shield.

Proof. First, the admissible strategy  $\rho^b$  is winning for all winning states of the Büchi game  $\mathcal{G}^b$ . Since winning strategies for Büchi games are subgame optimal, a shield that implements  $\rho^b$  ends deviations after the smallest number of steps possible, for all states of the design in which a finite number of deviations can be guaranteed. Second,  $\rho^b$  is cooperatively winning in the Büchi game  $\mathcal{G}^b$ . Therefore, in all states in which a finite number of deviation cannot be guaranteed, a shield that implements the strategy  $\rho^b$  recovers with the help of the design as soon as possible.

The standard algorithm for solving Büchi games contains the computation of attractors; the i-th attractor for the system contains all states from which the system can "force" a visit of an accepting state in i steps. For all states  $g \in G^b$  of the game  $\mathcal{G}^b$ , the attractor number of g corresponds to the smallest number of steps within which the recovery phase can be guaranteed to end, or can end with the help of the design if a finite number of deviation cannot be guaranteed.

**Theorem 8.2.** Let  $\varphi = \{Q, q_0, \Sigma, \delta, F\}$  be a safety specification and |Q|

be the cardinality of the state space of  $\varphi$ . An admissible shield with respect to  $\varphi$  can be synthesized in  $\mathcal{O}(2^{|Q|})$  time, if it exists.

*Proof.* The safety game  $\mathfrak{S}^s$  and Büchi game  $\mathfrak{S}^b$  have at most  $m=(2\cdot 2^{|Q|}+|Q|)\cdot 2\cdot |Q|$  states and at most  $n=m^2$  edges.

Safety games can be solved in  $\mathcal{O}(m+n)$  time and Büchi games in  $\mathcal{O}(m\cdot n)$  time (Mazala, 2001).

# 9

# Verifiable Learning-Based Synthesis

This section provides an overview of incorporating learning techniques in policy synthesis and discusses an approach that merges concepts from formal methods and machine learning. We consider the challenging task of computing a policy for a partially observable Markov decision process (POMDP) that satisfies certain classes of specifications, first introduced in Section 7. The approach discussed in this section represents a policy using recurrent neural networks (RNNs). We then show how such a representation can be integrated with formal methods by extracting a policy that is compatible with state-of-the-art verification tools.

RNNs offer an effective policy representation for POMDPs due to their ability to effectively process sequential data. As opposed to the conventional feed-forward architectures present in artificial neural networks, in which the nodes in each layer are only connected to nodes in subsequent layers, recurrent architectures allow for backward connections between nodes. Such backward connections allow RNNs to create internal memory states, such as those in long short-term memory (LSTM)

This section incorporates the results from the following publications (Carr, Jansen, Wimmer, Serban, Becker, and Topcu, 2019; Carr, Jansen, and Topcu, 2020; Carr, Jansen, and Topcu, 2021).

architectures (Hochreiter and Schmidhuber, 1997), which infer temporal behavior from sequences of data (Pascanu *et al.*, 2014). Reinforcement learning research has shown that RNNs used in environments modeled by POMDPs perform well as black-box functions for either state or value estimators (Wierstra *et al.*, 2007; Bakker, 2001) or as control policies (Hausknecht and Stone, 2015; Heess *et al.*, 2015b).

In POMDPs that model agents in safety-critical environments, policies that are guaranteed to prevent unsafe behavior are necessary. The agent's behavior may have to obey more complicated specifications than maximizing an expected reward, such as reachability, liveness or, more generally, specifications expressed in temporal logic, e.g. linear temporal logic (LTL), see Section 4.2.

Verifying whether an agent following an RNN-based policy in a POMDP satisfies temporal logic specifications is, in general, hard. RNNs are complex structures that capture non-linear input-output relations (Mulder et al., 2015). To formally analyze how RNNs interpret sequences of data, Sherstinsky (2020) suggest fixing a defined sequence length for analysis and performing an unrolling procedure, which converts the RNN to a feed-forward neural network with the same number of layers as that defined length (Goodfellow et al., 2016). Checking whether the agent's behavior satisfies the specification for the set of all possible sequences of data with a given length in the POMDP is intractable (Meuleau et al., 1999).

The approach discussed in this section combines the representation power of RNNs from machine learning with the provable guarantees that are at the heart of formal verification. The latter can efficiently verify whether an agent following a given policy, typically in the form of an finite state controller (FSC) (Poupart and Boutilier, 2003; Junges et al., 2018), adheres to a temporal logic specification (Baier and Katoen, 2008). However, directly synthesizing a policy requires—in general—memory of exponential size in the number of POMDP states (Baier et al., 2012). Machine learning, on the other hand, provides an efficient approach, in the form of training RNN-based policy representations from sequences of data, to find candidate policies that might ensure an

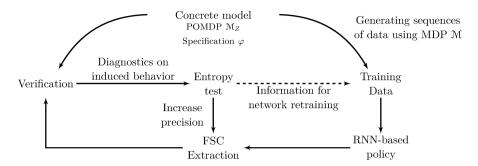


Figure 9.1: High-level iterative policy improvement process.

agent in a POMDP satisfies a temporal logic specification (Heess  $et\ al.$ , 2015a).

There remains a central gap: How to close the loop between training an RNN-based policy and efficiently verifying for a candidate policy? The approach closes this gap by tightly integrating formal verification and machine learning towards three key steps: (1) extracting an FSC from an RNN-based policy, (2) verifying this candidate FSC for the POMDP against a temporal logic specification, and (3) if needed, either refining the FSC or generating more training data for the RNN, see Figure 9.1.

# 9.1 Verifiable Recurrent Neural Network-Based Policies for Partially Observable Markov Decision Processes

In this section, we formulate a problem that is similar to that introduced in Problem 7.1. We restate it here without the robustness considerations:

**Problem 9.1** (POMDP synthesis). For a POMDP  $\mathcal{M}_Z$  and a specification  $\varphi$ , where either  $\varphi = \mathbb{P}_{\sim \lambda}(\psi)$  for  $\sim \in \{<, \leq, \geq, >\}$  and  $\lambda \in [0, 1]$  with  $\psi$  an LTL specification, or  $\varphi = \mathbb{E}_{\sim \lambda}(\lozenge a)$ , the problem is to determine a (finite-memory) policy  $\sigma \in \Sigma_Z^{\mathcal{M}}$  such that  $\mathcal{M}_Z^{\sigma} \models \varphi$ .

In Problem 9.1, if no (finite-memory) policy exists, the problem is infeasible. Note that, in general, Problem 9.1 is undecidable (Madani et al., 1999) and each method is necessarily incomplete.

Figure 9.1 outlines the learning-based overall approach to Problem 9.1. A trained RNN serves as an efficient representation of a POMDP policy. As safety-critical scenarios necessitate a *sound* notion of correctness, the approach evaluates such an RNN-based policy using formal verification against LTL specifications. There are four main building blocks towards that approach: (1) *Training* the RNN, (2) *extracting* FSCs as a tractable representation of the RNN-based policy, (3) *evaluating* the policy, and (4) *improving* the policy.

#### 9.1.1 Training a Recurrent Neural Network-Based Policy

We first define how RNNs can be used as a representation of POMDP policies.

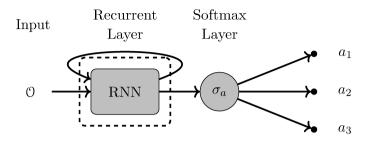
**Definition 9.1** (RNN-based policy). An RNN-based policy for a POMDP is a function  $\hat{\sigma} : \mathsf{ObsSeq}_{fin}^{\mathcal{M}_Z} \mapsto Distr(Act)$ , where  $\mathsf{ObsSeq}_{fin}^{\mathcal{M}_Z}$  is the set of all sequential observation-action inputs and Distr(Act) is the set of all distributions over actions. To be more precise, we identify the main components of such a network. An RNN-based policy  $\hat{\sigma}$  is sufficiently described by a  $hidden\text{-state update function } \hat{\delta} : \mathbb{R} \times Z \times Act \to \mathbb{R}$  and an  $action\text{-mapping } \sigma_h : \mathbb{R} \to Distr(Act)$ .

Consider the following observation-action sequence:

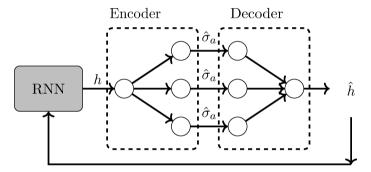
$$O(\pi) = O(s_0) \xrightarrow{Act_0} O(s_1) \xrightarrow{Act_1} \cdots O(s_i)$$
 (9.1)

The RNN-based policy receives an observation-action sequence and returns a distribution over the action choices. Throughout the execution of the sequence, the RNN holds a continuous hidden state  $h \in \mathbb{R}$ , typically described as an internal memory state, which captures previous information. On each transition, this hidden state is updated to include the information of the current state and the last action taken under the hidden-state update function  $\hat{\delta}$ . From the prior observation sequence in (9.1), the corresponding hidden state sequence would be defined as:

$$\hat{\delta}(\pi) = h_0 \xrightarrow{Act_0, \ O(s_1)} h_1 \xrightarrow{Act_1, \ O(s_2)} \cdots h_i$$



**Figure 9.2:** RNN-based Policy  $\hat{\sigma}$  with softmax layer activation  $\sigma_a$ .



**Figure 9.3:** RNN-based policy structure with a QBN. RNN block and associated QBN of  $B_h = 3$  with quantized activation  $\hat{\sigma}_a : \mathbb{R} \to \{-1, 0, 1\}$ .

Additionally, the output of the RNN-based policy is expressed by the action-distribution function  $\sigma_h$ , which maps the value of hidden state to an action mapping  $\sigma_h$ . At an internal memory state  $h_i$ , we have  $\hat{\delta}(h_i, O(s_i), a_i) = h_{i+1}$  and  $\sigma_h(h_{i+1}) = \mu(Act)$  for state  $s_i$  on path  $\pi$ .

The approach constructs an RNN-based policy using a three-layer network, shown in Figure 9.2. The policy network uses an LSTM architecture (Hochreiter and Schmidhuber, 1997) for the recurrent layer  $\hat{\delta}$  and then a softmax layer for the output action mapping  $\sigma_h$  (see Definition 9.1). To fit the RNN model to the sequences of training data, we use the Adam optimizer Kingma and Ba, 2015 with a categorical cross-entropy error function Goodfellow et al., 2016.

In practice, a machine learning-based approach requires a method of generating data. One such method involves first computing a policy  $\sigma \in \Sigma^{\mathcal{M}}$  of the underlying Markov decision process (MDP)  $\mathcal{M}$  that satisfies  $\varphi$  using the STORM probabilistic model checker (Dehnert  $et\ al.,\ 2017$ ). Then it samples an initial state uniformly over the initial belief support  $s_0 \in supp(b)$  and generate finite observation paths  $\pi_o$ , thereby creating multiple trajectory trees (Kearns  $et\ al.,\ 1999$ ). When generating sequences of data, selecting one of the trees and following it to a leaf, which forms either at a pre-defined maximum sequence length or a deadlock, gives a finite path  $\pi \in \mathsf{Paths}^{\mathcal{M}}_{fin}$ . From this path  $\pi$ , the method generates one possible observation-action sequence  $\pi_o \in \mathsf{ObsSeq}^{\mathcal{M}_Z}_{fin}$ .

**Example 9.1.** Consider the POMDP in Example 7.1 and Figure 7.1: a sample set of sequences of data would be:

 $\mathcal{D} = \{\pi_o^0 = (blue, \text{up}, blue, \text{down}, s_3), \ \pi_o^1 = (blue, \text{down}, s_3), \ \pi_o^2 = (blue, \text{up}, s_3)\}.$  An example RNN policy  $\sigma$  trained on these sequences would yield a policy for observation-action sequence  $\pi_{o,0} = (\text{blue})$  as  $\sigma(\pi_{o,0}) = \{0.67 : \text{up}, 0.33 : \text{down}\}$ , which has a categorical cross-entropy loss of approximately 0.585. Similarly, the same RNN policy for a longer observation-action sequence such as  $\pi_{o,1} = (\text{blue}, up, \text{blue})$  yields a policy  $\sigma(\pi_{o,1}) = \{1.0 : \text{down}\}$ , for a cross-entropy loss of 0.

## 9.1.2 Finite-State Policy Extraction

The discussed approach adapts a method called quantized bottleneck insertion (Koul et al., 2019) to extract an FSC from a given RNN-based policy. Let us first explain the relationship between the main components of an RNN-based policy  $\hat{\sigma}$  (Definition 9.1) and an FSC  $\mathcal{A}$  (Definition 7.2). In particular, the hidden-state update function  $\hat{\delta}$  takes as input a real-valued hidden state of the policy network, while the FSC's memory update function  $\delta$  takes a memory node from the finite set N. Figure 9.2 describes a simplified architecture for the former since its recurrent component acts as the hidden-state update function  $\hat{\delta}$ . The key for linking the two is therefore a mechanism that encodes the continuous hidden state h into a set N of discrete memory nodes. We outline such a mechanism in the sequel and in Figure 9.3 in which

we show the modified activation function (formed using an encoder and a decoder).

The approach leverages an autoencoder (Goodfellow et al., 2016) in the form of a quantized bottleneck network (QBN) (Koul et al., 2019). This QBN, consisting of an encoder and a decoder, is inserted into the RNN-based policy directly before the softmax layer, see Figure 9.3. In the encoder, the continuous hidden-state value  $h \in \mathbb{R}$  is mapped to an intermediate real-valued vector  $\mathbb{R}^{B_h}$  of pre-allocated size  $B_h$ . The decoder then maps this intermediate vector into a discrete vector space defined by  $\{-1,0,1\}^{B_h}$ . This process, illustrated in Figure 9.3, provides a mapping of the continuous hidden state h into  $3^{B_h}$  possible discrete values. We denote the discrete state for h by  $\hat{h}$  and the set of all such discrete states by  $\hat{H}$ . Note, that  $|\hat{H}| \leq 3^{B_h}$  since not all values of the hidden state may be reached in an observation sequence.

After the QBN insertion, we simulate a series of executions querying the modified RNN for action choices on the POMDP. These executions form a dataset of consecutive pairs  $(\hat{h}_t, \hat{h}_{t+1})$  of discrete hidden states, the action  $Act_t$  and the observation  $\mathcal{O}_{t+1}$  that led to the transition  $\{\hat{h}_t, Act_t, \mathcal{O}_{t+1}, \hat{h}_{t+1}\}$  at each time t during the execution of the RNN-based policy. The number of accessed memory nodes  $N \subseteq \hat{H}$  corresponds to the number of different discrete states  $\hat{h} \in \hat{H}$  in this dataset. The deterministic memory update rule  $\delta(n_t, Act_t, \mathcal{O}_{t+1}) = n_{t+1}$  is obtained by constructing a  $N \times (|Z| \times |Act|)$  transaction table (Koul et al., 2019). We can additionally construct the action-mapping  $\alpha \colon N \times Z \to Distr(Act)$  with  $\alpha(n_t, \mathcal{O}_t) = \mu \in Distr(Act)$  by querying the softmaxoutput layer (see Figure 9.2) for each memory node and observation.

### 9.1.3 Evaluating the Extracted Policy

We assume that for POMDP  $\mathcal{M}_Z = (\mathcal{M}, Z, O)$  and specification  $\varphi$ , we have an extracted FSC  $\mathcal{A}_{\hat{\sigma}} \in \Sigma^{\mathcal{M}_Z}$  as in Definition 7.2. The approach applies the policy  $\mathcal{A}_{\hat{\sigma}}$  to obtain an induced discrete-time Markov chain (DTMC)  $\mathcal{M}_Z^{\mathcal{A}_{\hat{\sigma}}}$ . For this DTMC, formal verification through model checking checks whether  $\mathcal{M}_Z^{\mathcal{A}_{\hat{\sigma}}} \models \varphi$  and thereby provides hard guarantees about the quality of the extracted FSC  $\mathcal{A}_{\hat{\sigma}}$  regarding  $\varphi$ . In particular,

(probabilistic) model checking provides the probability – or the expected reward – to satisfy a specification for *all states*  $s \in S$  via solving linear equation systems, see Section 6 and (Baier and Katoen, 2008).

**Example 9.2.** Consider the case in the 1-FSC  $\mathcal{A}_1$  from Example 7.1 (Figure 7.1b) where the parameter p=1 and the probability of reaching the state  $s_3$  in the induced DTMC is  $\Pr(\diamondsuit s_3) = \frac{1}{3}$ . The behavior induced by this 1-FSC violates the specification and we obtain two counterexamples of critical memory-state pairs for this policy  $\mathcal{A}_{\hat{\sigma}}$ :  $(0, s_0)$  and  $(0, s_1)$ .

If the specification does not hold, the policy may require refinement. On the one hand we can increase the number of memory nodes  $B_h$  to extract a new FSC. On the other hand, we may decide via a formal entropy check whether new data need to be generated.

#### 9.1.4 Improving the Policy with Counterexample Data

The approach attempts to determine whether an RNN-based policy requires more training data  $\mathcal{D}$  or not. Existing approaches in supervised learning methods leverage benchmark comparisons between a train-test set using a loss function (Baum and Wilczek, 1987). Loss visualization, proposed by Goodfellow and Vinyals (2015) provides a set of analytical tools to show model convergence. However, such approaches aim at continuous functions instead of the discrete representations as in the FSC. More importantly, the method leverages the information gained from a model-based approach.

We first determine a set of states that are critical for the satisfaction of the specification under the current policy. Consider a sequence of memory nodes and observations  $(n_0, \mathcal{O}_0) \xrightarrow{a_0} \cdots \xrightarrow{a_{t-1}} (n_t, \mathcal{O}_t)$  from the POMDP  $\mathcal{M}_Z$  under the FSC  $\mathcal{A}_{\hat{\sigma}}$ . For each of these sequences, we collect the states  $s \in S$  underlying the observations, e.g.,  $O(s) = \mathcal{O}_i$  for  $0 \le i \le t$ . As we know the probability or expected reward for these states to satisfy the specification from previous model checking, we can now directly assess their criticality regarding the specification. We collect all pairs of memory nodes and states from  $N \times S$  that contain

critical states and build the set  $Crit_{\mathcal{A}_{\sigma}}^{\mathcal{M}_{Z}} \subseteq N \times S$  that serves us as a counterexample. These pairs carry the joint information of critical states and memory nodes from the policy applied to the DTMC  $\mathcal{M}_{Z}^{\mathcal{A}_{\sigma}}$ .

**Entropy measure.** The average entropy across the distributions over actions at the choices induced by the counterexample set  $Crit_{A_{\hat{\sigma}}}^{M_Z}$  is our measure of choice to determine the level of training for the RNN-based policy. Specifically, for each pair  $(n,s) \in Crit_{A_{\hat{\sigma}}}^{M_Z}$ , we collect the distribution  $\mu \in Distr(Act)$  over actions that  $A_{\hat{\sigma}}$  returns for the observation O(s) when it is in memory node n. Then, we define the evaluation function H using the entropy  $\mathcal{H}(\mu)$  of the distribution  $\mu$ :

$$H \colon Crit_{\mathcal{A}_{\hat{\sigma}}}^{\mathcal{M}_Z} \to [0,1] \text{ with } H(n,s) = \mathcal{H}(\mu).$$

For high values of H, the distribution is uniform across all actions and the associated RNN-based policy is likely extrapolating from unseen inputs.

We observe that when there are fewer samples and higher memory nodes, the extracted FSC tends to perform arbitrarily, see Section 9.1.5 and Figure 9.7 for a detailed empirical analysis. We lift the function H to the full set  $Crit_{A,z}^{\mathcal{M}_Z}$ :

$$H(Crit_{\mathcal{A}_{\hat{\sigma}}}^{\mathcal{M}_{Z}}) = \frac{1}{|Crit_{\mathcal{A}_{\hat{\sigma}}}^{\mathcal{M}_{Z}}|} \sum_{(n,s) \in Crit_{\mathcal{A}_{\hat{\sigma}}}^{\mathcal{M}_{Z}}} H(n,s).$$
(9.2)

We compare the average entropy over all decision-points of the counterexample against a constant threshold  $\eta \in [0,1]$ , that is, if  $H(Crit_{\mathcal{A}_{\sigma}}^{\mathcal{M}_{Z}}) > \eta$ , we will provide more training data. Vice versa, if  $H(Crit_{\mathcal{A}_{\sigma}}^{\mathcal{M}_{Z}}) \leq \eta$ , we increase the upper bound on the number of memory nodes in the FSC.

**Example 9.3.** Under the working Example 7.1, the policy  $\mathcal{A}_1$  was the 1-FSC with p=1 (Figure 7.1b), which produces two counterexample memory and state pairs:  $Crit_{\mathcal{A}_1}^{\mathcal{M}_Z} = \{(0,s_0),(0,s_1)\}$ . The procedure would then examine the policy's average entropy at these critical components  $(n,s) \in Crit_{\mathcal{A}_1}^{\mathcal{M}_Z}$ , which in this trivial example is given by

 $H(Crit_{\mathcal{A}_1}^{\mathcal{M}_Z}) = -p \log_2(p) - (1-p) \log_2(1-p) = 0$  from (9.2). The average entropy is below a prescribed threshold,  $\eta = 0.5$ , and thus we increase the number of memory nodes, which results in the satisfying FSC  $\mathcal{A}_2$  in Figure 7.1c.

# 9.1.5 Performance on Partially Observable Markov Decision Process Benchmarks

We evaluate the RNN-based synthesis on benchmark examples that are subject to either LTL specifications or expected reward specifications. For the former, we compare to the tool PRISM-POMDP (Norman *et al.*, 2017), and for the latter to PRISM-POMDP and the point-based solver SolvePOMDP (Walraven and Spaan, 2017).

For a fair comparison, instead of terminating the synthesis procedure once the policy satisfies the specification, we always iterate 10 times, where one iteration encompasses the (re-)training of the RNN-based policy using counterexamples, the FSC extraction, the evaluations, and the policy improvement. For instance, for a specification  $\varphi = \mathbb{P}_{\geq \lambda}(\psi)$ , we leave the  $\lambda$  open and seek to compute  $\mathbb{P}_{\max}(\psi)$ , that is, we compute the minimal probability of satisfying  $\psi$  to obtain a policy that satisfies  $\varphi$ . We cannot guarantee to reach that optimum, but we rather improve as far as possible within the predefined 10 iterations.

We created the following Python toolchain to realize the full RNN-based procedure, combining the state-of-the-art tools from deep learning with those from formal verification. First, we use the deep learning library Keras (Ketkar, 2017) to train the RNN-based policy from sequences of data. To evaluate policies, we employ the probabilistic model checkers PRISM (Norman et al., 2017) and STORM (Dehnert et al., 2017) for LTL and undiscounted expected reward respectively. We evaluated on a 2.3 GHz machine with a 12 GB memory limit and a specified maximum computation time of  $10^5$  seconds. In Table 9.1 TO/MO denote violations of the time/memory limit, respectively and Res. refers to the output value of the induced DTMC.

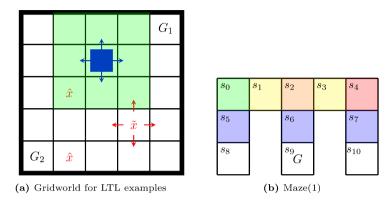


Figure 9.4: Physical environments for presented examples.

Problem settings with temporal logic specifications. We examined three problem settings involving motion planning with LTL specifications. For each of the settings, we use a square gridworld of length c with 4 action choices (cardinal directions of movement). The motivation for gridworld examples is that they provide a minimal safety check: a policy that fails to behave safely in such simple environments is also unlikely to behave safely in the real world (Leike  $et\ al.,\ 2017$ ). Inside this environment, there are a set of static  $(\hat{x})$  and moving  $(\tilde{x})$  obstacles as well as possible target cells  $G_1$  and  $G_2$ .

The agent has a limited visibility region, indicated by the green area, and can infer its state from observations and knowledge of the environment. We define observations as Boolean functions that take as input the positions of the agent and moving obstacles, see Figure 9.4a. Intuitively, the functions describe the 8 possible relative positions of the obstacles with respect to the agent inside its viewing range. The three problem settings are as follows:

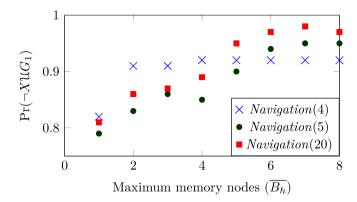
- 1. Navigation with moving obstacles—An agent and a single stochastically moving obstacle. The agent task is to maximize the probability to navigate to a goal state A while not colliding with obstacles (both static and moving):  $\varphi_1 = \mathbb{P}_{\max} (\neg X U G_1)$  with  $X = \hat{x} \cup \tilde{x}$ ,
- 2. Delivery without obstacles—An agent and static objects (land-

- marks). The task is to deliver an object from  $G_1$  to  $G_2$  in as few steps as possible:  $\varphi_2 = \mathbb{E}_{\max}(\lozenge(G_1 \land \lozenge G_2))$ .
- 3. Slippery delivery with static obstacles—An agent where the probability of moving perpendicular to the desired direction is 0.1 in each direction. It attempts to maximize the probability to travel from location  $G_1$  and to  $G_2$  without colliding with the static obstacles  $\hat{x}$ :  $\varphi_3 = \mathbb{P}_{\text{max}} (\Box \diamondsuit G_1 \land \Box \diamondsuit G_2 \land \neg \diamondsuit X)$ , with  $X = \hat{x}$ .

Problems settings for maximizing expected reward. For comparison to existing benchmarks, we extend a well-known POMDP example Maze(c) for an arbitrary-sized structure. These problems are quite different to the LTL examples, in particular the significantly smaller observation spaces, see (Carr  $et\ al.,\ 2021$ ) for details.

- 1.  $Grid-based\ Maze(c)\ with\ c+2\ rows$ —The agent can only detect its neighboring walls and attempts to reach a goal state G in as few steps as possible, see Figure 9.4b for Maze(1). Extra rows add uncertainty over the agent's position in the corridors, see the blue observations in Figure 9.4b.
- 2. **Grid(c)** with restricted vision—A square grid with length c where the agent attempts to reach a goal state G at the top right of the square. The agent is placed in the grid according to a uniform distribution of the states and can only observe its exact location when it reaches the goal state.

Increasing the number of memory nodes improves performance. In Figure 9.5, we show that increasing the number of memory nodes in the FSC produces higher performing policies, both in the form of higher probabilities of satisfying the specification and higher undiscounted expected rewards. A noticeable characteristic is that for each FSC in Figure 9.5, there is a point of diminishing returns where the additional memory does not produce higher quality policies. In most cases, this point falls between 6 and 8 memory nodes. As a consequence for the set of benchmarks, unless otherwise specified, we set the upper bound for the number of memory nodes at  $\overline{B_h} = 8$ .



**Figure 9.5:** Probability of satisfying the specification for agent employing k-FSCs for *Navigation* as the number of memory nodes increases. Each point represents the best performing policy over 10 iterations of the in Figure 9.1. Each FSC begins to experience diminishing returns at values of  $B_h = 6$  and higher.

Using counterexample data improves the quality of the extracted policies. Figure 9.6 compares the number of critical states in a set of counterexamples in relation to the probability of satisfying an LTL specification in each iteration of re-training for the proposed method. In particular, we depict the size of the set of critical states  $Crit_{\mathcal{A}_{\gamma}}^{\mathcal{M}Z} \subset S$  regarding the specification  $\varphi$ . In Figure 9.6, as the satisfaction probability and the expected reward increases, the number of the critical states identified by the verification decreases. In particular, the retraining of the RNN-based policy on the sequences of data generated using the local improvement step is effective in improving the policy with each iteration.

**Using counterexample data generates policies that make less arbitrary decisions.** In Figure 9.7, we ignore the decision at the entropy check, fix the memory precision, and iteratively add more sequences of data generated using the counterexamples. Each point in Figure 9.7 represents one instance of verification in the loop in Figure 9.1. As the RNN-based policy iteratively trains on additional sequences of data, the subsequent extracted policy makes less arbitrary decisions, shown in Figure 9.7 by the decrease in entropy of the FSC as the RNN-based policy is trained

			PRISM-	POMDP	RNN-ba	ased Policy
Problem	States	Type, $\varphi$	Res.	Time (s)	Res.	Time (s)
Navigation (3)	333	$\mathbb{P}_{\max}^{\mathcal{M}_Z},  \varphi_1$	0.84	73.88	0.80	123.14
Navigation (4)	1088	$\mathbb{P}_{\max}^{\mathcal{M}_Z},  \varphi_1$	$0.93^{\dagger}$	1034.64	0.92	160.32
Navigation (5)	2725	$\mathbb{P}_{\max}^{\mathcal{M}_Z},  \varphi_1$	MO	MO	0.95	311.65
Navigation $(10)$	49060	$\mathbb{P}_{\max}^{\mathcal{M}_Z}, \varphi_1$	MO	MO	0.85	2561.02
Navigation $(20)$	798040	$\mathbb{P}_{\max}^{\mathcal{M}_Z}, \varphi_1$	MO	MO	0.98*	8173.03*
Navigation $(30)$	4045840	$\mathbb{P}_{\max}^{\mathcal{M}_Z}, \varphi_1$	MO	MO	0.97*	61350.34*
Navigation $(40)$	_	$\mathbb{P}_{\max}^{\mathcal{M}_Z}, \varphi_1$	MO	MO	ТО	TO
Delivery (4)	80	$\mathbb{E}_{\max}^{\mathbb{M}_Z},  \varphi_2$	-6.0	28.53	-6.04	94.32
Delivery (5)	125	$\mathbb{E}_{\max}^{MZ}, \varphi_2$	-8.0	102.41	-8.13	150.44
Delivery (10)	500	$\mathbb{E}_{\max}^{\mathbb{M}_Z},  \varphi_2$	MO	MO	-18.13	347.98
Slippery (4)	460	$\mathbb{P}_{\max}^{\mathcal{M}_Z},  \varphi_3$	0.90	5.10	0.80	180.15
Slippery (5)	730	$\mathbb{P}_{\max}^{\mathcal{M}_Z}, \varphi_3$	0.93	83.24	0.89	212.79
Slippery (10)	2980	$\mathbb{P}_{\max}^{\mathcal{M}Z}, \varphi_3$	MO	MO	0.98	280.55
Slippery (20)	11980	$\mathbb{P}_{\max}^{\mathcal{M}Z}, \varphi_3$	MO	MO	0.99	2384.56
Maze (1)		$\mathbb{E}_{\max}^{\mathcal{M}_Z}$	-4.30	0.09	-4.33	80.31
Maze (2)		$\mathbb{E}_{\max}^{\mathcal{M}_Z}$	-5.23	2.176	-5.34	114.23
Maze (5)		$\mathbb{E}_{\max}^{\mathcal{M}Z}$	$-13.00^{\dagger}$	4110.50	-13.29	160.12
Maze (10)		$\mathbb{E}_{\max}^{MZ}$	MO	MO	-23.02	210.01
Grid (3)		$\mathbb{E}_{\max}^{\mathcal{M}_{Z}}$	-2.88	2.332	-2.90	87.31
Grid (4)		$\mathbb{E}_{ ext{max}}^{\mathfrak{M}Z}$	-4.13	1032.53	-4.20	124.31
Grid (5)		$\mathbb{E}_{ ext{max}}^{\mathfrak{M}_{Z}}$	MO	MO	-5.91	250.14
Grid (10)		$\mathbb{E}^{M_Z}$	MO	MO	-12.92	1031.21
Grid (25)		$\mathbb{E}_{\max}^{\mathbb{M}_Z}$	МО	MO	-35.32	6514.30

Table 9.1: Computing policies for examples with LTL specifications.

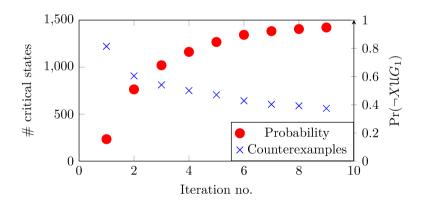
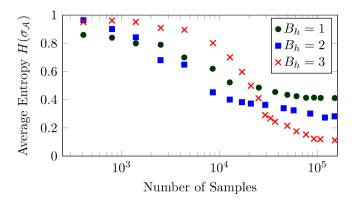


Figure 9.6: Progression of the number of critical states and the probability of satisfying LTL specification  $\varphi_1$  in Navigation(5)



**Figure 9.7:** Entropy of the extracted FSCs for *Navigation*(5) from an RNN as it is trained with more samples. Each point represents an extracted FSC, for color sequence we fix the discretization and add more samples guided by the counterexamples.

on larger sets of training sequences.

Increasing the number of memory nodes generates policies that In Figure 9.7, we also compare how make less arbitrary decisions. the number of memory nodes in the extracted FSC correlates to the entropy of the decisions at critical states. When trained on a large set of sequences of training data, the FSCs with a higher number of memory nodes have a lower entropy than those without. This behavior is likely due the fact that extracted FSC with more memory nodes can better approximate the RNN-based policy, which itself is making less arbitrary decisions due to the larger training set. Meanwhile, when the extracted FSCs are approximating RNN-based policies trained on smaller sets of training sequences, they generally make arbitrary decisions (see top left of Figure 9.7). In these cases, the FSC with more memory nodes tend to make more arbitrary decisions than those with less, which is likely a function of an under-defined hidden state update  $\hat{\delta}$  in the RNN-based policy.

Limiting the number of memory nodes creates a precision-performance trade-off. Increasing the number of memory states in the FSC produces policies with higher probabilities of satisfying the specification

and greater expected rewards. In Table 9.1, we include the sizes of the FSCs for the handcrafted procedure to demonstrate the trade-off between computational tractability and expressivity: a larger FSC means that the policy can store more information, which may lead to better decisions. However, larger FSCs require more computational effort and may require more sequences of data for training the RNN-based policy. Figure 9.8 shows the automatically extracted FSC for the Maze(1) environment. Note that a 2-FSC can represent the optimal Maze(1) policy. The FSC shown in Figures 9.8a and 9.8b is very close to this optimal policy. The stochastic action choices at  $(n_0, blue)$  and at  $(n_1, yellow)$  create the suboptimality in this example with the optimal policy taking the respective up and right actions at these points.

#### 9.2 Case Study: Autonomous Driving in Traffic Lights

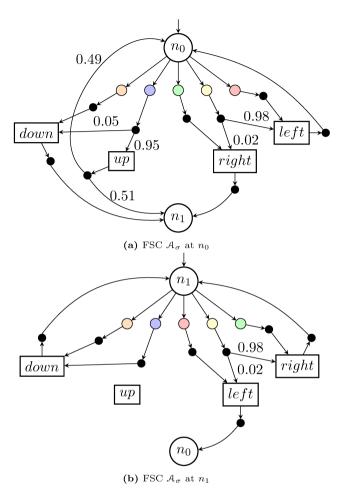
Consider the scenario, pictured in Figure 9.9, of an autonomous vehicle operating in a city with the following sensors:

- 1. Navigation System (GPS),
- 2. Optical Camera (Traffic Light Identification),
- 3. Lidar (Proximity),
- 4. Radar (Speed).

At each intersection is a set of traffic lights that restrict the available action choices of the vehicle:

- Red:  $\{Stop\}$ ,
- Yellow:  $\{Stop, Straight\},\$
- Green:  $\{Stop, Straight, Left, Right\}$ .

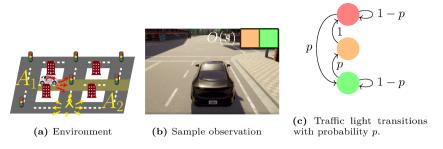
The state of each traffic light can be modeled as a Markov chain where the system switches from red to green and green to orange with probability p=0.1, see Fig. 9.9c. In this environment there is a pedestrian, who moves stochastically but with reduced speed. Each turn, the vehicle takes two actions while the pedestrian can move one



**Figure 9.8:** The corresponding FSC  $A_{\gamma}$  for Maze(1) example. The agent's initial state  $s_I \in S \setminus \{s_9\}$  is allocated over a uniform distribution and each color represents a different observation. The FSC  $A_{\gamma}$  has two memory nodes  $(n_0 \text{ and } n_1)$ , we prune action mappings and memory updates with low probabilities from 9.8a and 9.8b.

intersection at a time. While the vehicle must move according to the status of the lights, the pedestrian is under no such restriction and will ignore them.

**State space description.** A state in this system is  $s = (x, y, \theta, obs_x, obs_y, L)$ , which is the location and direction of the vehicle and the location of



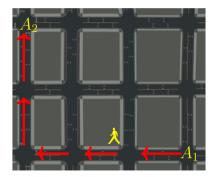
**Figure 9.9:** Autonomous car operating in an urban environment modeled as a gridworld and traffic lights. Each intersection has a decision point for the car who moves twice for each pedestrian move.

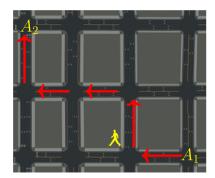
the pedestrian as well as the status of the lights L at each intersection.

Partial observability for the vehicle. The car is only able to see the traffic lights that are directly in front of it. For example in Figure 9.9a, the vehicle facing to the right can observe the lights of two intersections (the middle and the center right). The vehicle uses its optical sensors to determine the status of the lights. The task of navigating from an initial position  $A_1$  to a goal location  $A_2$  is modeled as a POMDP, with a belief on each status of the unseen intersection lights.

**Specifications for the vehicle.** We demonstrate the effect of different specifications on feature splitting. In particular, we describe three scenarios where the POMDP problem differs based on the safety task and the sensors that autonomous vehicle uses to measure the relevant features.

- Fastest trip to goal—The vehicle attempts to navigate the changing lights to ensure the vehicle reaches the goal location  $A_2$   $(\psi = \diamondsuit A_2)$ .
- Pedestrian never follows car—The vehicle attempts to ensure the pedestrian cannot read the number plate of the back of the car. To read the plate the pedestrian must be behind the car (for example when the car faces east the label is defined by  $B = ((x obs_x) = 1 \land (y obs_y) = 0))$  for two consecutive





- (a) A violating path in Carla environment
- (b) A satisfying path in Carla environment

**Figure 9.10:** Example paths in the autonomous vehicle case study. The vehicle must avoid speeding past the pedestrian which would violate specification  $\varphi_G$ .

turns—described by  $\psi_F = \neg \diamondsuit (B \land \bigcirc B) \mathcal{U} A_2$ . The vehicle uses the lidar proximity sensor to construct the relevant feature that determines its relative position to the pedestrian.

• Vehicle cannot speed within vicinity of pedestrian—The vehicle moves twice for every one action cycle. Accordingly, we define "speeding" as moving straight  $a_1$  twice through the green lights. For this specification, we utilize a lidar to sense the relative position as well as the radar to test speed.  $\psi_G = \neg \diamondsuit ((a_1 \land \bigcirc a_1) \land \text{ped} \mathcal{U} A_2)$  where ped is when the pedestrian within one space of the vehicle, defined by  $\text{ped} = |x - obs_x| + |y - obs_y| \le 1$ .

Figure 9.10 shows the difference between a vehicle path that violates  $\psi_G$  and one that satisfies  $\psi_G$ . In Figure 9.10a, the vehicle at the intersection with the pedestrian, having previously taken straight, takes the same action again. This action causes the label  $((a_1 \wedge \bigcirc a_1) \wedge \text{ped})$  to be True and since it occurs prior to the vehicle reaching  $A_2$ , this path violates specification  $\psi_G$ . In Figure 9.10b, an alternate satisfying path has the vehicle turning right at this intersection.

### 9.2.1 Experiments on the CARLA High-Fidelity Driving Simulator

For the high-fidelity autonomous vehicle simulation, we implement the computed policies on the open-source simulator CARLA (Dosovitskiy

Case (grid size)	Spec.	Size (states)	Time (s)	Result $V^{\sigma}(b_0)$
Autonomous Vehicle (3)	$\psi$	2916	50.65	3.75
Autonomous Vehicle (3)	$\psi_F$	26244	350.23	6.42
Autonomous Vehicle (3)	$\psi_G$	26244	475.23	4.5
Autonomous Vehicle (4)	$\psi$	708588	-TO-	-TO-
Autonomous Vehicle (4)	$\psi_F$	6377292	-TO-	-TO-
Autonomous Vehicle (4)	$\psi_G$	6377292	-TO-	-TO-

Table 9.2: Synthesis times for autonomous vehicle POMDP



**Figure 9.11:** Output policy for vehicle at s = (2, 0, west, 2, 0, L) (or at t = 1 in Figure 9.10). The straight action would result in a violation of specification  $\varphi_G$  while taking right leads to a path like that in Figure 9.10b.

et al., 2017). We run CARLA 0.9.11 on a 3.1 GHz machine with a GeForce RTX2060 graphics and 32GB of memory. Table 9.2 contains the model sizes, synthesis compute times and the expected values for city-grids of size 3 and 4.

In Figure 9.11, we see the synthesized policy for the situation that occurs at the pedestrian intersection in Figure 9.10. The computed policy rules out the choices of straight and left.

## 10

#### Some Future Directions

We have seen how formal methods can assist with validating system models and generating provable guarantees. However, the increasingly complex structures necessary for implementing today's systems open new problems. In this section, we examine some paths for further developing formal methods that still contain a rich set of questions to answer. In fact, we show that formal methods must intersect with other areas to produce formalisms and tools for better validating and confidently assuring engineered systems. These intersections include but are not limited to fields such as learning, security, regulation, and certification.

## 10.1 Formal Methods for Reinforcement Learning

Reinforcement learning (RL) algorithms search for policies that are optimal with respect to user-specified objectives. These algorithms allow for goal-oriented descriptions of complex behaviors and they provide a high degree of flexibility; they can be applied even when the system dynamics are high-dimensional, stochastic, and unknown (Sutton and Barto, 2018; Bertsekas, 2019; Powell, 2022).

Recently, deep RL algorithms—which use neural networks to parameterize value and policy functions—have demonstrated empirical success in a variety of applications: e.g., controlling plasma configurations in a nuclear fusion reactor (Degrave et al., 2022) and playing Chess, Shogi, and Go at superhuman levels (Silver et al., 2018). In these examples, the application of neural networks enables approximate solutions to problems in decision and control that would otherwise not be possible; the state and action spaces of these examples are too large for exact solution computed via dynamic programming algorithms.

However, there remain barriers to the deployment of RL algorithms in many engineering applications. Autonomous vehicles, power systems management, and robotic systems are examples of complex application domains that require strict adherence of the system's behavior to stakeholder requirements. However, the verification of RL policies is difficult. This is particularly true for deep RL algorithms, which typically only output the learned policy and its estimated value function, rendering their resulting behaviors opaque to further verification and analysis. The introduction of techniques borrowing from formal methods is necessary for the development of RL algorithms that yield behaviors with verifiable properties. To achieve this aim, we require frameworks to incorporate verifiable models into RL algorithms.

For example, given some temporal logic specification Alshiekh *et al.* (2018) synthesize reactive systems called *shields*, which prevent reinforcement learning systems from taking unsafe actions with respect to the specification. Meanwhile, a number of works have studied the use of RL as a method of controller synthesis for temporal logic objectives (Hahn *et al.*, 2019; Hasanbeig *et al.*, 2019b; Bozkurt *et al.*, 2020; Hasanbeig *et al.*, 2019a). While Wen *et al.* (2017) and Djeumou *et al.* (2021) use linear temporal logic formulas as *side information* that constrains the outputs of inverse reinforcement learning algorithms.

A general framework that brings verifiable models into reinforcement learning algorithms is that of the *reward machine*—a type of finite state machine used to encode reward functions in RL problems (Toro Icarte *et al.*, 2022). These structured task representations can encode

non-markovian tasks, and they can be exploited to improve learning efficiency (Icarte et al., 2018; Xu et al., 2020; Icarte et al., 2019) and to automatically translate specifications given in formal languages, such as linear temporal logic, into reward functions (Camacho et al., 2019). In the context of multiagent RL, Neary et al. (2021) use reward machines to specify multi-agent tasks, and to decompose these tasks into subtasks to be learned by individual agents through decentralized RL algorithms. The authors use the structure of the reward machines to prove conditions under which the resulting learned behavior is guaranteed to accomplish the original task.

Beyond the framework of reward machines, compositional design of RL systems can be used to greatly reduce the complexity of, and to more easily verify, individual subsystems. By creating well-defined interfaces between subsystems, system-level requirements may be decomposed into component-level ones. Conversely, each component may be developed and tested independently, and the satisfaction of component-level requirements may then be used to place assurances on the behavior of the system as a whole. Towards these ends, Neary et al. (2022) propose a framework to verify compositional RL systems against probabilistic task specifications. The framework builds a high-level system model, represented as a Markov decision process, which is used for high-level planning and to automatically generate subtask specifications for the low-level subsystems, each of which is implemented as an independent RL agent. Jothimurugan et al. (2021) similarly build compositional RL systems by using a graph-based representation of high-level tasks, and by using RL-based controllers to accomplish all necessary subtasks.

The above references provide examples of ways in which properties of RL-trained policies may be verified. Future work that continues to develop such frameworks and that applies them in experiments is necessary for the eventual deployment of trustworthy autonomous systems that incorporate RL-trained components.

# 10.2 Operating Under Limited Information and Concealing Information

The operation of autonomous systems relies on the information flow both within the components of the system and between the system and its environment. Most of the formal methods for autonomous systems discussed in this review inherently assume that this information flow is perfect. For example, numerical operations can be carried out with full accuracy, the sensor inputs are not quantized, and the communication channels are not noisy. However, in reality, the information is distorted in many ways due to both the system's internal design and the environmental factors. For example, cyber-physical systems typically have bandwidth limitations that require the sensor and controller outputs to be quantized (Franceschetti and Minero, 2014). In addition to these naturally occurring distortions, the information flow to the system may be adversarially modified.

The dependency on information brings multiple questions: What are the possible sources of distortion? What is the lowest amount of information that can ensure the safe operation of a system? How can we modify the existing formal methods for autonomous systems to account for these distortions? As an answer to these questions, early works from control theory show that the stability of a dynamical system imposes a lower bound on the information rate of a system (Nair and Evans, 2004; Franceschetti and Minero, 2014). For the discrete dynamical systems, Tanaka et al. (2021) and Eysenbach et al. (2021) provided methods that limit the information flow between the components of a system that is modeled with a Markov decision process. In the context of multi-agent systems, Wang et al. (2020) and Karabag et al. (2022) minimize the information shared between the agents to improve the performance under communication loss. In the case of an adversarial corruption of information flow in deterministic systems, we can represent the adversary as an additional player, represent the synthesis problem as a two-player game, and utilize the reactive synthesis methods mentioned in Section 5.1.

On the flip side, an autonomous system must not leak critical informa-

tion about itself to maintain the success of its operation. There is a growing literature on minimizing the information leakage of a system by considering concepts such as opacity (Saboori and Hadjicostis, 2007; Tong et al., 2018; Bérard et al., 2015), privacy (Such, 2017; Gohari et al., 2020), deceptiveness (Zhang and Zhu, 2018; Karabag et al., 2021), and estimation error (Farokhi and Sandberg, 2017; Karabag et al., 2019). There is a trade-off between minimizing the information leakage by considering such concepts and the performance of the system. However, minimizing the information leakage gives robustness to the system against its adversaries, thereby maximizing the performance in the long run.

### 10.3 Safeguarding Information Privacy in Autonomous Decision-Making Systems

In this survey, we mainly addressed the safety implications of autonomous decision-making systems. In particular, we surveyed papers that generally study the problems of policy synthesis and verification with respect to certain mathematically and formally specified objectives; for example, we frequently revisited the case study of an autonomous vehicle that must avoid crashing into obstacles and adhere to certain traffic rules formulated via temporal logic. While such developments are crucial to safeguarding the safety of the individuals whose daily lives will be affected by the deployment of autonomous decision-making systems, the societal impacts of these systems may extend far beyond the matters discussed in this survey; these systems often incorporate confidential, proprietary, operational, personal, or otherwise sensitive information in their decision-making algorithms, which raises privacy concerns.

Markov decision processes have been a major part of the formulation of many policy synthesis and verification problems discussed in this paper. Gohari et al. (2020) study a policy-synthesis problem in which it is within the privacy interests of a decision-maker to keep the transition probabilities of the underlying Markov decision process confidential while publicly taking actions according to synthesized policy. The paper

uses the framework of differential privacy to obfuscate the transition probabilities and then synthesizes a policy based on the obfuscated transition probabilities using dynamic programming. Then, the differential privacy of the overall synthesis algorithm becomes immediate due to differential privacy's immunity to post-processing. Although the proposed policy-synthesis algorithm addresses some of the named privacy concerns, it is not clear how the algorithm can satisfy a set of safety specifications, especially that differentially private algorithms are known to often trade off utility. Future work that incorporates the approaches discussed in this survey may be a potential solution to this issue.

Recall that, in Section 9, we reviewed an approach for policy synthesis and verification in which the policy is represented via recurrent neural networks. There exists mounting evidence that processing data in the form of training neural networks has privacy ramifications for the training data (Mireshghallah et al., 2020). Yang et al. (2021) show that the privacy ramifications of training recurrent neural networks can be worse than conventional feed-forward neural networks, especially in the task of deep reinforcement learning which is closely related to the policy-synthesis problem discussed in Section 9. The authors further study mitigation methods that leverage the promise of differential privacy by perturbing the trainable parameters of the neural network under protection. In this case, it must be further studied how the perturbations in the name of differential privacy affect a policy's ability to satisfy a given set of safety specifications.

### 10.4 Explainability in Verification and Synthesis

Formal methods such as model checking (Baier and Katoen, 2008) are capable of verifying human-generated robotic mission plans against a set of requirements (Humphrey and Patzek, 2013). In cases in which the plans may violate the requirements, such techniques generate *counterexamples* that illustrate requirement violations and provide valuable diagnostic information (Wimmer *et al.*, 2014; Feng *et al.*, 2016a). Nevertheless, these artifacts may be too complex for humans to understand,

because existing notions of counterexamples are defined as either a set of finite paths or an automaton typically with large number of states and transitions.

Counterexample generation for model checking Markov decision processes has been studied in several works using different representations of counterexamples: Han et al., 2009 computes the smallest number of paths in a Markov decision process whose joint probability mass exceeds the threshold and formulates the counterexample generation as a k-shortest path problem; (Wimmer et al., 2014) computes a critical subsystem of a Markov decision process with the minimal number of states and proposes solutions based on mixed-integer linear programming (MILP) and SAT-modulo-theories. There are several attempts to generate human-readable counterexamples: (Wimmer et al., 2013) computes a minimal fragment of model description in some high-level modeling language (e.g., probabilistic guarded commands), while (Feng et al., 2016a) computes structured probabilistic counterexamples as a sequence of "plays" that capture the high-level objectives in UAV mission planning. Feng et al. (2018) define a notion of explainable counterexample, which includes a set of structured natural language sentences to describe the system behavior that lead to a requirement violation in an MDP model of a robotic mission plan. They propose an approach based on MILP for generating explainable counterexamples that are minimal, sound and complete.

The study of natural language for robotic applications has mostly focused on translating human instructions expressed in natural language to robotic control commands. For example, Hayes and Shah (2017) considers the problem of synthesizing natural language descriptions of robotic control policy. Lignos et al. (2015) present an integrated system for synthesizing reactive controllers using natural language specifications which are translated into linear temporal logic formulas. If unsynthesizable, the minimal unsynthesizable core is returned as a subset of the natural language input specifications.

#### 10.5 Regulation

Regulatory requirements have to be incorporated when deriving specifications for safety-critical systems. Autonomous driving, for instance, is one of the fast-growing domains where regulation plays a significant role. Regulatory requirements, in this case, include the rules of the road and traffic regulations. Some rules, however, are ambiguously stated and include implicit assumptions related to the capability and rationality of human drivers. Such ambiguity potentially leads to inconsistent interpretation of the rules among different developers, complicates the design process, and ultimately jeopardizes the safety and efficiency of the system.

The lack of precise specifications has been acknowledged both by industry and academia (Shalev-Shwartz et al., 2017; Phan-Minh et al., 2019; Censi et al., 2019; Harel et al., 2022) and is a significant impediment towards formal methods realizing their full potential in this domain. Furthermore, without a clear description of how autonomous vehicles should behave, especially around dynamic objects such as humans and other vehicles, a common practice is to start with an implementation of specific behaviors and evaluate it based on simulation and testing. This approach is costly, as evidenced by the delay of deploying autonomous vehicles from the original estimate of 2020 (Anderson, 2020). Having a precise specification not only speeds up the development process but also ensures transparency and predictability of the system while giving the developers the flexibility of picking an implementation that leads to admissible behaviors as defined by the specification. Finally, it equips the regulators with an objective measure to validate autonomous vehicles.

Censi et al. (2019) make an initial attempt to tackle this problem by introducing frameworks for describing these specifications. The focus of these initial works is on the structure of the specifications that allows trajectories to be evaluated based on the violation of individual rules, which incorporate different factors such as safety, law, ethics, culture, etc. Systematic approaches to derive the specifications, taking into account these factors, however, remain an open problem.

#### 10.6 Dynamic Certification

In addition to more conventional versions of regulation that focus on regulatory requirements, autonomous systems require a more proactive approach to certification. Such dynamic approaches use formal methods as a form of design guide, to give a quick and inexpensive way of circling through different scenarios and contexts where the system ought to behave as expected. This type of dynamic certification is the iterative revision of permissible  $\langle use, context \rangle$  pairs for a system—rather than prespecified tests that a system must pass to be certified (Bakirtzis et al., 2022b), meaning that we can only certify for particular operational requirements and environments and can never fully guarantee any system-level requirements preserve in any deployment scenario.

Specifically, dynamic certification is based on applying testing that is informative to the different stages of a systems lifecycle, from early to transitional to finally confirmatory. These testing phases are not executed linear and testing is not expected to certify for *any* context. Rather, the expectation of testing, compared to more traditional methods within certification where the goal is *to be* certified, is continual and inputs of new contexts must again be tested through a series of living documents (e.g., capturing changing requirements within formal models), simulations (e.g., congested environments vs. open areas), and controlled environments before deployment (e.g., testing car tracks).

The output of dynamic certification is acceptable scenarios of operation, including but not limited to identified and potentially mitigated hazardous states, admissible environments of operation (e.g., congested vs. non-congested streets for an autonomous car), and modified system-level requirements. While the word output might suggest that the process terminates, the expectation is closer to "terminates yet starts again": Given that autonomous systems operate in changing environments it is natural to suggest that dynamic certification must apply to those changes. Having done dynamic certification for a (use, context) pair, however, only needs to document and test the changed parameters and do a form of regression testing in models, simulations, and controlled environments to make sure those changes do not break our overall

201

assumptions and guarantees.

Autonomous systems rely on two types of decision-making: design and deployment decisions. Design decisions concern the structure and intended operation of the autonomous system, which can include criteria and metrics of correct behavior and revisions to system requirements, behaviors, and architectures (reflected to formal models, testings procedures, and controlled environment parameters). Deployment decisions concern the contexts and uses for the autonomous system, including what hazardous states ought to be eliminated or mitigated against. Dynamic certification considers both.

Further research in autonomy is paramount to achieve a level of granularity that accounts for design and deployment decisions and their associated models and testing, particularly when dynamic certification is partially implemented with formal methods. Such thrusts can include finding more economical ways to compute formal models, addressing uncertainty to better match environmental parameters, and accounting for learning algorithms that are composed with traditional control-theoretic methods.

## Acknowledgements

The work this article reported has been carried out in collaboration with a large group of researchers. We are thankful to all of them.

The preceding sections draw material from several of our earlier publications heavily. We would like to express our appreciation to our co-authors on those publications: Erdem Arinc Bulgur, Kai-Wei Chang, Rayna Dimitrova, Lu Feng, Nils Jansen, Sebastian Junges, Joost-Pieter Katoen, Hadas Kress-Gazit, Ahmadreza Marandi, Richard Murray, Marnix Suilen.

Over the course of preparing this article, the authors have been supported by the following grants: NASA 80NSSC21M0071, NSF CNS-1836900, NSF 1646522, NSF 1652113, NSF 2141153, ARL ACC-APG-RTP W911NF1920333, AFRL FA9550-19-1-0169, and AFOSR FA9550-19-1-0005.

#### References

- Abadi, M. and L. Lamport. (1994). "An Old-Fashioned Recipe for Real Time". ACM Transactions on Programming Languages and Systems. 16(5): 1543–1571.
- Agha, G. and K. Palmskog. (2018). "A survey of statistical model checking". *ACM Transactions on Modeling and Computer Simulation* (TOMACS). 28(1): 1–39.
- Ahmadi, M., M. Cubuktepe, N. Jansen, and U. Topcu. (2018). "Verification of Uncertain POMDPs Using Barrier Certificates". In: *Allerton*. IEEE. 115–122.
- Ahmed, A., P. Varakantham, M. Lowalekar, Y. Adulyasak, and P. Jaillet. (2017). "Sampling Based Approaches for Minimizing Regret in Uncertain Markov Decision Processes (MDPs)". *J. Artif. Intell. Res.* 59: 229–264.
- Akash, K., G. McMahon, T. Reid, and N. Jain. (2020). "Human Trust-based Feedback Control: Dynamically Varying Automation Transparency to Optimize Human-Machine Interactions". *IEEE Control Systems Magazine*. 40(6): 98–116.
- Almagor, S., U. Boker, and O. Kupferman. (2016). "Formally Reasoning About Quality". *Journal of the ACM*. 63(3): 24:1–24:56.

Alshiekh, M., R. Bloem, R. Ehlers, B. Könighofer, S. Niekum, and U. Topcu. (2018). "Safe Reinforcement Learning via Shielding". In: Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 32. No. 1.

- Alur, R., A. Kanade, and G. Weiss. (2008). "Ranking Automata and Games for Prioritized Requirements". In: *Proc. International Conference on Computer-Aided Verification*. Vol. 5123. *LNCS*.
- Alur, R. and S. La Torre. (2004). "Deterministic generators and games for LTL fragments". *ACM Transactions on Computational Logic*. 5(1): 1–25.
- Amato, C., D. S. Bernstein, and S. Zilberstein. (2010). "Optimizing Fixed-Size Stochastic Controllers for POMDPs and Decentralized POMDPs". *AAMAS*. 21(3): 293–320.
- Anderson, M. (2020). "Surprise! 2020 Is Not the Year for Self-Driving Cars". IEEE Spectrum. URL: https://spectrum.ieee.org/transportation/self-driving/surprise-2020-is-not-the-year-for-selfdriving-cars.
- Arming, S., E. Bartocci, K. Chatterjee, J.-P. Katoen, and A. Sokolova. (2018). "Parameter-Independent Strategies for pMDPs via POMDPs". In: *QEST*. Springer. 53–70.
- Asarin, E., O. Maler, A. Pnueli, and J. Sifakis. (1998). "Controller synthesis for timed automata". In: *IFAC Symposium on System Structure and Control*. Elsevier. 469–474.
- Bacci, G., M. Hansen, and K. G. Larsen. (2019). "Model Checking Constrained Markov Reward Models with Uncertainties". In: *QEST*. 37–51.
- Badings, T., M. Cubuktepe, N. Jansen, S. Junges, J.-P. Katoen, and U. Topcu. (2022). "Scenario-Based Verification of Uncertain Parametric MDPs". *International Journal on Software Tools for Technology Transfer*. 24(5): 803–819.
- Bai, H., D. Hsu, and W. S. Lee. (2014). "Integrated perception and planning in the continuous space: A POMDP approach". *The International Journal of Robotics Research*. 33(9): 1288–1302.
- Baier, C., M. Größer, and N. Bertrand. (2012). "Probabilistic  $\omega$ -automata". Journal of the ACM. 59(1): 1:1–1:52.

References 205

Baier, C., C. Hensel, L. Hutschenreiter, S. Junges, J.-P. Katoen, and J. Klein. (2020). "Parametric Markov Chains: PCTL Complexity and Fraction-Free Gaussian Elimination". *Inf. Comput.* 272: 104504.

- Baier, C. and J.-P. Katoen. (2008). *Principles of Model Checking*. MIT Press.
- Bakirtzis, G., T. Sherburne, S. C. Adams, B. M. Horowitz, P. A. Beling, and C. H. Fleming. (2022a). "An Ontological Metamodel for Cyber-Physical System Safety, Security, and Resilience Coengineering". Software & Systems Modeling. DOI: 10.1007/s10270-021-00892-z.
- Bakirtzis, G., S. Carr, D. Danks, and U. Topcu. (2022b). "Dynamic Certification for Autonomous Systems". DOI: 10.48550/ARXIV.2203. 10950.
- Bakker, B. (2001). "Reinforcement Learning with Long Short-Term Memory". In: Advances in Neural Information Processing Systems (NIPS). MIT Press. 1475–1482.
- Basic Theory of Driving: The Official Handbook. (2018). 10th ed. Singapore Traffic Police.
- Basu, S., R. Pollack, and M. Roy. (2010). Algorithms in Real Algebraic Geometry. Springer. ISBN: 3540330984.
- Baum, E. B. and F. Wilczek. (1987). "Supervised Learning of Probability Distributions by Neural Networks". In: *Advances in Neural Information Processing Systems (NIPS)*. American Institute of Physics. 52–61.
- Benenson, R., S. Petti, T. Fraichard, and M. Parent. (2006). "Integrating perception and planning for autonomous navigation of urban vehicles". In: 2006 IEEE/RSJ International Conference on Intelligent Robots and Systems. IEEE. 98–104.
- Bérard, B., J. Mullins, and M. Sassolas. (2015). "Quantifying opacity". Mathematical Structures in Computer Science. 25(2): 361–403.
- Bertsekas, D. (2019). Reinforcement learning and optimal control. Athena Scientific.
- Biere, A., M. Heule, and H. van Maaren. (2009). *Handbook of satisfia-bility*. Vol. 185. IOS press.

Bloem, R., K. Chatterjee, T. A. Henzinger, and B. Jobstmann. (2009). "Better Quality in Synthesis through Quantitative Objectives". In: *Proc. International Conference on Computer-Aided Verification*. Vol. 5643. *LNCS*. 140–156.

- Bloem, R., B. Könighofer, R. Könighofer, and C. Wang. (2015). "Shield Synthesis: Runtime Enforcement for Reactive Systems". In: Proceedings of the 21st International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS 2015). Lecture Notes in Computer Science. Springer. DOI: 10.1007/978-3-662-46681-0\ 51.
- Boff, K. R. and J. E. Lincoln. (1988). Engineering data compendium: Human perception and performance. Vol. 2. Harry G. Armstrong Aerospace Medical Research Laboratory.
- Borkar, V. and R. Jain. (2014). "Risk-Constrained Markov Decision Processes". *IEEE Transactions on Automatic Control.* 59(9): 2574–2579.
- Bortolussi, L. and S. Silvetti. (2018). "Bayesian Statistical Parameter Synthesis for Linear Temporal Properties of Stochastic Models". In: *TACAS*. 396–413.
- Bouma, W., W. Levelt, A. Melisse, K. Middelburg, and L. Verhaard. (1994). "Formalization of properties for feature interaction detection: Experience in real-life situation". In: *Towards a Pan-European Telecommunication Servcie Infrastructure*. Ed. by H.-J. Kugler, A. Mullery, and N. Niebert. *Lecture Notes in Computer Science*. No. 851. Springer-Verlag. 393–405.
- Bouton, M., J. Tumova, and M. J. Kochenderfer. (2020). "Point-Based Methods for Model Checking in Partially Observable Markov Decision Processes". In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 34. 10061–10068.
- Boyd, S. and L. Vandenberghe. (2004). *Convex Optimization*. New York, NY, USA: Cambridge University Press. ISBN: 0521833787.
- Bozkurt, A. K., Y. Wang, M. M. Zavlanos, and M. Pajic. (2020). "Control synthesis from linear temporal logic specifications using model-free reinforcement learning". In: 2020 IEEE International Conference on Robotics and Automation (ICRA). IEEE. 10349–10355.

References 207

Büchi, J. R. (1990). "On a decision method in restricted second order arithmetic". In: *The collected works of J. Richard Büchi*. Springer. 425–435.

- Burdick, J. W., N. E. DuToit, A. Howard, C. Looman, J. Ma, R. M. Murray, and T. Wongpiromsarn. (2007). "Sensing, Navigation and Reasoning Technologies for the DARPA Urban Challenge". *Tech. rep.* DARPA Urban Challenge Final Report.
- Burns, B. and O. Brock. (2007). "Sampling-Based Motion Planning with Sensing Uncertainty". In: *ICRA*. IEEE. 3313–3318.
- Calafiore, G. C. and M. C. Campi. (2006). "The Scenario Approach to Robust Control Design". *IEEE Trans. Automat. Contr.* 51(5): 742–753.
- Calinescu, R., K. Johnson, and C. Paterson. (2016). "FACT: A Probabilistic Model Checker for Formal Verification with Confidence Intervals". In: Tools and Algorithms for the Construction and Analysis of Systems (TACAS). Vol. 9636. Lecture Notes in Computer Science. Springer. 540–546.
- Camacho, A., R. T. Icarte, T. Q. Klassen, R. A. Valenzano, and S. A. McIlraith. (2019). "LTL and Beyond: Formal Languages for Reward Function Specification in Reinforcement Learning." In: *IJCAI*. Vol. 19. 6065–6073.
- Campbell, M., M. Egerstedt, J. P. How, and R. M. Murray. (2010). "Autonomous driving in urban environments: approaches, lessons and challenges". *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences.* 368(1928): 4649–4672.
- Campi, M. C. and S. Garatti. (2011). "A Sampling-and-Discarding Approach to Chance-Constrained Optimization: Feasibility and Optimality". *Journal of Optimization Theory and Applications*. 148(2): 257–280.
- Campi, M. C. and S. Garatti. (2008). "The Exact Feasibility of Randomized Solutions of Uncertain Convex Programs". SIAM Journal on Optimization. 19(3): 1211–1230.
- Campi, M. C., S. Garatti, and F. A. Ramponi. (2018). "A General Scenario Theory for Nonconvex Optimization and Decision Making". *IEEE Trans. Automat. Contr.* 63(12): 4067–4078.

Carr, S., N. Jansen, and U. Topcu. (2020). "Verifiable RNN-Based Policies for POMDPs Under Temporal Logic Constraints". In: *International Joint Conference on Artificial Intelligence (IJCAI)*. IJCAI.org. 4121–4127.

- Carr, S., N. Jansen, and U. Topcu. (2021). "Task-Aware Verifiable RNN-Based Policies for Partially Observable Markov Decision Processes". J. Artif. Intell. Res. 72: 819–847.
- Carr, S., N. Jansen, R. Wimmer, A. C. Serban, B. Becker, and U. Topcu. (2019). "Counterexample-Guided Strategy Improvement for POMDPs Using Recurrent Neural Networks". In: *International Joint Conference on Artificial Intelligence (IJCAI)*. IJCAI.org. 5532–5539.
- Cassandra, A., M. L. Littman, and N. L. Zhang. (1997). "Incremental Pruning: A Simple, Fast, Exact Method for Partially Observable Markov Decision Processes". In: *UAI*. 54–61.
- Castro, L. I. R., P. Chaudhari, J. Tumova, S. Karaman, E. Frazzoli, and D. Rus. (2013). "Incremental sampling-based algorithm for minimum-violation motion planning". In: 52nd IEEE Conference on Decision and Control. 3217–3224.
- Censi, A., K. Slutsky, T. Wongpiromsarn, D. Yershov, S. Pendleton, J. Fu, and E. Frazzoli. (2019). "Liability, Ethics, and Culture-Aware Behavior Specification using Rulebooks". In: 2019 International Conference on Robotics and Automation (ICRA). 8536–8542. DOI: 10.1109/ICRA.2019.8794364.
- Cerrito, S. and M. C. Mayer. (1998). "Using Linear Temporal Logic to Model and Solve Planning Problems". In: AIMSA. 141–152.
- Chamie, M. E. and H. Mostafa. (2018). "Robust Action Selection in Partially Observable Markov Decision Processes with Model Uncertainty". In: *CDC*. IEEE. 5586–5591.
- Chao, H., Y. Cao, and Y. Chen. (2010). "Autopilots for small unmanned aerial vehicles: a survey". *International Journal of Control, Automation and Systems*. 8(1): 36–44. DOI: 10.1007/s12555-010-0105-z.
- Chatterjee, K., M. Chmelík, R. Gupta, and A. Kanodia. (2016). "Optimal Cost Almost-sure Reachability in POMDPs". *Artificial Intelligence*. 234: 26–48.

Chen, J. Y. and M. J. Barnes. (2012). "Supervisory control of multiple robots: Effects of imperfect automation and individual differences". *Human Factors*. 54(2): 157–174. DOI: 10.1177/0018720811435843.

- Chen, M., S. Nikolaidis, H. Soh, D. Hsu, and S. Srinivasa. (2020). "Trust-Aware Decision Making for Human-Robot Collaboration: Model Learning and Planning". *ACM Transactions on Human-Robot Interaction (THRI)*. 9(2): 1–23.
- Chen, T., T. Han, and M. Kwiatkowska. (2013a). "On the Complexity of Model Checking Interval-Valued Discrete Time Markov Chains". *Information Processing Letters.* 113(7): 210–216.
- Chen, X., L. Niu, and Y. Yuan. (2013b). "Optimality Conditions and a Smoothing Trust Region Newton Method for Non-Lipschitz Optimization". SIAM Journal on Optimization. 23(3): 1528–1552.
- Cimatti, A., E. M. Clarke, E. Giunchiglia, F. Giunchiglia, M. Pistore, M. Roveri, R. Sebastiani, and A. Tacchella. (2002). "NuSMV 2: An OpenSource Tool for Symbolic Model Checking". In: CAV. Ed. by E. Brinksma and K. G. Larsen. Vol. 2404. Lecture Notes in Computer Science. Springer. 359–364. ISBN: 3-540-43997-8.
- Cimatti, A., M. Roveri, V. Schuppan, and A. Tchaltsev. (2008). "Diagnostic Information for Realizability". In: *Proc. International Conference on Verification, Model Checking, and Abstract Interpretation. LNCS.*
- Cimatti, A., M. Roveri, V. Schuppan, and S. Tonetta. (2007). "Boolean Abstraction for Temporal Logic Satisfiability". In: *Proc. International Conference on Computer-Aided Verification*. Vol. 4590. *LNCS*. 532–546.
- Clarke, E. M., E. A. Emerson, and A. P. Sistla. (1986). "Automatic verification of finite-state concurrent systems using temporal logic specifications". *ACM Transactions on Programming Languages and Systems (TOPLAS)*. 8(2): 244–263. ISSN: 0164-0925. DOI: http://doi.acm.org/10.1145/5397.5399.
- Conner, D., H. Kress-Gazit, H. Choset, A. Rizzi, and G. Pappas. (2007). "Valet parking without a valet". In: *Proc. of IEEE/RSJ International Conference on Intelligent Robots and Systems.* 572–577. DOI: 10.1109/IROS.2007.4399374.

Cooke, N. J. and H. K. Pedersen. (2009). "Unmanned aerial vehicles". Handbook of Aviation Human Factors.

- Cubuktepe, M., N. Jansen, M. Alsiekh, and U. Topcu. (2020a). "Synthesis of Provably Correct Autonomy Protocols for Shared Control". In: *IEEE Transactions on Automatic Control*, Accepted. IEEE.
- Cubuktepe, M., N. Jansen, S. Junges, J.-P. Katoen, I. Papusha, H. A. Poonawala, and U. Topcu. (2017). "Sequential Convex Programming for the Efficient Verification of Parametric MDPs". In: TACAS (2). Vol. 10206. LNCS. 133–150.
- Cubuktepe, M., N. Jansen, S. Junges, J.-P. Katoen, and U. Topcu. (2018). "Synthesis in pMDPs: A Tale of 1001 Parameters". In: *ATVA*. Vol. 11138. *LNCS*. Springer. 160–176.
- Cubuktepe, M., N. Jansen, S. Junges, J.-P. Katoen, and U. Topcu. (2020b). "Scenario-Based Verification of Uncertain MDPs". In: *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. Springer. 287–305.
- Cubuktepe, M., N. Jansen, S. Junges, J.-P. Katoen, and U. Topcu. (2020c). "Scenario-Based Verification of Uncertain MDPs". In: TACAS. Vol. 12078. Lecture Notes in Computer Science. Springer. 287–305.
- Cubuktepe, M., N. Jansen, S. Junges, J.-P. Katoen, and U. Topcu. (2021a). "Convex Optimization for Parameter Synthesis in MDPs". In: *IEEE Transactions on Automatic Control*, Under Revision. IEEE.
- Cubuktepe, M., N. Jansen, S. Junges, A. Marandi, M. Suilen, and U. Topcu. (2021b). "Robust Finite-State Controllers for Uncertain POMDPs". In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI). AAAI Press. 11792–11800.
- Daws, C. (2004). "Symbolic and Parametric Model Checking of Discrete-Time Markov Chains". In: *International Colloquium on Theoretical* Aspects of Computing (ICTAC). Vol. 3407. Lecture Notes in Computer Science. Springer. 280–294. ISBN: 3-540-25304-1.
- Degrave, J., F. Felici, J. Buchli, M. Neunert, B. Tracey, F. Carpanese, T. Ewalds, R. Hafner, A. Abdolmaleki, D. de Las Casas, *et al.* (2022). "Magnetic control of tokamak plasmas through deep reinforcement learning". *Nature*. 602(7897): 414–419.

References 211

Dehnert, C., S. Junges, N. Jansen, F. Corzilius, M. Volk, H. Bruintjes, J.-P. Katoen, and E. Ábrahám. (2015). "PROPhESY: A PRObabilistic ParamEter SYnthesis Tool". In: *CAV* (1). Vol. 9206. *Lecture Notes in Computer Science*. Springer. 214–231.

- Dehnert, C., S. Junges, J.-P. Katoen, and M. Volk. (2017). "A Storm is Coming: A Modern Probabilistic Model Checker". In: *CAV* (2). Vol. 10427. *LNCS*. Springer. 592–600.
- Dimitrova, R., M. Ghasemi, and U. Topcu. (2018). "Maximum Realizability for Linear Temporal Logic Specifications". In: *Proc. Automated Technology for Verification and Analysis*. Springer. 458–475.
- Ding, X. C. D., S. L. Smith, C. Belta, and D. Rus. (2011). "LTL control in uncertain environments with probabilistic satisfaction guarantees". *IFAC Proceedings Volumes*. 44(1): 3515–3520.
- Djeumou, F., M. Cubuktepe, C. Lennon, and U. Topcu. (2021). "Task-Guided Inverse Reinforcement Learning under Partial Information". In: *International Conference on Neural Information Processing Systems*, Submitted.
- Donmez, B., C. E. Nehme, and M. L. Cummings. (2010). "Modeling Workload Impact in Multiple Unmanned Vehicle Supervisory Control". *IEEE Trans. Syst. Man Cybern. Part A.* 40(6): 1180–1190. DOI: 10.1109/TSMCA.2010.2046731.
- Donzé, A. and O. Maler. (2010). "Robust Satisfaction of Temporal Logic over Real-Valued Signals". In: Formal Modeling and Analysis of Timed Systems. Ed. by K. Chatterjee and T. A. Henzinger. Berlin, Heidelberg: Springer Berlin Heidelberg. 92–106. ISBN: 978-3-642-15297-9.
- Dosovitskiy, A., G. Ros, F. Codevilla, A. M. López, and V. Koltun. (2017). "CARLA: An Open Urban Driving Simulator". In: *CoRL*. Vol. 78. *Proceedings of Machine Learning Research*. PMLR. 1–16.
- DuToit, N. E., T. Wongpiromsarn, J. W. Burdick, and R. M. Murray. (2008). "Situational Reasoning for Road Driving in an Urban Environment". In: *International Workshop on Intelligent Vehicle Control Systems (IVCS)*.

Ehlers, R. and V. Raman. (2014). "Low-Effort Specification Debugging and Analysis". In: *Proc. Workshop on Synthesis*. Vol. 157. *EPTCS*. 117–133.

- Emerson, E. A. (1990). "Temporal and modal logic". In: *Handbook of theoretical computer science (vol. B): formal models and semantics*. Cambridge, MA: MIT Press. 995–1072. ISBN: 0-444-88074-7.
- Eysenbach, B., R. R. Salakhutdinov, and S. Levine. (2021). "Robust predictable control". *Advances in Neural Information Processing Systems*. 34.
- Faella, M. (2009). "Admissible Strategies in Infinite Games over Graphs".
  In: Mathematical Foundations of Computer Science 2009, 34th International Symposium, MFCS 2009, Novy Smokovec, High Tatras, Slovakia, August 24-28, 2009. Proceedings. Vol. 5734. Lecture Notes in Computer Science. Springer. 307–318. DOI: 10.1007/978-3-642-03816-7√27.
- Farokhi, F. and H. Sandberg. (2017). "Fisher information as a measure of privacy: Preserving privacy of households with smart meters using batteries". *IEEE Transactions on Smart Grid.* 9(5): 4726–4734.
- Faymonville, P., B. Finkbeiner, M. N. Rabe, and L. Tentrup. (2017). "Encodings of Bounded Synthesis". In: *Proc. International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. Vol. 10205. *LNCS*. 354–370.
- Feng, L., M. Ghasemi, K.-W. Chang, and U. Topcu. (2018). "Counterexamples for robotic planning explained in structured language". In: 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE. 7292–7297.
- Feng, L., L. Humphrey, I. Lee, and U. Topcu. (2016a). "Human-interpretable diagnostic information for robotic planning systems".
  In: Intelligent Robots and Systems (IROS), 2016 IEEE/RSJ International Conference on. IEEE. 1673–1680.
- Feng, L., C. Wiltsche, L. Humphrey, and U. Topcu. (2015). "Controller synthesis for autonomous systems interacting with human operators". In: *Proceedings of the acm/ieee sixth international conference on cyber-physical systems.* 70–79.

Feng, L., C. Wiltsche, L. R. Humphrey, and U. Topcu. (2016b). "Synthesis of Human-in-the-Loop Control Protocols for Autonomous Systems". *IEEE Trans Autom. Sci. Eng.* 13(2): 450–462. DOI: 10.1109/TASE.2016.2530623.

- Filieri, A., G. Tamburrelli, and C. Ghezzi. (2016). "Supporting Self-Adaptation via Quantitative Verification and Sensitivity Analysis at Run Time". *IEEE Trans. Software Eng.* 42(1): 75–99.
- Final Theory of Driving: The Official Handbook. (2017). 9th ed. Singapore Traffic Police.
- Finkbeiner, B. and S. Schewe. (2013). "Bounded synthesis". *International Journal on Software Tools for Technology Transfer*. 15(5-6).
- Forejt, V., M. Kwiatkowska, G. Norman, D. Parker, and H. Qu. (2011). "Quantitative multi-objective verification for probabilistic systems". In: International Conference on Tools and Algorithms for the Construction and Analysis of Systems. Springer. 112–127.
- Franceschetti, M. and P. Minero. (2014). "Elements of information theory for networked control systems". In: *Information and Control in Networks*. Springer. 3–37.
- Frey, G. R., C. D. Petersen, F. A. Leve, I. V. Kolmanovsky, and A. R. Girard. (2017). "Constrained Spacecraft Relative Motion Planning Exploiting Periodic Natural Motion Trajectories and Invariance". Journal of Guidance, Control, and Dynamics. 40(12): 3100–3115.
- Fu, J., N. Atanasov, U. Topcu, and G. J. Pappas. (2016). "Optimal temporal logic planning in probabilistic semantic maps". In: 2016 IEEE International Conference on Robotics and Automation (ICRA). IEEE. 3690–3697.
- Fu, J. and U. Topcu. (2015). "Integrating active sensing into reactive synthesis with temporal logic constraints under partial observations". In: ACC. IEEE. 2408–2413.
- Fu, J. and U. Topcu. (2016). "Synthesis of joint control and active sensing strategies under temporal logic constraints". *IEEE Transactions on Automatic Control*. 61(11): 3464–3476.
- Gabbay, D. M., C. J. Hogger, and J. A. Robinson. (1995). Handbook of Logic in Artificial Intelligence and Logic Programming (Vol. 4): Epistemic and Temporal Reasoning. Oxford, UK: Oxford University Press. ISBN: 0-19-853791-3.

Galton, A., ed. (1987). Temporal Logics and Their Applications. San Diego, CA: Academic Press Professional, Inc. ISBN: 0-12-274060-2.

- Ghasemi, M., E. Bulgur, and U. Topcu. (2020). "Task-oriented active perception and planning in environments with partially known semantics". In: *International Conference on Machine Learning*. PMLR. 3484–3493.
- Ghasemi, M. and U. Topcu. (2019a). "Online Active Perception for Partially Observable Markov Decision Processes with Limited Budget". In: 2019 IEEE 58th Conference on Decision and Control (CDC). IEEE. 6169–6174.
- Ghasemi, M. and U. Topcu. (2019b). "Perception-Aware Point-Based Value Iteration for Partially Observable Markov Decision Processes". In: *Proceedings of the International Joint Conference on Artificial Inteligence (IJCAI)*. 2371–2377.
- Girard, A. and G. J. Pappas. (2009). "Hierarchical control system design using approximate simulation". *Automatica*. 45(2): 566–571. ISSN: 0005-1098. DOI: http://dx.doi.org/10.1016/j.automatica.2008.09.016.
- Givan, R., S. Leach, and T. Dean. (2000). "Bounded-Parameter Markov Decision Processes". *Artificial Intelligence*. 122(1-2): 71–109.
- Gluck, P. and G. Holzmann. (2002). "Using SPIN model checking for flight software verification". In: *Proc. of IEEE Aerospace Conference*. Vol. 1. 1–105–1–113.
- Gohari, P., M. Hale, and U. Topcu. (2020). "Privacy-preserving policy synthesis in Markov decision processes". In: 2020 59th IEEE Conference on Decision and Control (CDC). IEEE. 6266–6271.
- Goodfellow, I. J., Y. Bengio, and A. C. Courville. (2016). *Deep Learning.*Adaptive computation and machine learning. MIT Press.
- Goodfellow, I. J. and O. Vinyals. (2015). "Qualitatively characterizing neural network optimization problems". In: *International Conference on Learning Representations (ICLR)*.
- Goodwin, G. C., M. M. Seron, and J. A. D. Doná. (2004). Constrained Control and Estimation: An Optimisation Approach. Springer.
- Gradel, E. and W. Thomas. (2002). "Automata, logics, and infinite games: a guide to current research".

Gunter, E. L. and D. Peled. (2002). "Temporal Debugging for Concurrent Systems". In: Proceedings of the Eighth International Conference on Tools and Algorithms for the Construction and Analysis of Systems. Vol. 2280. Lecture Notes in Computer Science. Springer International Publishing. 431–444. ISBN: 3-540-43419-4. DOI: 10.1007/3-540-46002-0 30.

- Guo, M., K. H. Johansson, and D. V. Dimarogonas. (2013). "Revising motion planning under linear temporal logic specifications in partially known workspaces". In: 2013 IEEE International Conference on Robotics and Automation. IEEE. 5025–5032.
- Hahn, E. M., V. Hashemi, H. Hermanns, M. Lahijanian, and A. Turrini. (2017). "Multi-Objective Robust Strategy Synthesis for Interval Markov Decision Processes". In: QEST. Vol. 10503. Lecture Notes in Computer Science. Springer. 207–223.
- Hahn, E. M., H. Hermanns, and L. Zhang. (2010). "Probabilistic reachability for parametric Markov models". Software Tools for Technology Transfer. 13(1): 3–19.
- Hahn, E. M., Y. Li, S. Schewe, A. Turrini, and L. Zhang. (2014). "iscasMc: A Web-Based Probabilistic Model Checker". In: *FM*. Vol. 8442. *Lecture Notes in Computer Science*. Springer. 312–317.
- Hahn, E. M., M. Perez, S. Schewe, F. Somenzi, A. Trivedi, and D. Wojtczak. (2019). "Omega-regular objectives in model-free reinforcement learning". In: *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. Springer. 395–412.
- Han, T., J.-P. Katoen, and D. Berteun. (2009). "Counterexample generation in probabilistic model checking". *IEEE Transactions on Software Engineering*. 35(2): 241–257.
- Harel, D., A. Marron, and J. Sifakis. (2022). "Creating a Foundation for Next-Generation Autonomous Systems". *IEEE Design & Test.* 39(1): 49–56. DOI: 10.1109/MDAT.2021.3069959.
- Hasanbeig, M., Y. Kantaros, A. Abate, D. Kroening, G. J. Pappas, and I. Lee. (2019a). "Reinforcement Learning for Temporal Logic Control Synthesis with Probabilistic Satisfaction Guarantees". In: 2019 IEEE 58th Conference on Decision and Control (CDC). 5338–5343. DOI: 10.1109/CDC40024.2019.9028919.

Hasanbeig, M., A. Abate, and D. Kroening. (2019b). "Logically-Constrained Neural Fitted Q-Iteration". In: Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems. AAMAS '19. Montreal QC, Canada: International Foundation for Autonomous Agents and Multiagent Systems. 2012–2014. ISBN: 9781450363099.

- Hausknecht, M. J. and P. Stone. (2015). "Deep Recurrent Q-Learning for Partially Observable MDPs". In: *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*. AAAI Press. 29–37.
- Havelund, K., M. Lowry, and J. Penix. (2001). "Formal analysis of a space-craft controller using SPIN". *IEEE Transactions on Software Engineering*. 27(8): 749–765. DOI: 10.1109/32.940728.
- Hayes, B. and J. A. Shah. (2017). "Improving Robot Controller Transparency Through Autonomous Policy Explanation". In: *Proceedings of the 2017 ACM/IEEE International Conference on Human-Robot Interaction*. ACM. 303–312.
- Heess, N., J. J. Hunt, T. P. Lillicrap, and D. Silver. (2015a). "Memory-based control with recurrent neural networks". *CoRR*. 1512.04455.
- Heess, N., G. Wayne, D. Silver, T. Lillicrap, T. Erez, and Y. Tassa. (2015b). "Learning continuous control policies by stochastic value gradients". In: *Advances in Neural Information Processing Systems* (NIPS). MIT Press. 2944–2952.
- Ho, C. P. and M. Petrik. (2018). "Fast Bellman Updates for Robust MDPs". In: *ICML*.
- Hobbs, K. L. and E. M. Feron. (2020). "A Taxonomy for Aerospace Collision Avoidance with Implications for Automation in Space Traffic Management". In: AIAA Scitech 2020 Forum. 0877.
- Hochreiter, S. and J. Schmidhuber. (1997). "Long short-term memory". Neural Computation. 9(8): 1735–1780.
- Hoeffding, W. (1994). "Probability inequalities for sums of bounded random variables". In: *The collected works of Wassily Hoeffding*. Springer. 409–426.
- Holzmann, G. (1994). "The Theory and Practice of A Formal Method: NewCoRe". In: Proc. of the IFIP World Computer Congress. Vol. 1. North-Holland Publ. 35–44.
- Holzmann, G. J. (2004). The Spin Model Checker. Addison-Wesley.

Holzmann, G. J. (2014). "Mars Code". Commun. ACM. 57(2): 64–73. ISSN: 0001-0782. DOI: 10.1145/2560217.2560218. URL: https://doi.org/10.1145/2560217.2560218.

- Hosseini, H., J. Hoey, and R. Cohen. (2014). "A Coordinated MDP Approach to Multi-Agent Planning for Resource Allocation, with Applications to Healthcare". arXiv preprint arXiv:1407.1584.
- Humphrey, L. and M. Patzek. (2013). "Model Checking Human-Automation UAV Mission Plans". In: *Proc. AIAA Guidance, Navigation, and Control Conf.*
- Humphrey, L. R., E. M. Wolff, and U. Topcu. (2014). "Formal specification and synthesis of mission plans for unmanned aerial vehicles". In: *AAAI Spring Symposium Series*.
- Humphrey, L. R., B. Könighofer, R. Könighofer, and U. Topcu. (2016). "Synthesis of Admissible Shields". In: Proceedings of the 12th International Conference on Hardware and Software: Verification and Testing (HVC 2016). Vol. 10028. Lecture Notes in Computer Science. 134–151. DOI: 10.1007/978-3-319-49052-6\\_9.
- Huth, M. and M. Ryan. (2004). Logic in Computer Science: Modelling and Reasoning about Systems. 2nd. Cambridge University Press.
- Icarte, R. T., T. Klassen, R. Valenzano, and S. McIlraith. (2018). "Using reward machines for high-level task specification and decomposition in reinforcement learning". In: *International Conference on Machine Learning*. 2107–2116.
- Icarte, R. T., E. Waldie, T. Klassen, R. Valenzano, M. Castro, and S. McIlraith. (2019). "Learning reward machines for partially observable reinforcement learning". In: *Proceedings of Advances in Neural Information Processing Systems*. 15523–15534.
- Itoh, H. and K. Nakamura. (2007). "Partially Observable Markov Decision Processes with Imprecise Parameters". *Artificial Intelligence*. 171(8): 453–490.
- Jadbabaie, A. (2000). "Nonlinear Receding Horizon Control: A Control Lyapunov Function Approach". *PhD thesis*. California Institute of Technology.

Jagadeesan, L. J., C. Puchol, and J. E. Von Olnhausen. (1996). "A formal approach to reactive systems software: a telecommunications application in ESTEREL". Formal Methods in System Design. 8(2): 123–151. ISSN: 0925-9856. DOI: http://dx.doi.org/10.1007/BF00122418.

- Jansen, N., B. Könighofer, S. Junges, A. Serban, and R. Bloem. (2020). "Safe Reinforcement Learning Using Probabilistic Shields". In: 31st International Conference on Concurrency Theory (CONCUR 2020). Schloss Dagstuhl-Leibniz-Zentrum für Informatik.
- Jiang, S. and R. Kumar. (2001). "Failure Diagnosis of Discrete Event Systems with Linear-time Temporal Logic Fault Specifications". In: *IEEE Transactions on Automatic Control.* 128–133.
- Jones, A., M. Schwager, and C. Belta. (2013). "Distribution temporal logic: Combining correctness with quality of estimation". In: 52nd IEEE Conference on Decision and Control. IEEE. 4719–4724.
- Jothimurugan, K., S. Bansal, O. Bastani, and R. Alur. (2021). "Compositional reinforcement learning from logical specifications". *Advances in Neural Information Processing Systems*. 34.
- Julian, K. D., M. J. Kochenderfer, and M. P. Owen. (2019). "Deep Neural Network Compression for Aircraft Collision Avoidance Systems". *Journal of Guidance, Control, and Dynamics.* 42(3): 598–608.
- Junges, S. (2020). "Parameter Synthesis in Markov Models". *PhD thesis*. RWTH Aachen University.
- Junges, S., E. Ábrahám, C. Hensel, N. Jansen, J. Katoen, T. Quatmann, and M. Volk. (2019). "Parameter Synthesis for Markov Models". CoRR. abs/1903.07993.
- Junges, S., N. Jansen, R. Wimmer, T. Quatmann, L. Winterer, J.-P. Katoen, and B. Becker. (2018). "Finite-state controllers of POMDPs via parameter synthesis". In: Conference on Uncertainty in Artificial Intelligence (UAI). AUAI Press. 519–529.
- Karabag, M. O., C. Neary, and U. Topcu. (2022). "Planning Not to Talk: Multiagent Systems that are Robust to Communication Loss". In: Autonomous Agents and Multi-Agent Systems.
- Karabag, M. O., M. Ornik, and U. Topcu. (2019). "Least inferable policies for Markov decision processes". In: 2019 American Control Conference (ACC). IEEE. 1224–1231.

Karabag, M. O., M. Ornik, and U. Topcu. (2021). "Deception in Supervisory Control". *IEEE Transactions on Automatic Control*.

- Karaman, S. and E. Frazzoli. (2009). "Sampling-based Motion Planning with Deterministic  $\mu$ -Calculus Specifications". In: *Proc. of the IEEE Conference on Decision and Control (CDC)*.
- Kearns, M. J., Y. Mansour, and A. Y. Ng. (1999). "Approximate Planning in Large POMDPs via Reusable Trajectories". In: *Advances in Neural Information Processing Systems (NIPS)*. MIT Press. 1001–1007.
- Ketkar, N. (2017). "Introduction to keras". In: Deep learning with Python. Springer. 97–111.
- Kim, K., G. E. Fainekos, and S. Sankaranarayanan. (2015). "On the minimal revision problem of specification automata". *International Journal of Robotics Research*. 34(12).
- Kim, S. C., S. W. Shepperd, H. L. Norris, H. R. Goldberg, and M. S. Wallace. (2007). "Mission Design and Trajectory Analysis for Inspection of a Host Spacecraft by a Microsatellite". In: 2007 IEEE Aerospace Conference. IEEE. 1–23.
- Kingma, D. P. and J. Ba. (2015). "Adam: A Method for Stochastic Optimization". In: *International Conference on Learning Representations* (ICLR).
- Klein, J. and C. Baier. (2007). "On-the-Fly Stuttering in the Construction of Deterministic  $\omega$ -Automata". In: *Implementation and Application of Automata*. Ed. by J. Holub and J. Žďárek. Berlin, Heidelberg: Springer Berlin Heidelberg. 51–61.
- Kloetzer, M. and C. Belta. (2008a). "A Fully Automated Framework for Control of Linear Systems from Temporal Logic Specifications". *IEEE Transactions on Automatic Control.* 53(1): 287–297.
- Kloetzer, M. and C. Belta. (2008b). "Dealing with Nondeterminism in Symbolic Control". In: *Proceedings of the 11th international workshop on Hybrid Systems: Computation and Control.* Berlin, Heidelberg: Springer-Verlag. 287–300.
- Könighofer, B., M. Alshiekh, R. Bloem, L. R. Humphrey, R. Könighofer, U. Topcu, and C. Wang. (2017). "Shield synthesis". Formal Methods Syst. Des. DOI: 10.1007/s10703-017-0276-9.

Koul, A., A. Fern, and S. Greydanus. (2019). "Learning Finite State Representations of Recurrent Policy Networks". In: *International Conference on Learning Representations (ICLR)*.

- Kozen, D. (1983). "Results on the propositional  $\mu$ -calculus". Theoretical Computer Science. 27(3): 333–354.
- Kress-Gazit, H., G. Fainekos, and G. Pappas. (2007). "Where's Waldo? Sensor-Based Temporal Logic Motion Planning". In: *Proc. of IEEE International Conference on Robotics and Automation*. 3116–3121. DOI: 10.1109/ROBOT.2007.363946.
- Kress-Gazit, H., G. E. Fainekos, and G. J. Pappas. (2009). "Temporal-logic-based reactive mission and motion planning". *IEEE transactions on robotics*. 25(6): 1370–1381.
- Kress-Gazit, H., T. Wongpiromsarn, and U. Topcu. (2011). "Correct, Reactive, High-Level Robot Control". *IEEE Robotics Automation Magazine*. 18(3): 65–74.
- Kupferman, O. and M. Y. Vardi. (2001). "Model checking of safety properties". Formal Methods in System Design. 19(3): 291–314.
- Kupferman, O. and M. Y. Vardi. (2005). "Safraless Decision Procedures". In: *Proc. IEEE Annual Symposium on Foundations of Computer Science*. 531–542.
- Kwiatkowska, M., G. Norman, and D. Parker. (2011). "PRISM 4.0: Verification of Probabilistic Real-Time Systems". In: Computer Aided Verification (CAV). Vol. 6806. Lecture Notes in Computer Science. Springer. 585–591.
- Lahijanian, M., S. Almagor, D. Fried, L. E. Kavraki, and M. Y. Vardi. (2015). "This Time the Robot Settles for a Cost: A Quantitative Approach to Temporal Logic Planning with Partial Satisfaction". In: *Proc. Association for the Advancement of Artificial Intelligence*.
- Lahijanian, M. and M. Z. Kwiatkowska. (2016). "Specification revision for Markov decision processes with optimal trade-off". In: *Proc. IEEE Conference on Decision and Control.* 7411–7418.
- Lahijanian, M., M. R. Maly, D. Fried, L. E. Kavraki, H. Kress-Gazit, and M. Y. Vardi. (2016). "Iterative temporal planning in uncertain environments with partial satisfaction guarantees". *IEEE Transactions* on Robotics. 32(3): 583–599.

Lamport, L. (2002). Specifying systems: The TLA+ language and tools for hardware and software engineers. Addison-Wesley. ISBN: 0-3211-4306-X.

- Lamport, L. (1983). "Specifying concurrent program modules". ACM Transactions on Programming Languages and Systems. 5(2): 190–222. ISSN: 0164-0925. DOI: http://doi.acm.org/10.1145/69624.357207.
- Lamport, L. (1994). "The Temporal Logic of Actions". ACM Transactions on Programming Languages and Systems. 16(3): 872–923.
- Lanotte, R., A. Maggiolo-Schettini, and A. Troina. (2007). "Parametric Probabilistic Transition Systems for System Design and Analysis". Formal Aspects Comput. 19(1): 93–109.
- Lecarpentier, E. and E. Rachelson. (2019). "Non-Stationary Markov Decision Processes, a Worst-Case Approach using Model-Based Reinforcement Learning, Extended Version". arXiv preprint arXiv:1904.10090.
- Leike, J., M. Martic, V. Krakovna, P. A. Ortega, T. Everitt, A. Lefrancq, L. Orseau, and S. Legg. (2017). "AI Safety Gridworlds". *CoRR*. abs/1711.09883.
- Leucker, M. and C. Schallhart. (2009). "A brief account of runtime verification". J. Log. Algebraic Methods Program. 78(5): 293–303. DOI: 10.1016/j.jlap.2008.08.004.
- Lignos, C., V. Raman, C. Finucane, M. Marcus, and H. Kress-Gazit. (2015). "Provably correct reactive control from natural language". *Autonomous Robots.* 38(1): 89–105.
- Lin, F. (1993). "Analysis and synthesis of discrete event systems using temporal logic". Control Theory and Advanced Technologies. 9(1): 341–350.
- Liu, X., Q. Zhou, H. Ren, and C. Sun. (2018). "Reinforcement Learning for Robot Navigation in Nondeterministic Environments". In: 2018 5th IEEE International Conference on Cloud Computing and Intelligence Systems (CCIS). IEEE. 615–619.
- Livingston, S. C., R. M. Murray, and J. W. Burdick. (2012). "Backtracking temporal logic synthesis for uncertain environments". In: 2012 IEEE International Conference on Robotics and Automation. IEEE. 5163–5170.

Llerena, Y. R. S., M. Böhme, M. Brünink, G. Su, and D. S. Rosenblum. (2018). "Verifying the Long-run Behavior of Probabilistic System Models in the Presence of Uncertainty". In: ESEC/SIGSOFT FSE. ACM. 587–597.

- Madani, O., S. Hanks, and A. Condon. (1999). "On the Undecidability of Probabilistic Planning and Infinite-Horizon Partially Observable Markov Decision Problems". In: *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*. AAAI Press. 541–548.
- Manna, Z. and A. Pnueli. (1992). The Temporal Logic of Reactive and Concurrent Systems. Springer-Verlag. ISBN: 0-387-97664-7.
- Mao, Y., M. Szmuk, X. Xu, and B. Acikmese. (2018). "Successive Convexification: A Superlinearly Convergent Algorithm for Non-convex Optimal Control Problems". arXiv preprint arXiv:1804.06539.
- Mason, G. R., R. C. Calinescu, D. Kudenko, and A. Banks. (2017). "Assured Reinforcement Learning with Formally Verified Abstract Policies". In: 9th International Conference on Agents and Artificial Intelligence (ICAART). York.
- Mayne, D., J. Rawlings, C. Rao, and P. Scokaert. (2000). "Constrained model predictive control: Stability and optimality". *Automatica*. 36: 789–814.
- Mazala, R. (2001). "Infinite Games". In: Automata, Logics, and Infinite Games: A Guide to Current Research [outcome of a Dagstuhl seminar, February 2001]. Vol. 2500. Lecture Notes in Computer Science. Springer. 23–42. DOI: 10.1007/3-540-36387-4\\_2.
- McCulloch, W. S. and W. Pitts. (1943). "A logical calculus of the ideas immanent in nervous activity". The bulletin of mathematical biophysics. 5(4): 115–133.
- McMillan, K. L. (1993). Symbolic Model Checking. Kluwer Academic Publishers.
- Mehdipour, N., C.-I. Vasile, and C. Belta. (2019). "Average-based Robustness for Continuous-Time Signal Temporal Logic". In: 2019 IEEE 58th Conference on Decision and Control (CDC). 5312–5317. DOI: 10.1109/CDC40024.2019.9029989.
- Meuleau, N., K.-E. Kim, L. P. Kaelbling, and A. R. Cassandra. (1999). "Solving POMDPs by Searching the Space of Finite Policies". In: *UAI*. 417–426.

Michaud, T. and A. Duret-Lutz. (2015). "Practical Stutter-Invariance Checks for ω-Regular Languages". In: Proceedings of the 22Nd International Symposium on Model Checking Software - Volume 9232. SPIN 2015. Stellenbosch, South Africa: Springer-Verlag. 84–101. ISBN: 978-3-319-23403-8.

- Mireshghallah, F., M. Taram, P. Vepakomma, A. Singh, R. Raskar, and H. Esmaeilzadeh. (2020). "Privacy in deep learning: A survey". arXiv preprint arXiv:2004.12254.
- Montana, F. J., J. Liu, and T. J. Dodd. (2017). "Sampling-based reactive motion planning with temporal logic constraints and imperfect state information". In: *Critical Systems: Formal Methods and Automated Verification*. Springer. 134–149.
- Mulder, W. D., S. Bethard, and M. Moens. (2015). "A survey on the application of recurrent neural networks to statistical language modeling". Computer Speech & Language. 30(1): 61–98.
- Murray, R. M., J. Hauser, A. Jadbabaie, M. B. Milam, N. Petit, W. B. Dunbar, and R. Franz. (2003). "Online Control Customization via Optimization-Based Control". In: Software-Enabled Control: Information Technology for Dynamical Systems. Ed. by G. B. Tariq Samad. Wiley-IEEE Press.
- Nair, G. N. and R. J. Evans. (2004). "Stabilizability of stochastic linear systems with finite feedback data rates". SIAM Journal on Control and Optimization. 43(2): 413–436.
- Neary, C., C. Verginis, M. Cubuktepe, and U. Topcu. (2022). "Verifiable and Compositional Reinforcement Learning Systems". In: *Proceedings of the International Conference on Automated Planning and Scheduling*. Vol. 32.
- Neary, C., Z. Xu, B. Wu, and U. Topcu. (2021). "Reward Machines for Cooperative Multi-Agent Reinforcement Learning". In: *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems. AAMAS '21.* 934–942.
- Newcombe, C., T. Rath, F. Zhang, B. Munteanu, M. Brooker, and M. Deardeuff. (2015). "How Amazon web services uses formal methods". *Communications of the ACM*. DOI: 10.1145/2699417.

Nilim, A. and L. El Ghaoui. (2005). "Robust Control of Markov Decision Processes with Uncertain Transition Matrices". *Operations Research*. 53(5): 780–798.

- Norman, G., D. Parker, and X. Zou. (2017). "Verification and Control of Partially Observable Probabilistic Systems". *Real-Time Systems*. 53(3): 354–402.
- O'Donoghue, B., E. Chu, N. Parikh, and S. Boyd. (2016). "Conic Optimization via Operator Splitting and Homogeneous Self-Dual Embedding". *Journal of Optimization Theory and Applications*. 169(3): 1042–1068.
- Omidshafiei, S., A.-.-A. Agha–Mohammadi, C. Amato, S.-.-Y. Liu, J. P. How, and J. Vian. (2017). "Decentralized Control of Multi-Robot Partially Observable Markov Decision Processes Using Belief Space Macro-Actions". The International Journal of Robotics Research. 36(2): 231–258.
- Osogami, T. (2015). "Robust Partially Observable Markov Decision Process". In: *ICML*. Vol. 37, 106–115.
- Paden, B., M. Čáp, S. Z. Yong, D. Yershov, and E. Frazzoli. (2016). "A Survey of Motion Planning and Control Techniques for Self-Driving Urban Vehicles". *IEEE Transactions on Intelligent Vehicles*. 1(1): 33–55. DOI: 10.1109/TIV.2016.2578706.
- Pascanu, R., C. Gulcehre, K. Cho, and Y. Bengio. (2014). "How to Construct Deep Recurrent Neural Networks". In: *International Conference on Learning Representations (ICLR)*.
- Peled, D. and T. Wilke. (1997). "Stutter-invariant temporal properties are expressible without the next-time operator". *Inf. Process. Lett.* 63(5): 243–246. ISSN: 0020-0190. DOI: http://dx.doi.org/10.1016/S0020-0190(97)00133-6.
- Phan-Minh, T., K. X. Cai, and R. M. Murray. (2019). "Towards Assume-Guarantee Profiles for Autonomous Vehicles". In: 2019 IEEE 58th IEEE Conference on Decision and Control (CDC).
- Pineau, J., G. Gordon, and S. Thrun. (2003). "Point-based value iteration: An anytime algorithm for POMDPs". In: *International Joint Conference on Artificial Intelligence (IJCAI)*. 1025–1032.

Piterman, N., A. Pnueli, and Y. Sa'ar. (2006). "Synthesis of Reactive(1) Designs". In: Verification, Model Checking and Abstract Interpretation. Vol. 3855. Lecture Notes in Computer Science. Springer-Verlag. 364–380. DOI: DOI:10.1007/11609773.

- Pnueli, A. (1986). "Applications of temporal logic to the specification and verification of reactive systems: a survey of current trends". Current Trends in Concurrency. Overviews and Tutorials: 510–584.
- Pnueli, A. and R. Rosner. (1989). "On the synthesis of a reactive module". In: *Proceedings of the 16th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM. 179–190.
- Pnueli, A. (1977). "The Temporal Logic of Programs". In: Annual Symposium on Foundations of Computer Science (FOCS). IEEE Computer Society. 46–57.
- Poupart, P. and C. Boutilier. (2003). "Bounded Finite State Controllers". In: Advances in Neural Information Processing Systems (NIPS). MIT Press. 823–830.
- Powell, W. B. (2022). Reinforcement Learning and Stochastic Optimization: A unified framework for sequential decisions. John Wiley & Sons.
- Puggelli, A., W. Li, A. L. Sangiovanni-Vincentelli, and S. A. Seshia. (2013). "Polynomial-Time Verification of PCTL Properties of MDPs with Convex Uncertainties". In: Computer Aided Verification (CAV). Vol. 8044. Lecture Notes in Computer Science. Springer. 527–542.
- Puterman, M. L. (2014). Markov Decision Processes: Discrete Stochastic Dynamic Programming. John Wiley & Sons.
- Quatmann, T., C. Dehnert, N. Jansen, S. Junges, and J.-P. Katoen. (2016). "Parameter Synthesis for Markov Models: Faster than Ever". In: *ATVA*. Vol. 9938. *Lecture Notes in Computer Science*. Springer. 50–67.
- Rabin, M. O. and D. Scott. (1959). "Finite automata and their decision problems". *IBM J. Res. Dev.* DOI: 10.1147/rd.32.0114.

Radoszycki, J., N. Peyrard, and R. Sabbadin. (2015). "Solving  $F^3$  MDPs: Collaborative Multiagent Markov Decision Processes with Factored Transitions, Rewards and Stochastic Policies". In: *International Conference on Principles and Practice of Multi-Agent Systems*. Springer. 3–19.

- Raman, V. and H. Kress-Gazit. (2013). "Towards minimal explanations of unsynthesizability for high-level robot behaviors". In: *Proc. IEEE/RSJ International Conference on Intelligent Robots and Systems*. 757–762.
- Saboori, A. and C. N. Hadjicostis. (2007). "Notions of security and opacity in discrete event systems". In: 2007 46th IEEE Conference on Decision and Control. IEEE. 5056–5061.
- Scheftelowitsch, D., P. Buchholz, V. Hashemi, and H. Hermanns. (2017). "Multi-Objective Approaches to Markov Decision Processes with Uncertain Transition Parameters". In: *VALUETOOLS*. 44–51.
- Schewe, S. and B. Finkbeiner. (2007). "Bounded Synthesis". In: *Proc. Automated Technology for Verification and Analysis*. Vol. 4762. *LNCS*. 474–488.
- Schneider, F., S. Easterbrook, J. Callahan, and G. Holzmann. (1998). "Validating requirements for fault tolerant systems using model checking". In: *Proceedings of IEEE International Symposium on Requirements Engineering: RE '98.* 4–13. DOI: 10.1109/ICRE.1998.667803.
- Schuppan, V. (2012). "Towards a notion of unsatisfiable and unrealizable cores for LTL". Science of Computer Programming. 77(7-8): 908–939.
- Sen, K., M. Viswanathan, and G. Agha. (2006). "Model-Checking Markov Chains in the Presence of Uncertainties". In: *TACAS*. Vol. 3920. *Lecture Notes in Computer Science*. Springer. 394–410.
- Seow, K. T. and R. Devanathan. (1994). "A temporal framework for assembly sequence representation and analysis". *IEEE Transactions on Robotics and Automation*. 10(2): 220–229.
- Shah, S., D. Dey, C. Lovett, and A. Kapoor. (2017). "AirSim: High-Fidelity Visual and Physical Simulation for Autonomous Vehicles". In: *Field and Service Robotics*. eprint: arXiv:1705.05065. URL: https://arxiv.org/abs/1705.05065.

Shalev-Shwartz, S., S. Shammah, and A. Shashua. (2017). "On a Formal Model of Safe and Scalable Self-driving Cars". *ArXiv*. abs/1708.06374.

- Sherstinsky, A. (2020). "Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network". *Physica D: Nonlinear Phenomena*. 404: 132306.
- Silva, R. R. da, V. Kurtz, and H. Lin. (2019). "Active Perception and Control From Temporal Logic Specifications". *IEEE Control Systems Letters*. 3(4): 1068–1073.
- Silver, D., T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, et al. (2018). "A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play". Science. 362(6419): 1140–1144.
- Spel, J., S. Junges, and J. Katoen. (2019). "Are Parametric Markov Chains Monotonic?" In: ATVA. Vol. 11781. Lecture Notes in Computer Science. Springer. 479–496.
- Srinivasan, M., S. Coogan, and M. Egerstedt. (2018). "Control of Multi-Agent Systems with Finite Time Control Barrier Certificates and Temporal Logic". In: 57th IEEE Conference on Decision and Control, CDC 2018, Miami, FL, USA, December 17-19, 2018. IEEE. 1991–1996. DOI: 10.1109/CDC.2018.8619113. URL: https://doi.org/10.1109/CDC.2018.8619113.
- Steimle, L. N., D. L. Kaufman, and B. T. Denton. (2018). "Multi-Model Markov Decision Processes". *Optimization Online*.
- Such, J. M. (2017). "Privacy and autonomous systems". In: Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17. 4761–4767.
- Suilen, M., N. Jansen, M. Cubuktepe, and U. Topcu. (2020). "Robust Policy Synthesis for Uncertain POMDPs via Convex Optimization". In: *IJCAI*. ijcai.org. 4113–4120.
- Sutton, R. S. and A. G. Barto. (2018). Reinforcement Learning: An Introduction. MIT press.
- Tabuada, P. and D. Neider. (2016). "Robust Linear Temporal Logic". In: *Proc. Computer Science Logic*. Vol. 62. *LIPIcs.* 10:1–10:21.
- Tabuada, P. and G. J. Pappas. (2006). "Linear Time Logic control of linear systems". *IEEE Transactions on Automatic Control.* 51(12): 1862–1877. ISSN: 0018-9286.

Tanaka, T., H. Sandberg, and M. Skoglund. (2021). "Transfer-entropy-regularized markov decision processes". *IEEE Transactions on Automatic Control*.

- Thomas, P. S., G. Theocharous, M. Ghavamzadeh, I. Durugkar, and E. Brunskill. (2017). "Predictive Off-Policy Policy Evaluation for Nonstationary Decision Problems, with Applications to Digital Marketing". In: AAAI. 4740–4745.
- Tomita, T., A. Ueno, M. Shimakawa, S. Hagihara, and N. Yonezaki. (2017). "Safraless LTL synthesis considering maximal realizability". *Acta Informatica*. 54(7).
- Tong, Y., Z. Li, C. Seatzu, and A. Giua. (2018). "Current-state opacity enforcement in discrete event systems under incomparable observations". *Discrete Event Dynamic Systems*. 28(2): 161–182.
- Toro Icarte, R., T. Q. Klassen, R. Valenzano, and S. A. McIlraith. (2022). "Reward Machines:: Exploiting Reward Function Structure in Reinforcement Learning". *Journal of Artificial Intelligence Research*. 73.
- Tumova, J., G. C. Hall, S. Karaman, E. Frazzoli, and D. Rus. (2013).
  "Least-violating control strategy synthesis with safety rules". In:
  Proc. ACM International Conference on Hybrid Systems: Computation and Control.
- Vasile, C. I. and C. Belta. (2013). "Sampling-based temporal logic path planning". In: 2013 IEEE/RSJ International Conference on Intelligent Robots and Systems. IEEE. 4817–4822.
- Vlassis, N., M. L. Littman, and D. Barber. (2012). "On the Computational Complexity of Stochastic Controller Optimization in POMDPs". *ACM Trans. on Computation Theory.* 4(4): 12:1–12:8.
- Walraven, E. and M. Spaan. (2017). "Accelerated Vector Pruning for Optimal POMDP Solvers". In: Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 31. No. 1.
- Wang, R., X. He, R. Yu, W. Qiu, B. An, and Z. Rabinovich. (2020). "Learning efficient multi-agent communication: An information bottleneck approach". In: *International Conference on Machine Learning*. PMLR. 9908–9918.

Wen, M., I. Papusha, and U. Topcu. (2017). "Learning from Demonstrations with High-Level Side Information". In: *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, *IJCAI-17*. 3055–3061. DOI: 10.24963/ijcai.2017/426. URL: https://doi.org/10.24963/ijcai.2017/426.

- Wierstra, D., A. Förster, J. Peters, and J. Schmidhuber. (2007). "Solving deep memory POMDPs with recurrent policy gradients". In: *International Conference on Artificial Neural Networks (ICANN)*. Springer. 697–706.
- Wiesemann, W., D. Kuhn, and B. Rustem. (2013). "Robust Markov Decision Processes". *Mathematics of Operations Research.* 38(1): 153–183.
- Wimmer, R., N. Jansen, E. Ábrahám, J. Katoen, and B. Becker. (2014). "Minimal counterexamples for linear-time probabilistic verification". *Theoretical Computer Science*. 549: 61–100.
- Wimmer, R., N. Jansen, A. Vorpahl, E. Ábrahám, J.-P. Katoen, and B. Becker. (2013). "High-level counterexamples for probabilistic automata". In: *International Conference on Quantitative Evaluation of Systems*. Springer. 39–54.
- Winkler, T., S. Junges, G. A. Pérez, and J. Katoen. (2019). "On the Complexity of Reachability in Parametric Markov Decision Processes". In: *CONCUR*. Vol. 140. *LIPIcs*. Schloss Dagstuhl Leibniz-Zentrum für Informatik. 14:1–14:17.
- Wolff, E. M., U. Topcu, and R. M. Murray. (2012). "Robust Control of Uncertain Markov Decision Processes with Temporal Logic Specifications". In: CDC. 3372–3379.
- Wongpiromsarn, T., S. Karaman, and E. Frazzoli. (2011a). "Synthesis of provably correct controllers for autonomous vehicles in urban environments". In: 2011 14th International IEEE Conference on Intelligent Transportation Systems (ITSC). 1168–1173.
- Wongpiromsarn, T. (2010). "Formal Methods for Design and Verification of Embedded Control Systems: Application to an Autonomous Vehicle". *PhD thesis*. California Institute of Technology.
- Wongpiromsarn, T. (2020). "The Journey of Autonomous Vehicles". https://tichakorn.dev/post/av-journey/.

Wongpiromsarn, T., S. Mitra, R. M. Murray, and A. Lamperski. (2009). "Periodically Controlled Hybrid Systems: Verifying A Controller for An Autonomous Vehicle". In: *Hybrid Systems: Computation and Control.* Ed. by R. Majumdar and P. Tabuada. Vol. 5469. *Lecture Notes in Computer Science*. Springer. 396–410. ISBN: 0302-9743.

- Wongpiromsarn, T. and R. M. Murray. (2008). "Distributed Mission and Contingency Management for the DARPA Urban Challenge". In: International Workshop on Intelligent Vehicle Control Systems (IVCS).
- Wongpiromsarn, T., K. Slutsky, E. Frazzoli, and U. Topcu. (2021).
  "Minimum-Violation Planning for Autonomous Systems: Theoretical and Practical Considerations". In: 2021 American Control Conference.
- Wongpiromsarn, T., U. Topcu, and R. M. Murray. (2010a). "Automatic Synthesis of Robust Embedded Control Software". In: AAAI Spring Symposium on Embedded Reasoning: Intelligence in Embedded Systems. 104–111.
- Wongpiromsarn, T., U. Topcu, and R. M. Murray. (2010b). "Receding horizon control for temporal logic specifications". In: *HSCC*. Ed. by K. H. Johansson and W. Yi. ACM ACM. 101–110. ISBN: 978-1-60558-955-8.
- Wongpiromsarn, T., U. Topcu, and R. M. Murray. (2012). "Receding Horizon Temporal Logic Planning". *IEEE Transactions on Automatic Control.* 57(11): 2817–2830.
- Wongpiromsarn, T., U. Topcu, and R. M. Murray. (2013). "Synthesis of Control Protocols for Autonomous Systems". *Unmanned Systems*. 1(1): 21–39.
- Wongpiromsarn, T., U. Topcu, N. Ozay, H. Xu, and R. M. Murray. (2011b). "TuLiP: A Software Toolbox for Receding Horizon Temporal Logic Planning". In: *Proceedings of the 14th International Conference on Hybrid Systems: Computation and Control. HSCC '11.* Chicago, IL, USA: Association for Computing Machinery. 313–314.
- Wu, D. and X. Koutsoukos. (2008). "Reachability Analysis of Uncertain Systems using Bounded-Parameter Markov Decision Processes". *Artificial Intelligence*. 172(8-9): 945–954.

Xu, Z., I. Gavran, Y. Ahmad, R. Majumdar, D. Neider, U. Topcu, and B. Wu. (2020). "Joint inference of reward machines and policies for reinforcement learning". In: Proceedings of the International Conference on Automated Planning and Scheduling. Vol. 30, 590– 598.

- Yang, Y., P. Gohari, and U. Topcu. (2021). "On the Privacy Risks of Deploying Recurrent Neural Networks in Machine Learning". arXiv preprint arXiv:2110.03054.
- Yu, Y., P. Manolios, and L. Lamport. (1999). "Model Checking TLA+ Specifications". In: *Conference on Correct Hardware Design and Verification Methods*. 54–66.
- Yuan, Y.-x. (2015). "Recent Advances in Trust Region Algorithms". Mathematical Programming. 151(1): 249–281.
- Zhang, T. and Q. Zhu. (2018). "Hypothesis testing game for cyber deception". In: *International Conference on Decision and Game Theory for Security*. Springer. 540–555.