

CPSGrader: Auto-Grading and Feedback Generation for Cyber-Physical Systems Education

Garvit Juniwal



Electrical Engineering and Computer Sciences
University of California at Berkeley

Technical Report No. UCB/EECS-2014-237

<http://www.eecs.berkeley.edu/Pubs/TechRpts/2014/EECS-2014-237.html>

December 21, 2014

Report Documentation Page

*Form Approved
OMB No. 0704-0188*

Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

| | | | | | |
|---|------------------------------------|---|---|----------------------------------|---------------------------------|
| 1. REPORT DATE 21 DEC 2014 | 2. REPORT TYPE | 3. DATES COVERED 00-00-2014 to 00-00-2014 | | | |
| 4. TITLE AND SUBTITLE CPSGrader: Auto-Grading and Feedback Generation for Cyber-Physical Systems Education | | 5a. CONTRACT NUMBER | | | |
| | | 5b. GRANT NUMBER | | | |
| | | 5c. PROGRAM ELEMENT NUMBER | | | |
| 6. AUTHOR(S) | | 5d. PROJECT NUMBER | | | |
| | | 5e. TASK NUMBER | | | |
| | | 5f. WORK UNIT NUMBER | | | |
| 7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) University of California at Berkeley,Electrical Engineering and Computer Sciences,Berkeley,CA,94720 | | 8. PERFORMING ORGANIZATION REPORT NUMBER | | | |
| 9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) | | 10. SPONSOR/MONITOR'S ACRONYM(S) | | | |
| | | 11. SPONSOR/MONITOR'S REPORT NUMBER(S) | | | |
| 12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited | | | | | |
| 13. SUPPLEMENTARY NOTES | | | | | |
| 14. ABSTRACT Formal methods and machine learning together have the potential to enhance technologies for education. In this thesis, we consider the problem of designing CPSGrader, an automatic grader for laboratory-based courses in the area of cyber-physical systems. The work is motivated by a UC Berkeley course in which students program a robot for speci ed navigation tasks. Given a candidate student solution (control program for the robot), CPSGrader rst checks whether the robot performs the task correctly under a representative set of environment conditions. If it does not, CPSGrader automatically generates feedback hinting at possible errors in the program. CPSGrader is based on a novel notion of constrained parameterized tests based on signal temporal logic (STL) that capture symptoms pointing to success or causes of failure in traces obtained from a realistic simulator. We de ne and solve the problem of synthesizing constraints on a parameterized test such that it is consistent with a set of reference solutions with and without the desired symptom. We also develop a clustering-based active learning technique that selects from a large database of unlabeled solutions, a small number of reference solutions for the expert to label. The goal is to achieve better accuracy of fault identi cation with fewer reference solutions as compared to random selection. We demonstrate the e ectiveness of CPSGrader using two data sets: one obtained from an on-campus laboratory-based course at UC Berkeley, and the other from a massive open online course (MOOC) o ering. In addition, CPSGrader was successfully deployed in the laboratory section of this MOOC on the edX platform. | | | | | |
| 15. SUBJECT TERMS | | | | | |
| 16. SECURITY CLASSIFICATION OF: | | | 17. LIMITATION OF ABSTRACT Same as Report (SAR) | 18. NUMBER OF PAGES 60 | 19a. NAME OF RESPONSIBLE PERSON |
| a. REPORT unclassified | b. ABSTRACT unclassified | c. THIS PAGE unclassified | | | |

Copyright © 2014, by the author(s).
All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

Acknowledgement

First, and foremost, I would like to thank my advisor, Professor Sanjit Seshia, from whom I have learned an incredible amount from him in all aspects of academic life. I would also like to thank my other reader, Professor Edward Lee, for feedback and comments. I would like to acknowledge Alexandre Donz{\e} and Jeff Jensen for their invaluable support in discussing and implementing ideas. I would like to thank Sakshi Jain for being a great partner in a class project that became a part of this thesis. This work was funded in part from NSF ExCAPE project (CCF-1139138), TerraSwarm, one of six centers of STARnet, a Semiconductor Research Corporation program sponsored by MARCO and DARPA, and National Instruments Inc.

**CPSGrader: Auto-Grading and Feedback Generation for Cyber-Physical
Systems Education**

by

Garvit Juniwal

B.Tech. Indian Institute of Technology Bombay 2012

A thesis submitted in partial satisfaction
of the requirements for the degree of

Master of Science

in

Electrical Engineering and Computer Sciences

in the

GRADUATE DIVISION

of the

UNIVERSITY OF CALIFORNIA, BERKELEY

Committee in charge:

Professor Sanjit A. Seshia, Chair
Professor Edward A. Lee

Fall 2014

The thesis of Garvit Juniwal is approved.

Chair

Date

Date

University of California, Berkeley
Fall 2014

CPSGrader: Auto-Grading and Feedback Generation for Cyber-Physical Systems
Education

Copyright © 2014

by

Garvit Juniwal

Abstract

CPSGrader: Auto-Grading and Feedback Generation for Cyber-Physical Systems
Education

by

Garvit Juniwal

Master of Science in Electrical Engineering and Computer Sciences

University of California, Berkeley

Professor Sanjit A. Seshia, Chair

Formal methods and machine learning together have the potential to enhance technologies for education. In this thesis, we consider the problem of designing CPSGrader, an automatic grader for laboratory-based courses in the area of cyber-physical systems. The work is motivated by a UC Berkeley course in which students program a robot for specified navigation tasks. Given a candidate student solution (control program for the robot), CPSGrader first checks whether the robot performs the task correctly under a representative set of environment conditions. If it does not, CPSGrader automatically generates feedback hinting at possible errors in the program. CPSGrader is based on a novel notion of constrained parameterized tests based on signal temporal logic (STL) that capture symptoms pointing to success or causes of failure in traces obtained from a realistic simulator. We define and solve the problem of synthesizing constraints on a parameterized test such that it is consistent with a set of reference solutions with and without the desired symptom. We also develop a clustering-based active learning technique that selects from a large database of unlabeled solutions, a small number of reference solutions for the expert to label. The goal is to achieve better accuracy of fault identification with fewer reference solutions as compared to random selection. We demonstrate the effectiveness of CPSGrader using two

data sets: one obtained from an on-campus laboratory-based course at UC Berkeley, and the other from a massive open online course (MOOC) offering. In addition, CPSGrader was successfully deployed in the laboratory section of this MOOC on the edX platform.

Professor Sanjit A. Seshia
Thesis Committee Chair

To my Mother, who serves as my primary source of inspiration.

Contents

| | |
|---|-----------|
| Contents | ii |
| Acknowledgements | iv |
| 1 Introduction | 1 |
| 1.1 Formal Methods in Education | 1 |
| 1.2 Target Laboratory Course | 2 |
| 1.3 Problem Motivation | 3 |
| 1.4 Contributions | 7 |
| 1.5 Related Work | 8 |
| 2 Background | 9 |
| 2.1 Signals, Controllers, and Environments | 9 |
| 2.2 Signal Temporal Logic | 10 |
| 2.3 Defects and Faults | 12 |
| 2.4 Dynamic Time Warping Distance (DTW) | 13 |
| 2.5 Density-Based Spatial Clustering (DBSCAN) | 13 |
| 3 Synthesis of Test Benches | 14 |
| 3.1 Constrained Parametrized Tests | 14 |
| 3.2 Synthesis of Test Bench Constraints | 17 |
| 3.3 Computing the Satisfaction Region | 19 |
| 3.4 Adequate Test Samples for Grading | 22 |
| 3.5 Related Work | 25 |
| 4 Clustering-Based Active Learning | 26 |

| | | |
|----------|--|-----------|
| 4.1 | Iterative Synthesis of Test Benches by Active Learning | 27 |
| 4.2 | Clustering with Precomputed Distances | 27 |
| 4.3 | Selection of Training Data from Clusters | 29 |
| 4.4 | Related Work | 30 |
| 5 | Evaluation | 31 |
| 5.1 | Obstacle Avoidance | 33 |
| 5.2 | Hill Climbing | 35 |
| 5.3 | Accuracy of Classification | 38 |
| 5.4 | Grade Correlation | 39 |
| 5.5 | Effectiveness of Iterative Synthesis | 40 |
| 5.6 | Investigating Unknown Faults Using Clustering | 41 |
| 5.7 | Discussion | 43 |
| 6 | Conclusion and Future Work | 44 |
| | Bibliography | 46 |
| | References | 46 |
| A | STL Semantics | 48 |

Acknowledgements

First, and foremost, I would like to thank my advisor, Professor Sanjit Seshia, from whom I have learned an incredible amount from him in all aspects of academic life. I would also like to thank my other reader, Professor Edward Lee, for feedback and comments. I would like to acknowledge Alexandre Donzé and Jeff Jensen for their invaluable support in discussing and implementing ideas. I would like to thank Sakshi Jain for being a great partner in a class project that became a part of this thesis. This work was funded in part from NSF ExCAPE project (CCF-1139138), TerraSwarm, one of six centers of STARnet, a Semiconductor Research Corporation program sponsored by MARCO and DARPA, and National Instruments Inc.

Chapter 1

Introduction

1.1 Formal Methods in Education

Massive open online courses (MOOCs) [1] and related technological advances promise to bring world-class education to anyone with Internet access. Additionally, MOOCs present a range of problems to which the field of formal methods has much to contribute. These include *automatic grading*, *automated exercise generation*, and *virtual laboratory environments*. In automatic grading, a computer program verifies that a candidate solution provided by a student is “correct”, i.e., that it meets certain instructor-specified criteria (the specification). In addition, and particularly when the solution is incorrect, the automatic grader (henceforth, *auto-grader*) should provide feedback to the student as to where he/she went wrong. Automatic exercise generation is the process of synthesizing problems (with associated solutions) that test students’ understanding of course material, often starting from instructor-provided sample problems. Finally, for courses involving laboratory assignments, a virtual laboratory (henceforth, *lab*) seeks to provide the remote student with an experience similar to that provided in a real, on-campus lab.

Lab-based courses that are not software-only pose a particular technical challenge. An example of such a course is *Introduction to Embedded Systems* at UC Berkeley [2]. In this course, students not only learn theoretical content on modeling, design, and analysis [3],

but also perform lab assignments on programming an embedded platform interfaced to a mobile robot [4]. What would an online lab assignment in embedded systems look like? In an ideal world, we would provide an infrastructure where students can log in remotely to a computer which has been preconfigured with all development tools and laboratory exercises; in fact, pilot projects exploring this approach have already been undertaken (e.g., see [5]). However, in the MOOC setting, the large numbers of students makes such a remotely-accessible physical lab expensive and impractical. A virtual lab environment, driven by simulation of real-world environments, appears to be the only solution at present. For example, the MIT circuits course (MITx 6.002x) uses rudimentary circuit simulation software [6].

In this thesis, we formalize the auto-grading problem for a virtual lab environment in the field of embedded and cyber-physical systems (CPS). The virtual lab under consideration is the one designed for EECS149.1x [7], a MOOC on Cyber-Physical Systems offered on the edX platform, based on the afore-mentioned on-campus course. Next, we give the details of this virtual laboratory under consideration.

1.2 Target Laboratory Course

The embedded systems laboratory course offered at University of California, Berkeley employs a custom mobile robotic platform called the Cal Climber [8], [9]. The Cal Climber is based on the commercially-available iRobot Create (derived from the iRobot-Roomba autonomous vacuum cleaner) (Fig. 1.1a), and the National Instruments myRIO embedded controller. This off-the-shelf platform is capable of driving, sensing bumps and cliffs, executing simple scripts, and communicating with an external controller. This configuration demonstrates the composition of cyber-physical systems, where a robotics platform is modeled as a sub-system and treated as a collection of sensors and actuators potentially located beyond a network boundary. The problem statement centers on model-based design and is given as follows (paraphrased from [4]):

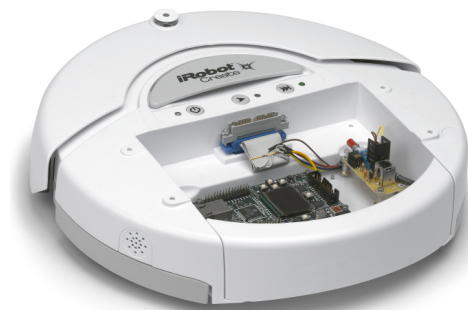
Design a StateChart to drive the Cal Climber. On level ground, your robot should drive straight. When an obstacle is encountered, such as a cliff or an object, your robot should navigate around the object and continue in its original orientation. On an incline, your robot should navigate uphill, while still avoiding obstacles. Use the accelerometer to detect an incline and as input to a control algorithm that maintains uphill orientation.

Source files distributed with the Cal Climber laboratory are structured such that students only need to implement a function that receives as arguments the most recent values of the accelerometer and robot sensors and returns desired wheel speeds. This function is called repeatedly at short regular intervals of time (60 ms in our case) with most recent sensor and accelerometer data. Students implement this function for controlling the Cal Climber. In the on-campus course, students first prototype their controller to work within a simulated environment (without any auto-grading) based on the LabVIEW Robotics Environment Simulator by National Instruments. The simulator is based on the Open Dynamics Engine [10] rigid body dynamics software that can simulate robots in a virtual environment(Fig. 1.1b). In EECS149.1x, the afore-mentioned online version of the course, the same simulator, extended with the auto-grader described in the present paper, has been used(Fig. 1.1c).

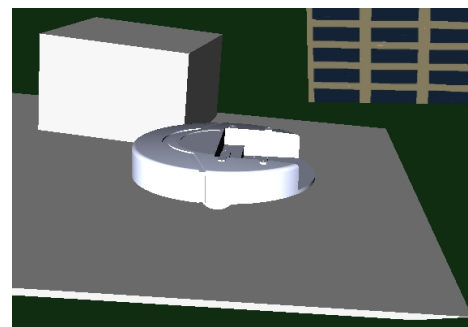
We refer to the functions implemented by students as *solutions* or *controllers*. A solution is evaluated in a collection of environments against a collection of goal and fault properties, forming *test benches* (a notion formalized in the following sections). For this purpose, the simulator allows to define customized environments (with walls, objects, obstacles, ramps, etc) described in XML files and we further extended its API to facilitate the exportation of simulation traces to external property monitoring tools.

1.3 Problem Motivation

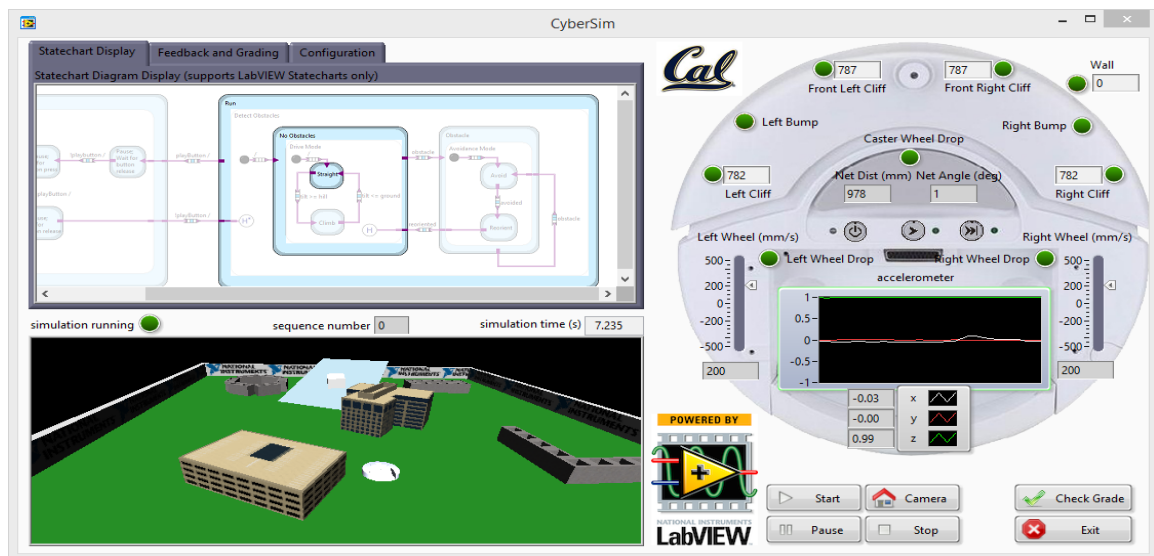
As mentioned before, in this thesis we tackle the problem of auto-grading a CPS lab. Auto-grading is a verification and debugging problem where the objective is to be able to check whether a student's solution meets the desired goals and also provide feedback pointing to possible causes of failure. The main point we make here is that the dynamical model



(a)



(b)



(c)

Figure 1.1: (a) Cal Climber laboratory platform. (b) Cal Climber in the LabVIEW Robotics Environment Simulator. (b) Simulator with auto-grading functionality used in EECS 149.1x

for the virtual lab is so complex that simulation is currently the only verification method that can be practically employed. Thus, the auto-grader is based on simulation-based verification. The high-level approach, previously hinted at in a position paper [11], is as follows. Correctness properties are formalized in *signal temporal logic* (STL) [12]. Simulation test benches are created by a combination of manual environment setup and simulation-based falsification implemented in tools such as Breach [13]. For each lab assignment, there is an *end-to-end correctness property*, hereafter referred to as the *goal* property. If the goal is satisfied, the student solution (hereafter referred to as a *controller*) is deemed correct. Otherwise, it is incorrect, and more analysis must be performed to identify the mistake (fault) and provide feedback. This latter analysis is based on monitoring simulation traces of the student controller on a library of known faults, also formalized in STL. If any of these “fault properties” hold for a student controller, they are provided to the student as feedback.

This approach, though straightforward on the surface, requires further technical advances to be effective. The first problem is that the STL properties that encode both goal and fault properties reference parameters that can vary over the set of environments and student controllers; in fact, such variation must be allowed. For example, in a real lab, students may program robots to move at different velocities while performing obstacle avoidance. If the goal of the lab is only to correctly avoid an obstacle, the speed at which it does so is irrelevant. However, given the variations in the controllers students design, setting a reasonable range for parameters such as time or velocity in STL properties can be tricky. Similarly, environments can also be parametric (for example, the location of obstacles) and tests should be synthesized in a manner that accounts for these variations. Thus, an effective approach to auto-grading CPS labs requires one to solve a certain *parameter synthesis* problem.

We formalize this parameter synthesis problem and give an algorithm to solve it. First, we define the notion of a *parametrized test* which is a combination of a parametrized environment and a parametrized STL (PSTL) property. A parametrized test is thus a collection of tests. However, as discussed above, one needs to impose a constraint on this collection

to capture “legal” variations in student solutions. Such a constraint, termed a *sub-domain*, defines the allowed set of parameter valuations. However, manually computing this sub-domain is tedious and error-prone. We therefore give an algorithmic approach to synthesize the sub-domain from reference controllers that should/should not pass the test bench. In practice, it is easier for instructors to provide such reference controllers than it is to manually compute sub-domains. In machine learning terminology, this can be thought of as the *training* phase. The resulting *constrained parameterized test bench* then becomes the “specification” that determines whether a student solution is correct, and, if not, which fault is present. In machine learning terminology, this would be the *classification* phase. Further, we identify a property, *monotonicity*, under which we can efficiently compute the sub-domain, and which holds for the lab of interest.

Another issue with this approach is the availability of “positive” and “negative” reference controllers. An instructor has to manually look at the simulation video to decide whether a particular controller is good or bad and then it can be used for training the test bench. In essence, labeling of controllers is an expensive manual process. We formulate the problem of obtaining labeled data as an active learning problem. We give a clustering-based active learning methodology that takes as input a large set of unlabeled controllers collected over various stages of development and chooses the controllers that an instructor should label to get high accuracy of classification with fewer number of training examples as compared to random selection. Since the simulation traces are timed sequences of multidimensional variables that capture environment and state data, we choose *dynamic time warping* distance [14] as a measure of dissimilarity between controller behavior and use that for clustering.

We believe clustering is useful because amongst many student solutions, the total number solution strategies are still few. Furthermore, any two solutions that follow the same strategy will likely have the same set of faults present or absent. Hence, if some clustering technique can identify each strategy as a separate cluster, then choosing one example from each cluster should account for a training set with good coverage and the synthesized test bench will have high accuracy.

Any auto-grader must have at least two desirable properties: *accuracy* and *efficiency*. The former means that the auto-grader must correctly classify right and wrong student solutions, and for wrong solutions, correctly explain the mistake (fault). The latter means that it must run efficiently in practice. For efficiency, we show how monotonicity can be exploited again to avoid the need to run the entire constrained parametric test bench. Instead, we define the notion of an *adequate* test sample and show that it is much smaller in practice than the entire constrained test bench. We also provide an experimental evaluation on the on-campus lab demonstrating that our approach is both accurate and efficient in practice. We also test our active learning approach and show that selection of training examples based on clustering leads to higher accuracy of classification as compared to random selection of the same number of training examples, and therefore it can lead to reduced overhead for instructors in providing labeled data.

1.4 Contributions

To summarize, the main novel contributions of this work are:

- A formalization of the auto-grading problem for simulation-based virtual laboratories in cyber-physical systems,
- A formalization of the problem of synthesizing a constrained parametric test bench for the auto-grader along with an efficient solution approach based on monotonicity,
- A novel clustering-based active learning approach to aid in generation of labeled training data for the synthesis algorithm, and
- An empirical evaluation demonstrating the accuracy and efficiency of CPSGrader, the auto-grader for the on-campus embedded systems lab, and also the effectiveness of the clustering-based active learning, on a database of student solutions from: (1) on-campus offering of the course EECS 149, and (2) the online edition of the same course on edX (EECS 149.1x.)

Note that a large part of this thesis is based on the EMSOFT 2014 paper [15], joint work with Alexandre Donzé, Jeff C. Jensen, and Sanjit A. Seshia.

1.5 Related Work

There is a growing number of efforts to incorporate formal methods into technologies for education. Singh et al. [16] present an approach to automatically generate problems in high-school algebra. Sadigh et al. [17] show how the problem of generating variants of exercises in an Embedded Systems textbook [3] can be mapped to standard problems in formal methods and apply some of these methods to classes of exercises. Singh et al. [18] present an auto-grader for a Python programming course, where, similar to the present paper, feedback is generated based on a library of common mistakes, but, differently, the technical approach uses an encoding to SAT-based program synthesis. Alur et al. [19] consider auto-grading DFA construction problems, providing a novel blend of three techniques for assigning partial grades for incorrect answers. This thesis proposes different formalisms and algorithms, and represents the first auto-grader for lab assignments in the area of embedded, cyber-physical systems.

Related work on parameter synthesis for temporal logic and use of clustering for active learning is covered in later chapters.

The outline of the thesis is as follows. We introduce basic terminology and background results in Ch. 2. In Ch. 3, we describe the main theoretical contributions, including our formalization of the problem of synthesis of test benches and solution approach. In Ch. 4, we describe the clustering-based active learning approach that serves as an aid to the synthesis algorithm. Experimental results are given in Sec. 5. We conclude with future directions Ch. 6.

Chapter 2

Background

2.1 Signals, Controllers, and Environments

Definition 1 (*Signal*) A (uni-dimensional) signal is a function mapping the time domain $\mathbb{T} = \mathbb{R}_{\geq 0}$ to the reals \mathbb{R} .

Boolean signals, used to represent discrete dynamics, are signals whose values are restricted to *false* (denoted \perp) and *true* (denoted \top). Vectors in \mathbb{R}^n with $n > 1$ are denoted in bold fonts and their components are indexed from 1 to n , for example, $\mathbf{p} = (p_1, \dots, p_n)$. Likewise, a *multi-dimensional signal* \mathbf{x} is a function from \mathbb{T} to \mathbb{R}^n such that $\forall t \in \mathbb{T}$, $\mathbf{x}(t) = (x_1(t), \dots, x_n(t))$. We will use the term “signal” to also refer to multi-dimensional signals.

Definition 2 (*Controller*) A controller C is a (deterministic) dynamical system that takes as input a signal $\mathbf{y}(t)$ and computes an output signal $\mathbf{u}(t)$. It is common to drop time, and say $\mathbf{u} = C(\mathbf{y})$.

Note that we make no assumption about how a controller computes its output. A controller can have discrete or continuous dynamics or it can be a hybrid system. As an example, a program running on the Cal Climber is a controller that takes bumps and

cliff sensors signals, and accelerometer data as input $\mathbf{y}(t) = (\text{bump}(t), \text{cliff}(t), \text{accel}(t))$, and responds with the desired left and right wheel speeds as output $\mathbf{u}(t) = (\text{lws}(t), \text{rws}(t))$.

Definition 3 (*Environment*) An environment E for a controller C is a dynamical system generating all inputs to C .

As before, we make no assumptions about the form of the environment. All we assume is the existence of a simulator that can take representations of E and C , compose them, and produce execution traces of the composite system. In other words, the simulator is an oracle that gives semantics to the composite system $E\|C$.

We only consider deterministic environments, i.e., the composition of a controller and an environment has deterministic behavior. For example, an arena composed of obstacles and hills on level ground is an environment for the Cal Climber controller. Formally, a trace $\text{sim}(C, E)$ is a multi-dimensional signal $(\mathbf{x}(t), \mathbf{y}(t), \mathbf{u}(t))$ consisting of the inputs \mathbf{y} and outputs \mathbf{u} of the controller and optionally other signals \mathbf{x} regarding the state of the environment. For example, the position and orientation (in the plane of the ground) of the robot in the arena $\mathbf{x}(t) = (\text{pos}(t), \text{angle}(t))$ are a part of the observable environment state. By varying the environment, or the property being verified on the composition (see Sec. 2.2), the instructor can test different features of the controller.

2.2 Signal Temporal Logic

Since propositional (linear) temporal logic was introduced by Amir Pnueli [20], variants have also been proposed. Temporal logics to reason about real-time signals, such as Timed Propositional Temporal Logic [21], and Metric Temporal Logic (MTL) [22] were introduced later to deal with dense-time signals. More recently, Signal Temporal Logic [12] was proposed in the context of analog and mixed-signal circuits to deal with dense-time signals taking values over both discrete and continuous domains. We use STL as the specification language for the Embedded Systems lab assignment. Goals that the robotic controller must achieve are expressed as STL properties.

The primitive constraints, or *predicates*, in STL take the form $\mu \doteq f(\mathbf{x}) \sim \pi$, where f is a scalar-valued function over the signal \mathbf{x} , $\sim \in \{<, \leq, \geq, >, =, \neq\}$, and π is a real number. Temporal formulas are formed using temporal operators, “always” (denoted as \square), “eventually” (denoted as \diamond) and “until” (denoted as \mathbf{U}). Each temporal operator is indexed by intervals of the form (a, b) , $(a, b]$, $[a, b)$, $[a, b]$, (a, ∞) or $[a, \infty)$ where each of a, b is a non-negative real-valued constant. If I is an interval, then an STL formula is written using the grammar:

$$\varphi := \top \mid \mu \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \varphi_1 \mathbf{U}_I \varphi_2$$

The always and eventually operators are defined as special cases of the until operator in the standard way: $\square_I \varphi \triangleq \neg \diamond_I \neg \varphi$, $\diamond_I \varphi \triangleq \top \mathbf{U}_I \varphi$. When the interval I is omitted, we use the default interval of $[0, +\infty)$. The semantics of STL formulas are defined informally as follows. The signal \mathbf{x} satisfies $f(\mathbf{x}) > 10$ at time t (where $t \geq 0$) if $f(\mathbf{x}(t)) > 10$. It satisfies $\varphi = \square_{[0,2)} (x > -1)$ if for all time $0 \leq t < 2$, $x(t) > -1$. The signal x_1 satisfies $\varphi = \diamond_{[1,2)} x_1 > 0.4$ iff there exists time t such that $1 \leq t < 2$ and $x_1(t) > 0.4$. The two-dimensional signal $\mathbf{x} = (x_1, x_2)$ satisfies the formula $\varphi = (x_1 > 10) \mathbf{U}_{[2.3,4.5]} (x_2 < 1)$ iff there is some time u where $2.3 \leq u \leq 4.5$ and $x_2(u) < 1$, and for all time v in $[2.3, u)$, $x_1(v)$ is greater than 10. The formal semantics of STL can be found in [12] and is given in Appendix A.

Parametric Signal Temporal Logic (PSTL) is an extension of STL introduced in [23] to define *template formulas* containing unknown parameters. Syntactically speaking, a PSTL formula is an STL formula where numeric constants, either in the constraints given by the predicates μ or in the time intervals of the temporal operators, can be replaced by symbolic parameters.

An STL formula is obtained by pairing a PSTL formula with a valuation function that assigns a value to each symbolic parameter. For example, consider the PSTL formula $\varphi(\pi, \tau) = \square_{[0,\tau]} x > \pi$, with symbolic parameters π (scale) and τ (time). The STL formula $\square_{[0,10]} x > 1.2$ is an instance of φ obtained with the valuation $v = \{\tau \mapsto 10, \pi \mapsto 1.2\}$.

2.3 Defects and Faults

A controller is usually designed to meet certain *goals*. For example, the Cal Climber controller should be able to navigate around obstacles and climb hills. To talk about grading and feedback generation, we introduce some relevant terminology from the fault testing and diagnosis literature.

Definition 4 (*Defect, symptom and fault*) *Given a controller and an environment with some desired goals,*

- *A defect is a bug in the controller implementation that leads to failure in meeting goals;*
- *A symptom is an interesting pattern in a simulation trace of the controller-environment composition that can be characterized, for example, using STL, and*
- *A fault is a symptom that is present in a trace as a result of some defect in the controller.*

A general symptom, such as the inability to meet an end-to-end correctness goal (for example, obstacle avoidance), is a fault that could be the result of multiple defects in the controller. On the other hand, certain specific faults could be mapped to specific kinds of defects. As an example, consider an obstacle avoidance strategy for the Cal Climber controller, implemented in a language like C. The strategy states that every time the `bump` sensor signal indicates a bump, the robot backs up, moves some distance to either right or left and then re-orientes by turning in-place until the heading direction is same as the original direction angle_0 . A controller will check the guard $|\text{angle}(t) - \text{angle}_0| \leq \epsilon$ for some small $\epsilon > 0$ to determine when to stop turning in the re-orientation mode. A defect can be introduced by replacing this guard by the exact equality check $\text{angle}(t) == \text{angle}_0$. This modification usually leads to failure in practice, because the controller implementation polls its sensors at certain intervals, and therefore, it is highly unlikely that the sensor value at some polled time t , $\text{angle}(t)$, will be exactly angle_0 . The fault resulting from this defect is that in the re-orientation mode, the robot keeps turning in-place while making full circles multiple times. We call this the *circle* fault and will revisit it again in the paper.

The ability to classify traces that present a fault from those that do not is important for auto-grading. Using this classification, we can not only separate correct solutions from incorrect ones but also generate diagnostic feedback for failed traces by monitoring for relevant faults that will likely correspond to known defects.

2.4 Dynamic Time Warping Distance (DTW)

In time series analysis, dynamic time warping (DTW) [24] is an algorithm for measuring similarity between two temporal sequences which may vary in time or speed. For instance, similarities in movement patterns of two Cal Climber controllers could be detected using DTW, even if one was moving at a faster wheel speed than the other, or if the matching sub-patterns occur at different absolute times. DTW has been applied to temporal sequences of video, audio, and graphics data. In general, DTW is a method that calculates an optimal match between two given sequences (e.g. time series) with certain restrictions. The sequences are “warped” non-linearly in the time dimension to compute the optimal sequence alignment for which the two sequences match closely. Hence, the distance between two sequences is agnostic of shifting and scaling, making DTW suitable for our purpose. DTW can be extended to multi-dimensional timed sequences [14].

2.5 Density-Based Spatial Clustering (DBSCAN)

Density-based spatial clustering (DBSCAN) [25] clusters the samples based on provided estimation of the density of cluster nodes. It can take as input pre-computed pairwise distances between samples and does not need the feature vectors to be given explicitly. The number of clusters does not need to be specified in advance. DBSCAN starts off by finding small groups of points that are very close to each other and marks these groups as potential clusters. It then expands each cluster by including other close neighbours. It can find arbitrarily shaped clusters and is robust to outliers. These features make it a good fit for our application.

Chapter 3

Synthesis of Test Benches

In this chapter, we formally define the auto-grading problem, the technical challenge in synthesizing a constrained parametrized set of tests, and our approach to solve this problem.

For the purpose of examples in this chapter, we always assume the controller is a Cal Climber program and the environment is an arena with one robot, multiple obstacles and fixed inclines (flat rectangular planks) placed on level ground. Positions in the arena are given using x , y , and z coordinates (in meters). Orientation in the $x - y$ plane is given by the yaw angle varying from -180 to 180 degrees, increasing in counter-clockwise direction with 0 aligned with y -axis. The initial position and orientation of the robot is also a part of the environment.

3.1 Constrained Parametrized Tests

One of the fundamental notions for auto-grading is that of a test.

Definition 5 (*Test*) A pair (E, φ) of an environment E and an STL formula φ is called a test. A test passes for a controller C if and only if $\text{sim}(C, E) \models \varphi$.

Note that our definition of a test is different from the more common definition because in addition to controller inputs (provided in form of an environment), it also contains an assertion specified via STL.

For the end-to-end correctness property (goal), we will employ the convention that the STL formula φ in a test for this goal is the *negation* of the property that we want to hold. In other words, if a test “passes,” it actually means that the correctness property did not hold for that test case. The reason for this convention is that it allows us to treat STL formulas encoding correctness goals and fault symptoms in a symmetric fashion, something that is required for the main technical results of this paper. Hereafter we will treat the STL property as specifying a fault unless explicitly stated otherwise.

Example 1 Consider an environment E_0 with a square obstacle occupying the region $[4.5, 5.5] \times [5.0, 5.5]$. The initial position of the robot is $\langle 5.0, 4.9 \rangle$ and the initial orientation is 0. Consider the STL property $\varphi = \square(\text{pos.y} \leq 5.5)$ which states that the robot is never able to reach a point with y coordinate more than 5.5. If the test (E_0, φ) passes, we can assert that the robot did not meet the goal of being able to avoid the obstacle.

Consider a vector of symbolic parameters $\mathbf{p} = (p_1, p_2, \dots, p_n)$. A valuation function v maps each symbolic parameter to a concrete value (for example, in \mathbb{R}^n) and $v(p_i)$ denotes the value of parameter p_i in v . The set of all possible valuations of \mathbf{p} , its domain, is \mathfrak{U} .

Definition 6 (*Parametrized Test*) A parametrized environment is an environment with unknown parameters, denoted $E(\mathbf{p})$. A parametrized test $\Gamma(\mathbf{p}) = (E(\mathbf{p}), \varphi(\mathbf{p}))$ is a pair of a parametrized environment $E(\mathbf{p})$ and a PSTL formula $\varphi(\mathbf{p})$. Given any valuation $v \in \mathfrak{U}$, $\Gamma(v(\mathbf{p})) = (E(v(\mathbf{p})), \varphi(v(\mathbf{p})))$ is a concrete test.

Example 2 Consider the same environment E_0 from Example 1 except that the initial orientation of the robot is an unknown parameter θ_{init} that can take one of two possible values $\{-45, 45\}$. (See Figure 3.1a.) Consider the PSTL property $\varphi_0(\pi) = \square(\text{pos.y} > 5.5 \Rightarrow \pi_l < \text{pos.x} < \pi_u)$, where $\pi = (\pi_l, \pi_u)$, with unknown parameters π_l and π_u that can

take one of three possible values $\{-\infty, 5.0, \infty\}$ each. The property states that if the robot is able to get around the obstacle and reach a point $\text{pos.y} > 5.5$, then pos.x is always in the interval (π_l, π_u) . The pair $\Gamma_0(\theta_{init}, \pi) = (E_0(\theta_{init}), \varphi_0(\pi))$ is an example of a parameterized test.

Definition 7 (*Satisfaction Region*) The satisfaction region $\Omega(C, \Gamma(\mathbf{p}))$ of a controller C on a parametrized test $\Gamma(\mathbf{p})$ is the set of all valuations v of \mathbf{p} such that $\Gamma(v(\mathbf{p}))$ passes for C , i.e., $\Omega(C, \Gamma(\mathbf{p})) = \{v \in \mathfrak{U} \mid \Gamma(v(\mathbf{p})) \text{ passes for } C\}$.

Definition 8 (*Test Bench*) Given a parameterized test $\Gamma(\mathbf{p})$ and a set of valuations $\rho \subseteq \mathfrak{U}$, the pair $(\Gamma(\mathbf{p}), \rho)$ is called a constrained parametrized test, simply referred to as test bench. The set of valuations ρ is called the sub-domain of the test bench. We say that test bench $(\Gamma(\mathbf{p}), \rho)$ succeeds for a controller C iff there exists a $v \in \rho$ such that $\Gamma(v(\mathbf{p}))$ passes for C or equivalently, $\Omega(C, \Gamma(\mathbf{p})) \cap \rho$ is non-empty.

Since a test bench typically includes both the goal properties (determining whether a student controller is correct or not) and the fault properties (determining the mistakes the student made), the crux of the auto-grading problem is to *synthesize a test bench* that can *accurately* classify an “unlabeled” controller as correct/incorrect and with the fault(s), if any. Treating goal and fault properties uniformly, we seek to synthesize a test bench to classify whether an unlabeled controller exhibits faulty behaviors.

To auto-grade, for every known fault, we create a test bench. If the test bench succeeds for an unlabeled controller, we can conclusively label it as one exhibiting faulty behavior. The sub-domain of a test bench essentially identifies the set of tests that indicate the presence of the fault. As mentioned earlier, a test bench can also be used in a similar way to check if a given controller meets goal requirements by formulating the failure to meet the goal as a fault.

Example 3 Consider the parameterized test Γ_0 from Example 2. Consider the sub-domain $\rho_0 = \{[\theta_{init} \mapsto 45, \pi \mapsto (5.0, \infty)], [\theta_{init} \mapsto -45, \pi \mapsto (-\infty, 5.0)]\}$. For a controller, if either

of valuations in ρ_0 leads to a test that passes, it provides good evidence that the robot is either unable to avoid the obstacle or it is not able to proceed in the initial direction. (See Figure 3.1a.) So the test bench $(\Gamma_0(\theta_{init}, \pi), \rho_0)$ can be used to capture this failure to meet desired goals.

Example 4 Consider an environment E_1 with a fixed incline s.t. the uphill direction is along the orientation 0. The initial location of the robot is fixed at the center of the bottom boundary of the incline. The initial orientation of the robot is a parameter $\theta_{init} \in [-180, 180]$. We wish to determine whether a given controller (in an initial orientation pointing towards the incline) fails to climb within reasonable time. This can be expressed via the STL property $\varphi_1(h, \tau) = \square_{[0, \tau]}(\text{pos.z} \leq h)$, that states that the robot is not able to reach the height h , within time τ . The parametrized test bench $\Gamma_1(\theta_{init}, h, \tau) = (E_1(\theta_{init}), \varphi_1(h, \tau))$, combined with the sub-domain $\rho_1 = \{[\theta_{init} \mapsto v_{\theta_{init}}, h \mapsto v_h, \tau \mapsto v_\tau] \text{ s.t. } |v_{\theta_{init}}| < 90 \wedge v_\tau > 60 \wedge v_h \leq 0.4\}$ can reliably capture the failure to climb to a height above 0.4 m within 60 secs for some initial orientation pointing towards the hill.

3.2 Synthesis of Test Bench Constraints

Designing a test bench for a fault involves (i) creating a parametrized test bench, and (ii) finding a sub-domain of the parameters such that it reliably captures the fault. While creating a parametrized test bench by hand is easy, in our experience manually coming up with the sub-domain is tedious. It not only requires the instructor to be a relative expert in STL and run-time verification, but also requires careful observation of traces where the fault is known to be present and not present, and a number of iterations of trial and error. On the other hand, instructors can easily come up with a set of *reference controllers*: a set \mathcal{C}^+ of positive-labeled controllers that are all known to exhibit the faulty behavior, and a set \mathcal{C}^- of negative-labeled controllers that are all known to *not* exhibit the faulty behavior.

We define below the problem of synthesizing a sub-domain from a set \mathcal{C}^+ of positive-labeled controllers and a set \mathcal{C}^- of negative-labeled controllers.

Problem 1 Given the following: (1) a parameterized test $\Gamma(\mathbf{p})$ with a domain \mathfrak{U} for parameters \mathbf{p} , and (2) two sets \mathcal{C}^+ and \mathcal{C}^- of controllers. Synthesize a sub-domain $\rho \subseteq \mathfrak{U}$ s.t. test bench $(\Gamma(\mathbf{p}), \rho)$ does not succeed for any $C \in \mathcal{C}^-$ and succeeds for all $C \in \mathcal{C}^+$.

We can see that any sub-domain that does not intersect with $\Omega(C, \Gamma(\mathbf{p}))$ for any $C \in \mathcal{C}^-$ and has a non-empty intersection with $\Omega(C, \Gamma(\mathbf{p}))$ for every $C \in \mathcal{C}^+$ satisfies the requirements in Problem 1. From amongst all these possibilities, we choose the following (also illustrated in Figure 3.1b)

$$\rho = \bigcup_{C \in \mathcal{C}^+} \Omega(C, \Gamma(\mathbf{p})) \setminus \bigcup_{C \in \mathcal{C}^-} \Omega(C, \Gamma(\mathbf{p})) \quad (3.1)$$

For convenience, we use $\Omega(\mathcal{C}^+, \Gamma(\mathbf{p}))$ (and $\Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$) to refer to $\bigcup_{C \in \mathcal{C}^+} \Omega(C, \Gamma(\mathbf{p}))$ (and $\bigcup_{C \in \mathcal{C}^-} \Omega(C, \Gamma(\mathbf{p}))$). The rationale behind this choice of ρ is two-fold:

1. To increase *coverage* of fault detection for unlabeled controllers, we wish to include as much of $\Omega(\mathcal{C}^+, \Gamma(\mathbf{p}))$ in ρ as possible because every parameter valuation in that set corresponds to a test that passed on some positively-labeled controller, i.e. a controller that exhibits the faulty behavior.
2. For the tests corresponding to valuations that are not in either one of $\Omega(\mathcal{C}^+, \Gamma(\mathbf{p}))$ or $\Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$, we choose a *lenient* grading route and do not include them in ρ . This means that if an unlabeled controller does not pass on any test that lies in $\Omega(\mathcal{C}^+, \Gamma(\mathbf{p}))$, it will not be labeled as one exhibiting the fault. This is how instructors often grade labs in practice, i.e., if tests conclude that a solution may or may not be faulty, it is considered to be non-faulty, pending a more detailed manual review. Here we are also assuming that we have a good range of positive and negative labeled controllers that cover a wide variety of ways in which the fault may or may not be exhibited.

To generate ρ as in Eqn. 3.1, we compute Ω , as discussed next.

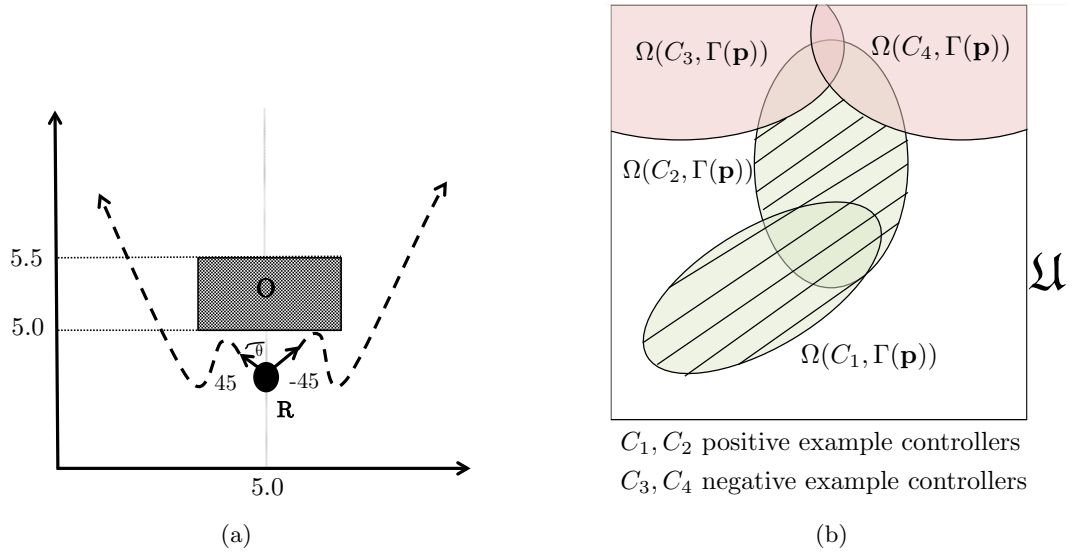


Figure 3.1: (a) Environment E_0 from Examples 1, 2, and 3 with robot R and obstacle O . The two trajectories shown by dotted lines meet the goals for the cases $\theta = 45$ and $\theta = -45$. (b) The hatched region is the sub-domain ρ obtained from satisfaction regions of positive and negative controller examples.

3.3 Computing the Satisfaction Region

Given a controller C and a parametrized test $\Gamma(\mathbf{p})$ with $\mathbf{p} = (p_1, p_2, \dots, p_k)$, we wish to compute $\Omega(C, \Gamma(\mathbf{p}))$. We assume that all parameters are numerical. Every parameter that is not finite valued is discretized by sampling uniformly at some granularity within reasonable lower and upper bounds. By this construction, the domain \mathfrak{U} is now a finite k -dimensional array and can be written as a Cartesian product of finite sets $\mathfrak{U}_1 \times \mathfrak{U}_2 \times \dots \times \mathfrak{U}_k$, where p_i takes values in the set \mathfrak{U}_i . We assume some indexing on each \mathfrak{U}_i such that $\mathfrak{U}[j_1, j_2, \dots, j_k]$ refers to the element of \mathfrak{U} formed by picking the j_i -th element from each \mathfrak{U}_i . Moreover, we assume that this indexing is consistent with the natural order defined over each \mathfrak{U}_i (i.e., a lower index implies a smaller value). Let $N = \max_i(|\mathfrak{U}_i|)$. The size of \mathfrak{U} is $\mathcal{O}(N^k)$. Given this representation of \mathfrak{U} , $\Omega(C, \Gamma(\mathbf{p}))$ can be represented by a k -dimensional bit-array, such that, $\Omega(C, \Gamma(\mathbf{p}))[j_1, j_2, \dots, j_k] = 1$ iff the test $\Gamma(\mathfrak{U}[j_1, j_2, \dots, j_k](\mathbf{p}))$ passes on the test $\Gamma(\mathfrak{U}[j_1, j_2, \dots, j_k](\mathbf{p}))$ passes on C . The most naïve way to compute $\Omega(C, \Gamma(\mathbf{p}))$ is

to perform the test $\Gamma(v(\mathbf{p}))$ for every valuation $v(\mathbf{p}) \in \mathfrak{U}$. We describe a more efficient approach to do this in cases where the test bench is monotonic in one or more parameters.

Definition 9 (*Monotonicity*) Given an order \preceq on a parameter p_i in the parameter vector $\mathbf{p} = (p_1, p_2, \dots, p_k)$, a parameterized test $\Gamma(\mathbf{p})$ is monotonic in p_i if for every controller C

$$\begin{aligned} \forall v, v' \quad v(p_i) \preceq v'(p_i), \forall j \neq i \cdot v'(p_j) = v(p_j) \\ \Gamma(v(\mathbf{p})) \text{ passes for } C \Rightarrow \Gamma(v'(\mathbf{p})) \text{ passes for } C \end{aligned} \quad (3.2)$$

Example 5 Consider the parameterized test $\Gamma_1(\theta_{init}, h, \tau)$ from Example 4. Consider the order \leq over h and two values $v_h \leq v'_h$. For any controller C , if $\Gamma_1(v_{\theta_{init}}, v_h, v_\tau)$ passes, it means that the `pos.z` always stays below v_h for the time interval $[0, v_\tau]$, which implies that it stays below v'_h as well and hence $\Gamma_1(v_{\theta_{init}}, v'_h, v_\tau)$ will pass. Thus $\Gamma_1(\theta_{init}, h, \tau)$ is monotonic in h .

Similarly, for the order \geq on the parameter τ and two values $v_\tau \geq v'_\tau$, if a test $\Gamma_1(v_{\theta_{init}}, v_h, v_\tau)$ passes for any controller, it means that the `pos.z` always stays below v_h for the time interval $[0, v_\tau]$, which implies that the same is true for the time interval $[0, v'_\tau]$ and hence the test $\Gamma_1(v_{\theta_{init}}, v_h, v'_\tau)$ will also pass.

We can extend the definition of monotonicity to sets of parameters by defining required orders on tuples of parameter values. For example, $\Gamma_1(\theta_{init}, h, \tau)$ is monotonic in (h, τ) if we consider \preceq as the order, where $(v_h, v_\tau) \preceq (v'_h, v'_\tau)$ iff $v_h \leq v'_h$ and $v_\tau \geq v'_\tau$. Note that we do not need separate monotonically increasing and decreasing parameterized tests since we can always invert the order on the parameter and keep the definition consistent.

Note that the definition of monotonicity allows a parameterized test to be monotonic in environment parameters but, so far in practice we have never encountered cases when this happens. Checking that a parameterized test is monotonic in certain parameters that only occur in the PSTL part of the test can be done by reduction to satisfiability modulo theories (SMT) as described in more detail by Jin et al. [26]. This is an offline step carried out at the time of design of a parameterized test.

Definition 10 (*Monotone Bit-Array*) For two indices $\mathbf{j} = [j_1, j_2, \dots, j_k]$ and $\mathbf{j}' = [j'_1, j'_2, \dots, j'_k]$ of a k -dimensional bit-array A , we say $\mathbf{j} \leq \mathbf{j}'$ iff $j_1 \leq j'_1, j_2 \leq j'_2, \dots, j_k \leq j'_k$. The array A is said to be monotone if for any indices \mathbf{j} and \mathbf{j}' s.t. $\mathbf{j} \leq \mathbf{j}'$, $A[\mathbf{j}] = 1$ implies that $A[\mathbf{j}'] = 1$.

We now describe how monotonicity proves to be a useful property to efficiently compute $\Omega(C, \Gamma(\mathbf{p}))$. First consider the case when $\Gamma(\mathbf{p})$ is monotonic in all k parameters p_1, p_2, \dots, p_k . Owing to monotonicity, we can index the valuations using their respective orders s.t. for any controller C , the k -dimensional bit-array representation of $\Omega(C, \Gamma(\mathbf{p}))$ is monotone. We describe an algorithm to compute $\Omega(C, \Gamma(\mathbf{p}))$ in three separate cases.

3.3.1 Case: $k=1$

For the single parameter p_1 we can perform a binary search within its domain to determine the index b such that $\Gamma(\mathcal{U}[j_1 = b](\mathbf{p}))$ does not pass on C while $\Gamma(\mathcal{U}[j_1 = b + 1])$ passes. We would have to perform $\mathcal{O}(\log N)$ tests.

3.3.2 Case: $k=2$

For two parameters p_1 and p_2 , say we have the 2-d array of indices $[1 \dots U] \times [1 \dots V]$. We start at the index $\langle row = 1, col = V \rangle$. At each step we perform the test $\Gamma(\mathcal{U}[j_1 = row, j_2 = col](\mathbf{p}))$ on C . If the test passes, we mark the complete column $\Omega(C, \Gamma(\mathbf{p}))[j_1 \geq row, j_2 = col]$ with 1s (we can do this because of monotonicity) and decrement col by 1. If the test does not pass, we mark the complete row $\Omega(C, \Gamma(\mathbf{p}))[j_1 = row, j_2 \leq col]$ with 0s and increment row by 1. We do this until we have covered the whole array. We would have to perform $\mathcal{O}(\max(U, V)) = \mathcal{O}(N)$ tests since we mark out a complete row or column after every test. Figure 3.2a shows an intermediate step in a run of this algorithm.

3.3.3 Case: $k \geq 3$

For more than 2 parameters, we enumerate over all possible valuations of first $k - 2$ parameters and use the case for $k = 2$ for the 2-d sub-array obtained by fixing p_1, p_2, \dots, p_{k-2} . We would have to perform $\mathcal{O}(N^{k-1})$ tests. We cannot hope to do (asymptotically) better than this as it is shown in [27] that searching in a monotone d -dimensional array where each dimension is of size at most n is lower bounded by $c_2(d)n^{d-1}$, where $c_2(d) = \mathcal{O}(d^{-\frac{1}{2}})$ for $d \geq 2$.

For the general case, let $\Gamma(\mathbf{p})$ be non-monotonic in the first $k - d$ parameters and monotonic in the d others. We enumerate over all possibilities of the first $k - d$ parameters and apply the algorithm for monotonic parameters to the d dimensional sub-array obtained by fixing p_1, p_2, \dots, p_{k-d} .

Using the above approach, we can compute $\Omega(\mathcal{C}^+, \Gamma(\mathbf{p}))$, $\Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$ and $\rho = \Omega(\mathcal{C}^+, \Gamma(\mathbf{p})) \setminus \Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$, all represented in the form of k -dimensional bit-arrays.

3.4 Adequate Test Samples for Grading

Checking whether a new controller C succeeds on a test bench $(\Gamma(\mathbf{p}), \rho)$ amounts to searching for a valuation in ρ such that $\Gamma(v(\mathbf{p}))$ passes for C . The naive approach to solve the search problem is to enumerate all valuations in ρ . We describe a more efficient search strategy when $\Gamma(\mathbf{p})$ is monotonic in one or more parameters.

Definition 11 (*Adequate Test Sample*) An adequate test sample $\alpha \subseteq \rho$ is a set of valuations s.t. for any controller C , $(\Gamma(\mathbf{p}), \rho)$ succeeds on C iff there is at least one $v \in \alpha$ for which $\Gamma(v(\mathbf{p}))$ passes for C .

Definition 12 (*Corner*) A corner in a monotone k -dimensional bit-array A is an index $\mathbf{j} = [j_1, j_2, \dots, j_k]$ s.t. $A[\mathbf{j}] = 0$ and $\forall 1 \leq l \leq k$, if the index $[j_1, j_2, \dots, j_l + 1, \dots, j_k]$ lies within the bounds of A , then $A[j_1, j_2, \dots, j_l + 1, \dots, j_k] = 1$.

First consider the case when a parameterized test $\Gamma(\mathbf{p})$ is monotonic in all parameters $\mathbf{p} = (p_1, p_2, \dots, p_k)$. Say we have computed $\Omega(\mathcal{C}^+, \Gamma(\mathbf{p}))$, $\Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$ and $\rho = \Omega(\mathcal{C}^+, \Gamma(\mathbf{p})) \setminus \Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$ in k -dimensional bit-array form.

Proposition 1 *The set α comprising of all valuations $\mathfrak{U}[\mathbf{j}]$ s.t. \mathbf{j} is a corner of $\Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$ and $\Omega(\mathcal{C}^+, \Gamma(\mathbf{p}))[\mathbf{j}] = 1$, is a minimal adequate test sample for $(\Gamma(\mathbf{p}), \rho)$.*

Proof. We first show that α is adequate then we show α is also minimal. For this proof, we refer to $\Omega(\mathcal{C}^+, \Gamma(\mathbf{p}))$ by Ω^+ and $\Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$ by Ω^- .

Assume $\Gamma(v(\mathbf{p}))$ passes for C for some $v \in \alpha$. Let the index of this valuation be \mathbf{j}_v . By definition of α , $\Omega^+[\mathbf{j}_v] = 1$ and \mathbf{j}_v is a corner of Ω^- implying $\Omega^-[\mathbf{j}_v] = 0$. From the way we have defined ρ , we can say that $\rho[\mathbf{j}_v] = 1$ or $v \in \rho$ which means $(\Gamma(\mathbf{p}), \rho)$ succeeds for C . For reverse implication, assume $(\Gamma(\mathbf{p}), \rho)$ succeeds for C , it means that it is possible to find an index $\mathbf{j} = [j_1, j_2, \dots, j_k]$ s.t. $\mathfrak{U}[\mathbf{j}] \in \rho$ (equivalently, $\rho[\mathbf{j}] = 1$) and $\Gamma(v(\mathbf{p}))$ passes for C (equivalently, $\Omega(C, \Gamma(\mathbf{p}))[\mathbf{j}] = 1$). Since $\mathbf{j} \in \rho$, we have $\Omega^+[\mathbf{j}] = 1$ and $\Omega^-[\mathbf{j}] = 0$. If \mathbf{j} is a corner of Ω^- , then we have $\mathfrak{U}[\mathbf{j}] \in \alpha$ and we are done. If not, then there exists $1 \leq l \leq k, \mathbf{j}' = [j_1, j_2, \dots, j_l + 1, \dots, j_k]$ s.t. $\Omega^-[\mathbf{j}'] = 0$. By monotonicity, we also have $\Omega^+[\mathbf{j}'] = 1$ and $\Omega(C, \Gamma(\mathbf{p}))[\mathbf{j}'] = 1$. If \mathbf{j}' is a corner of Ω^- , then $\mathfrak{U}[\mathbf{j}'] \in \alpha$ and we are done. Else we set \mathbf{j} to \mathbf{j}' and proceed again in the same way. Since \mathfrak{U} is finite, this procedure is guaranteed to terminate at a corner of Ω^- .

To show minimality, we remove some arbitrary valuation v from α and show that it becomes inadequate. Say \mathbf{j}_v is the index corresponding to v . Consider a controller C s.t. $\Omega(C, \Gamma(\mathbf{p}))[\mathbf{j}] = 1$ iff $\mathbf{j} \geq \mathbf{j}_v$. Since \mathbf{j}_v is a corner of Ω^- , for every index $\mathbf{j} \neq \mathbf{j}_v$ and $\mathbf{j} \geq \mathbf{j}_v$, we have that $\Omega^-[\mathbf{j}] = 1$. This means there is no corner of Ω^- in $\Omega(C, \Gamma(\mathbf{p}))$ apart from \mathbf{j}_v . Hence, we will not be able to find another $v' \in \alpha, v' \neq v$ s.t. $\Gamma(v'(\mathbf{p}))$ passes on C , even though $(\Gamma(\mathbf{p}), \rho)$ succeeds on C . This means α becomes inadequate if we remove any of its elements, thus making it minimal. ■

Figure 3.2b shows an example of a minimal adequate test sample for the 2-d case. To compute α , similar to Sec. 3.3; in case $k = 1$, we can do a binary search to find the corner; in case $k = 2$, we can find corners by starting at the boundary of the

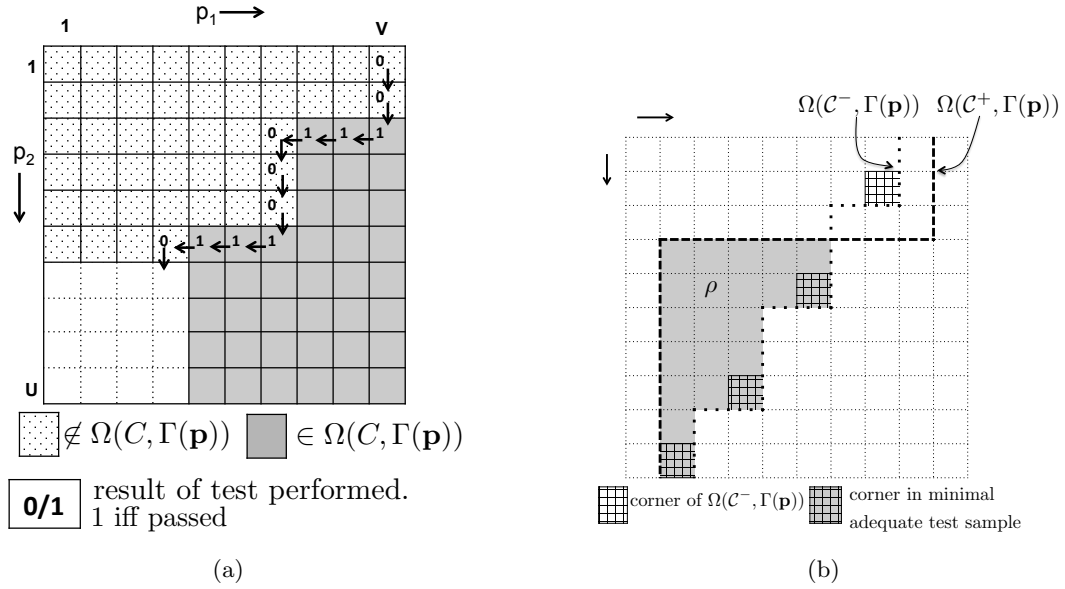


Figure 3.2: (a) An intermediate step in a run of the algorithm used to compute $\Omega(C, \Gamma(\mathbf{p}))$ for two monotonic parameters $\mathbf{p} = (p_1, p_2)$. The arrows indicate the tests that are performed. Monotonicity allows us to compute whole of $\Omega(C, \Gamma(\mathbf{p}))$ by performing $\mathcal{O}(\max(U, V))$ tests. (b) For the case of two monotonic parameters (increasing in the directions shown by arrows), the dashed (and dotted) lines represent the boundary between cells containing 0s and 1s for $\Omega(\mathcal{C}^+, \Gamma(\mathbf{p}))$ (and $\Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$). The shaded part is ρ . The hatched cells are corners of $\Omega(\mathcal{C}^-, \Gamma(\mathbf{p}))$ and the shaded hatched cells comprise the minimal adequate test sample.

2-d array and eliminating rows and columns; and in case $k \geq 3$, we can enumerate over first $k - 2$ parameters and apply the case for $k = 2$ on the rest. For the general case of $k - d$ non-monotonic and d monotonic parameters, we enumerate over all possibilities of first $k - d$ parameters, and keep accumulating the adequate test sample calculated for the d -dimensional monotone sub-array obtained by fixing the first $k - d$ parameters.

We conclude this section with a remark about an alternative mathematical formulation. If we treat a monotone bit-array as a partially-ordered set (poset) \mathfrak{D} , then, the satisfaction region $\Omega(C, \Gamma(\mathbf{p}))$ of some controller C is an upward closed subset of \mathfrak{D} . The sub-domain ρ is now the intersection of an upward-closed (Ω^+) and another downward-closed set ($\mathfrak{U} \setminus \Omega^-$). With some effort, we can show that the minimal adequate test sample α corresponds to the maximal elements of ρ . However, we find the monotone bit-array formulation more useful for our purposes because it is a special case of a poset that allows for efficient algorithms (as given in Sec. 3.3.1 and 3.3.2) for computation of α , which is not obvious with the general poset formulation.

3.5 Related Work

Parameter synthesis for PSTL formulas has been studied before [23], [26]. Unlike our work, these efforts seek to find specific parameter values rather than sub-domains, and are not directly usable in the auto-grading context of this paper. A symbolic approach to PSTL parameter synthesis has been discussed in [23], which reports that an enumerative approach outperforms the symbolic one.

We also note related work in the area of fault localization only using execution traces (black-box localization) [28], [29]. However, these techniques apply to digital systems and are not directly usable in our context of hybrid systems with continuous variables.

Chapter 4

Clustering-Based Active Learning

In this chapter, we give the details of an active learning algorithm based on clustering, that we use to select the set of controllers that should be labeled to serve as the training set. For ease of presentation in this chapter, we make a simplifying assumption that *parameterized tests* do not contain any environment parameters. For parameterized tests that contain environment parameters, we can apply the same technique by enumerating over the environment parameters. The synthesis algorithm described in Sec. 3.3 is used as a black box *training* module TRAIN, which takes as input a parameterized test $\Gamma(\mathbf{p}) = (\varphi(\mathbf{p}), E)$ (again without environment parameters), and two sets of controllers \mathcal{C}^+ and \mathcal{C}^- (positively and negatively labeled *training data*), and gives as output a test bench $(\Gamma(\mathbf{p}), \rho)$. A synthesized test bench $(\Gamma(\mathbf{p}), \rho)$ (also referred to as a *classifier* in this chapter), is then used by a *classification* module CLASSIFY that can label new solutions as being faulty or not. In other words, given a dataset \mathcal{D} of solutions, CLASSIFY will output a partition of \mathcal{D} into two sets \mathcal{D}_0 and \mathcal{D}_1 , of solutions labeled 0 and 1, corresponding to fault being present and not respectively.

Generating labeled data for the training module is expensive. An instructor would have to manually look at the simulation video to determine whether solution is faulty or not. This is the problem we tackle here. How can we make generation of training examples easier and more efficient?

4.1 Iterative Synthesis of Test Benches by Active Learning

Active Learning [30] is a form of machine learning where the learning algorithm is able to interactively query the user to get the correct labels for new data points. Our problem fits well within this definition. We extend the training module with another *selection* module `SELECT` that decides which new controller(s) to get a correct label for. The overall active learning procedure is iterative. Algorithm 1 takes a dataset of solutions, an expert labeling oracle (that generates true labels), and a parameterized test corresponding to a fault as input, and outputs a synthesized test bench. The algorithm works iteratively by using clustering to select the controllers to be added to the training data and using the synthesis procedure described in Sec. 3.3 at each step. The training module first generates a classifier based on some sets of training controllers. Depending on the results of the classifier, the controllers labeled as 0 and 1 are separately clustered by a clustering module `CLUSTER`. Using the clusters formed, the selection module `SELECT` chooses new controllers to get correct labels for. All the selected controllers that were incorrectly labeled by the classifier are now added to the training set and the classifier is trained again. This continues until no fresh training data is added. Details of the clustering algorithm and the selection module are given in the following sections.

4.2 Clustering with Precomputed Distances

`CLUSTER` performs density-based spatial clustering (DBSCAN) on a set of unlabeled controllers. DBSCAN only takes pairwise distances among the data points as input. There is no need to specify a feature vector or the number of cluster apriori. We use multi-dimensional dynamic time warping (DTW) (with point-wise Euclidean distance) as the measure of distance between two controllers for a given environment. More concretely, say for a parameterized test $\Gamma(\mathbf{p}) = (\varphi(\mathbf{p}), E)$, the set of variables that occur in the formula φ is V . Given a controller C , we can obtain a simulation trace $\text{sim}(C, E)$, which is a multi-dimensional timed sequence. Note however, that the classification would only depend on the

Algorithm 1: ITERATIVESYNTHESIS

Input: A dataset of student solutions \mathcal{D} , a true labeling oracle \mathcal{O} , and a parameterized test $\Gamma(\mathbf{p}) = (\varphi(\mathbf{p}), E)$ corresponding to some fault

Output: A classifier $(\Gamma(\mathbf{p}), \rho)$ for \mathcal{D}

```
1  $\mathcal{C}^+, \mathcal{C}^- \leftarrow \emptyset$ 
2 repeat
3    $(\Gamma(\mathbf{p}), \rho) \leftarrow \text{TRAIN}(\Gamma(\mathbf{p}), \mathcal{C}^+, \mathcal{C}^-)$ 
4    $\mathcal{D}_0, \mathcal{D}_1 \leftarrow \text{CLASSIFY}((\Gamma(\mathbf{p}), \rho), \mathcal{D})$ 
5    $\theta_0, \theta_1 \leftarrow \text{CLUSTER}(\mathcal{D}_0), \text{CLUSTER}(\mathcal{D}_1)$ 
6    $\mathcal{R}_0, \mathcal{R}_1 \leftarrow \text{SELECT}(\theta_0), \text{SELECT}(\theta_1)$ 
7    $\mathcal{C}_\Delta^+ = \{C \text{ s.t. } (C \in \mathcal{R}_0 \wedge \mathcal{O}(C) = \text{with\_fault})\}$ 
8    $\mathcal{C}^+ = \mathcal{C}^+ \cup \mathcal{C}_\Delta^+$ 
9    $\mathcal{C}_\Delta^- = \{C \text{ s.t. } (C \in \mathcal{R}_1 \wedge \mathcal{O}(C) = \text{without\_fault})\}$ 
10   $\mathcal{C}^- = \mathcal{C}^- \cup \mathcal{C}_\Delta^-$ 
11 until  $\mathcal{C}_\Delta^+$  or  $\mathcal{C}_\Delta^-$  is empty
12 return  $(\Gamma(\mathbf{p}), \rho)$ 
```

variables V , hence we project out rest of the variables from the simulation trace $\text{sim}(C, E)$. We compute DTW distance on the resulting multi-dimensional timed sequence.

4.3 Selection of Training Data from Clusters

The selection module implements the policy used for selecting data points to be added to the training set. This is done bearing in mind that the training algorithm works well if the training data is balanced in terms of number of positive and negative examples. Training data balancing is a standard technique in machine learning [31]. This is specially important in our context because some faults are rare, and other are very common (in non-final versions of solutions), and hence the occurrence of positive and negative examples is imbalanced.

For initialization of the training set during the first iteration, we cluster all the samples using the clustering module. We then select a randomly chosen sample from each cluster, look up its label and add to the training set. If the number of samples for positive and negative training is skewed, we continue picking more training instances until either a threshold upper bound is reached or we are unable to reduce the skew any further. To reduce the skew, we randomly pick a cluster from which a minority instance was obtained (positive or negative) and sample again hoping to obtain another instance of the minority class thereby reducing the skew.

Once the initialization step is complete, we move on to running the classifier and obtaining predicted labels on the test set. If the accuracy on the test set is not 100%, we try and improve our training set by adding examples of samples which were marked wrongly. In order to achieve that, we re-cluster all the samples (test and training) in each class separately, randomly pick a cluster which has not been already represented in the training set (i.e. the cluster and the training set has no sample in common) and pick a random sample from the same. We do this for both class and add the respective sample to the training set if the predicted label was not same as the actual label. This step is performed in a loop until we are unable to increase the size of the training set or 100% accuracy has been achieved.

4.4 Related Work

DTW has been previously used for classification of temporal sequences of video, audio, and graphics data [32], using an algorithm similar to k-nearest neighbours [33]. Active learning is a popular methodology for cases where obtaining training data is expensive, using strategies like uncertainty sampling, expected model change, expected error reduction, etc. [30] We have not seen past work that applies clustering based strategies for active learning.

Chapter 5

Evaluation

The design and initial experimental evaluation of CPSGrader was done using a collection of solutions implemented by 50 groups of students as part of the laboratory component of the Fall 2013 instance of the EECS 149 class at UC Berkeley.

The code was anonymized and collected automatically using post-build commands so that each group provided a variable number of versions, most of which being intermediate non-final solutions. The lab was organized in two sessions, one focusing on the obstacle avoidance problem, and another focusing on the hill climbing. In this section, we describe the set of test benches that we used to establish diagnostics with respect to each goal. For each test bench, we first manually label a set T of 100 randomly selected student solutions. We select 30 solutions out of the 100 while maintaining balance between the number of positive and negative examples which are input to the synthesis algorithm. To elaborate, if we have more than 15 each of positive and negative examples (say 45 positive and 55 negative) then we select some 15 examples of each type arbitrarily. If either one of positive or negative examples is less than 15 (say 5 positive and 95 negative), then we select all instances of the type of example that is scarce and select the remainder of the 30 from the other type (in the example, we will take 5 positive and 25 negative). This is a standard technique in machine learning done to improve coverage and reduce bias in case a fault is rare [31]. In Sec. 5.1 and 5.2, for each test bench, we describe (1) the fault symptom and the

corresponding PSTL formula, (2) environment and STL parameters, and their monotonic nature, (3) synthesized sub-domain and adequate test sample, and (4) synthesis time per training example.

In Sec 5.3.1, we measure *accuracy* of the grader by comparing labels generated by the auto-grader against another set of manually graded solutions (disjoint from T). We also demonstrate *efficiency* in terms of the average grading time per solution.

The auto-grader designed as described above was used in the MOOC version EECS 149.1x. Since we perform synthesis of test benches based on a training set obtained from the on-campus version of the course, in Sec. 5.3.2 we evaluate the *accuracy* of the grader on a set of student solutions collected from the MOOC to show robustness of the grader on a new data set that might have different kinds of variations. We also study the correlation of overall grades assigned by the auto-grader as compared to grades assigned by an expert manual grader in Sec. 5.4.

In Sec. 5.5, we evaluate the iterative synthesis algorithm Alg. 1 (referred to as ISYN in the rest of the chapter) by comparing it against the technique RANDOM where we randomly choose the training set and show that ISYN can obtain higher overall accuracy, with a smaller size of training set used.

In Sec. 5.6, we propose a semi-automated methodology for identifying new fault scenarios using solutions that do not pass the objectives but also do not exhibit any faults in our library. This methodology is based on clustering of simulation traces.

Experiments are performed using a single core of a 2.3 GHz processor with 8 GB of memory. Since more than one tests share the same environment configuration, we run simulations for all solutions in all the environment configurations as needed for our evaluation in a pre-processing step and store traces to files. Each simulation is run for 60 secs of virtual time with a step size of 5 ms which takes about 10 secs of system time. For each test bench, in Sec. 5.1 and 5.2, we report running times of the synthesis algorithm that computes the sub-domain and the adequate test sample, and in Sec. 5.3.1, we report running times of the auto-grader which checks for existence of a passing test in the adequate test sample.

These running times do not include time required for simulation since we are reading traces from files. When using the auto-grader in loop with the simulator, we need one simulation for every environment in each test bench per solution (the aggregate is lower in practice because more than one test benches share the same environment). All simulations are run using NI Robotics Simulator. STL monitoring is performed using Breach [13]. The synthesis modules and grading software with an extended library of faults is made available at the CPSGrader website [34].

5.1 Obstacle Avoidance

In assessing faults in obstacle avoidance, we use an environment $E_3(\theta_{init})$ which contains an obstacle occupying the region $[4.5, 5.5] \times [5.0, 5.5]$. Initial position of the robot is $(5.0, 4.9)$. The parameter θ_{init} encodes the initial orientation of the robot.

Failing simple obstacle avoidance (`avoid_front`)

This test bench checks whether the robot can get past the obstacle when started with the initial orientation $\theta_{init} = 0$, facing the obstacle directly.

- Parameterized Test: $(E_3(0), \varphi_{orient})$ with $\varphi_{orient} = \square_{[0, \tau]}(\text{pos.y} < y_{min})$. If φ_{orient} is satisfied for suitable values of τ and y_{min} , it indicates failure to avoid the obstacle.
- Parameters: (τ, y_{min})
- Domain:¹ $(\tau, y_{min}) \in \{60 : -5 : 10\} \times \{3.0 : 0.1 : 7.0\}$
- Monotonicity: τ monotonic for \geq and y_{min} monotonic for \leq .
- Synthesized sub-domain: See Figure 5.1a
- Adequate Test Sample: $\{(60, 5.7), (55, 4.9), (50, 4.6)\}$
- Average synthesis time per training example: 1.9 sec

¹The notation $\{a : d : b\}$ denotes the set $\{a, a + d, a + 2d, \dots, a + kd\}$, where k is the greatest integer s.t. if $d \geq 0$ then $a + kd \leq b$ else if $d < 0$ then $a + kd \geq b$

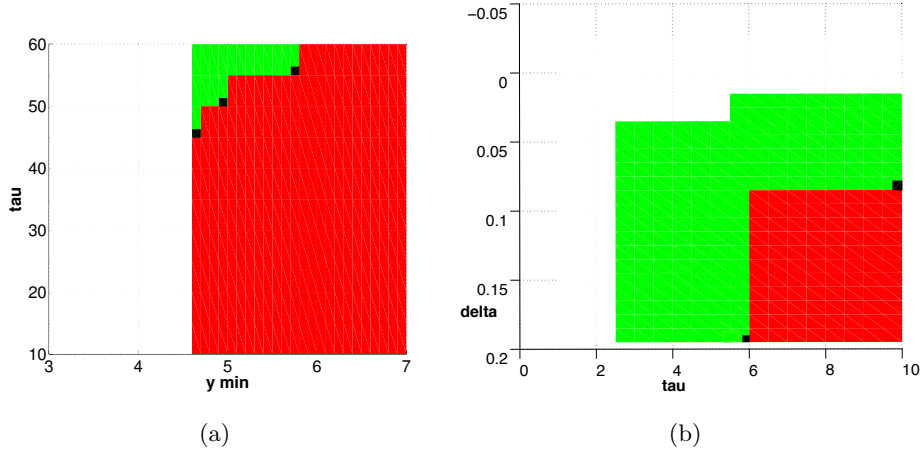


Figure 5.1: (a) Test bench `avoid_front`. Green (lightly shaded) region is the computed sub-domain. Red (dark shaded) region is the set of tests excluded from the sub-domain because they are triggered on at least one negative example. White (unshaded) region is the set of tests that are not triggered on any negative or positive example. Little black squares are the points in the adequate test sample. (b) Test bench `circle`.

Failing re-orienting after obstacle avoidance (`avoid_left/avoid_right`)

This test bench checks whether the robot can get past the obstacle and keep heading in the initial heading direction. We perform the test in two possible initial orientations; facing left ($\theta_{init} = 45$) or right ($\theta_{init} = -45$). We show details for the case $\theta_{init} = 45$.

- Parameterized Test: $(E_3(45), \varphi_{reorient})$ with $\varphi_{reorient} = \square_{[0,\tau]}(\text{pos}.y < y_{min} \vee \text{pos}.x > x_{max})$. If $\varphi_{reorient}$ is satisfied for suitable values of τ , x_{max} and y_{min} , it indicates either failure to avoid the obstacle or failure to re-orient in the correct heading direction.
- Parameters: (τ, y_{min}, x_{max})
- Domain: $(\tau, y_{min}, x_{max}) \in \{60 : -5 : 10\} \times \{3.0 : 0.1 : 7.0\} \times \{6.0 : -0.1 : 3.0\}$
- Monotonicity: τ monotonic for \geq ; y_{min} monotonic for \leq and x_{max} monotonic for \geq .
- Synthesized sub-domain: Due to more than 2 parameters, it is not possible to show it in a figure.
- Adequate Test Sample: $\{(60, 5.4, 4.2), (55, 5.4, 5.0), (50, 4.8, 5.8), (10, 4.4, 5.8)\}$

- Average synthesis time per training example: 26.2 sec

Strict equality check (circle)

This test bench investigates the circle fault mentioned in Section 2.3. The purpose of the test is to detect that at some time instant t_0 , the robot bumps into the obstacle, then turns about itself with a maximum period of τ , while remaining close to its position at t_0 with a margin of δ .

- Parameterized Test: $(E_3(0), \varphi_{\text{circle}})$

$$\varphi_{\text{circle}}(t_0, \delta, \tau) = \diamond(\varphi_{\text{bump}}(t_0) \wedge \diamond_{[0, 2\tau]}(\varphi_{\text{fullturn}}(t_0, \delta)))$$

Where $\varphi_{\text{bump}}(t_0) = \text{bump}(t_0) \equiv \text{TRUE}$ and $\varphi_{\text{fullturn}}$ is given by $\varphi_{\text{fullturn}}(t_0, \delta, \tau) = (\varphi_{\theta \sim 0} \wedge \varphi_{\text{close}}(t_0, \delta) \mathbf{U}_{[0, \tau]}(\varphi_{\theta \sim 180} \wedge \varphi_{\text{close}}(t_0, \delta) \mathbf{U}_{[0, \tau]} \varphi_{\theta \sim 0}))$ where $\varphi_{\text{close}}(t_0, \delta) = \text{dist}(\text{pos}(t_0), \text{pos}) < \delta$ for some distance function dist and $\varphi_{\theta \sim 0}$ and $\varphi_{\theta \sim 180}$ assess that `angle` is close to 0 degrees and 180 degrees, respectively. The suitable value for the parameter t_0 can be determined by the first collision instant with the obstacle, which is common to all solutions since they all start moving forward in the same direction (say this common value is t_0). We fix t_0 to t_0 .

- Parameters: (τ, δ)
- Domain: $(\tau, \delta) \in \{1 : 1 : 10\} \times \{-0.025 : 0.01 : 0.2\}$
- Monotonicity: τ monotonic for \leq and δ monotonic for \leq
- Synthesized sub-domain: See Figure 5.1b
- Adequate Test Sample: $\{(5.5, 0.195), (10.0, 0.075)\}$
- Average synthesis time per training example: 2.7 sec

5.2 Hill Climbing

To assess faults in the hill climbing part of the assignment, we use an environment $E_4(\beta)$ which contains a hill. The parameter β encodes the initial configuration of the robot. It can take two values B and M . In B the robot starts at the bottom of the hill facing 45

degrees rightwards of uphill and in M the robot starts on the hill (midway between bottom and top) facing downhill.

Failing simple hill climb (`hill_climb`)

This test bench checks whether the robot fails to reach near the top of the hill. We perform this test for both possible values of β .

- Parameterized Test: $(E_4, \varphi_{\text{hill}})$ with $\varphi_{\text{hill}} = \square_{[0,\tau]}(\text{pos}.z \leq h)$. If φ_{hill} is satisfied for suitable values of τ and h , it indicates failure to reach near top of the hill.
- Parameters: (β, τ, h)
- Domain: $(\beta, \tau, h) \in \{B, M\} \times \{60 : -5 : 10\} \times \{-0.1 : 0.01 : 0.7\}$
- Monotonicity: τ monotonic for \geq and h monotonic for \leq
- Synthesized sub-domain: See Figure 5.2
- Adequate Test Sample: $\{(M, 55, 0.41), (M, 50, 0.37), (M, 35, 0.35), (M, 15, 0.33), (M, 10, 0.31), (B, 55, 0.45), (B, 50, 0.34), (B, 45, 0.18), (B, 40, 0.07)\}$
- Average synthesis time per training example: 6.2 sec

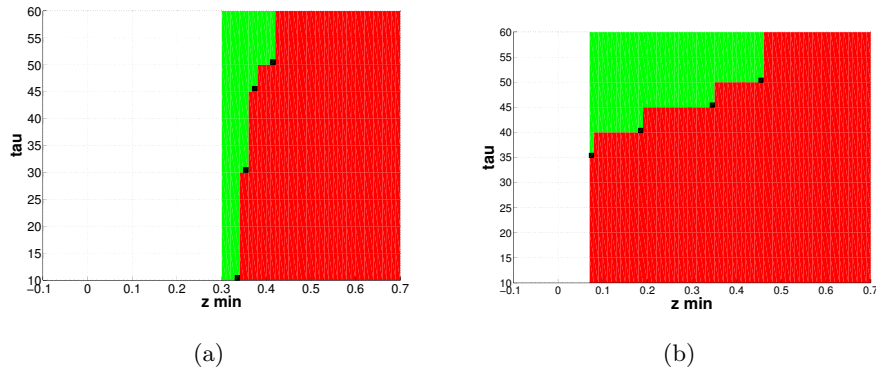


Figure 5.2: (a) Test bench `hill_climb` ($\beta = M$) (b) Test bench `hill_climb` ($\beta = B$)

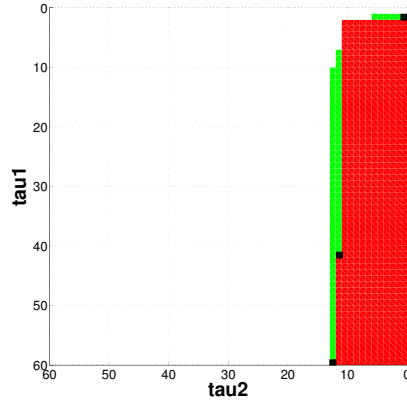


Figure 5.3: Test bench what_hill

Failure to detect hill (**what_hill**)

This test bench checks the failure of robot to detect when it is on a hill. This is a specific bug which leads to failure in hill climbing. We use the environment E_4 with $\beta = B$.

- Parameterized Test: $(E_4(B), \varphi_{\text{hilldet}})$ with $\varphi_{\text{hilldet}} = \diamond_{[0, \tau_1]}(\varphi_{\text{fwd}} \mathbf{U}_{[\tau_2, +\infty]} \varphi_{\text{cliff}})$, where φ_{fwd} assesses that the robot is moving forward and φ_{cliff} assess firing of cliff sensor. If this property is satisfied for suitable values of τ_1 and τ_2 , it means that the robot keeps driving straight until it hits a cliff even if it is on a hill instead of re-orienting towards uphill direction.
- Parameters: (τ_1, τ_2)
- Domain: $(\tau_1, \tau_2) \in \{0 : 1 : 60\} \times \{60 : -1 : 0\}$
- Monotonicity: τ_1 monotonic for \geq and τ_2 monotonic for \leq
- Synthesized sub-domain: See Figure 5.3
- Adequate Test Sample: $\{(1, 0), (41, 12), (60, 13)\}$
- Average synthesis time per training example: 8.3 sec

No filtering (filter)

This test bench checks whether the reason for a failure to climb a hill is the absence of a low-pass filter applied to the accelerometer data to smoothen it. We check this by performing the test `hill.climb` with E_4 but applying a low-pass filter to the accelerometer data externally (before it is fed into the controller). If the robot is able to climb the hill with an external filter but fails to do so without it, we can conclude that absence of the filter is the bug.

5.3 Accuracy of Classification

To measure accuracy we use the synthesized test benches to label a set of student solutions (disjoint from the training set) and compare the labels assigned by the auto-grader to manually assigned labels. We evaluate on the set of solutions collected from both the on-campus offering of the course as well as the MOOC version.

5.3.1 On-campus EECS 149

Table 5.1 shows obtained accuracy results and average running times for 8 test benches. The running times do not include time needed for simulation. For each solution, simulation in a total of 6 environment configurations is collectively needed for the 8 test benches (2 environments are shared). Note that we find a majority of solutions that are not able to meet goals but that is expected because our solution set has preliminary and intermediate versions of the solutions as well. We also find that accuracy is poorer in the hill climbing cases, which shows that variation in student solutions is higher in that part of the lab.

5.3.2 edX MOOC EECS 149.1x

Table 5.2 shows obtained accuracy results by running CPSGrader on student solutions collected from the MOOC. Here we find that a majority of solutions meet the goals and this is because most solutions are collected from the final assignment submissions. The overall

| Test Bench | N^+ | N^{++} | N^- | N^{--} | T_{avg} |
|----------------------------|-------|----------|-------|----------|-----------|
| avoid_front | 74 | 74 | 27 | 27 | 0.119 |
| avoid_left | 78 | 78 | 23 | 23 | 0.158 |
| avoid_right | 82 | 82 | 19 | 19 | 0.148 |
| circle | 2 | 2 | 99 | 99 | 0.382 |
| hill_climb ($\beta = B$) | 49 | 36 | 345 | 345 | 0.111 |
| hill_climb ($\beta = M$) | 35 | 32 | 359 | 359 | 0.120 |
| what_hill | 220 | 216 | 174 | 156 | 0.288 |
| filter | 8 | 7 | 354 | 339 | 0.412 |

Table 5.1: N^+ is the number of solutions with fault (manually labeled). N^{++} is the number of solutions that the auto-grader correctly labeled as faulty. N^- and N^{--} are defined similarly for solutions without fault. T_{avg} is the average labeling time per solution in seconds.

accuracy is poorer as compared to the on-campus dataset but that is expected because all test benches are synthesized using reference solutions chosen from within the on-campus data set. The test bench `filter` is excluded from this evaluation because accelerometer filtering was added as a default in the simulator for the MOOC offering.

| Test Bench | N^+ | N^{++} | N^- | N^{--} |
|----------------------------|-------|----------|-------|----------|
| avoid_front | 189 | 181 | 1018 | 1014 |
| avoid_left | 172 | 167 | 1035 | 1035 |
| avoid_right | 172 | 169 | 1035 | 960 |
| circle | 10 | 10 | 1197 | 1196 |
| hill_climb ($\beta = B$) | 360 | 304 | 234 | 230 |
| hill_climb ($\beta = M$) | 236 | 175 | 358 | 346 |
| what_hill | 314 | 312 | 280 | 194 |

Table 5.2: Notation is same as in Table 5.1

5.4 Grade Correlation

We study how the overall grades assigned by an expert are related to the grades assigned by CPSGrader on the MOOC data. Overall grades are calculated based on how many assignment goals (`avoid_front`, `avoid_right`, `avoid_left`, `hill_climb`) the solution meets and does not depend on the presence/absence of specific faults (`circle`, `what_hill`, `filter`). The faults are

only meant for feedback and debugging support. We assign 1 point for each goal met, thus grading on a scale of 0 to 5. Table 5.3 notes the number of solutions that achieved each grade bar. The correlation coefficient of expert grades v/s CPSGrader assigned grades is found to be 0.87. These results show that CPSGrader assigns grades that are highly correlated with expert grades. The grade distribution appears to be slightly skewed towards lower grades for CPSGrader as compared to expert grades. This is against the intuition that the test benches in CPSGrader are designed to be lenient and needs further investigation.

| Grade Bar | Expert | CPSGrader |
|------------------|--------|-----------|
| 0 | 14 | 11 |
| 1 | 11 | 12 |
| 2 | 11 | 32 |
| 3 | 182 | 216 |
| 4 | 185 | 209 |
| 5 | 192 | 115 |

Table 5.3: Number of solutions at each grader bar for the Expert grader and CPSGrader.

5.5 Effectiveness of Iterative Synthesis

To evaluate the active learning technique developed in Sec. 4.1, we compare our technique ISYN against the technique RANDOM where we choose our training set uniformly at random from the complete dataset. We evaluate the two techniques based on overall accuracy achieved, the size of training set used, and the balance of training labels. For each fault, we train the test bench using both ISYN and RANDOM and then test the accuracy of the obtained test bench on a disjoint set of solutions. To simplify the comparison, we set the upper bound on the number of training instances used in ISYN and total number of randomly chosen samples in RANDOM as 30. In some cases, ISYN may terminate with less than 30 examples in the training set if the clustering algorithm is not able to find enough number of clusters. To compare accuracy we note the True Positive Rate (TPR), True Negative Rate (TNR), and F-score for both techniques. The F-score is specifically insightful for our current application auto-grading, because the classifier is inherently lenient and a better classifier would be the one that can identify at least a few cases on existence of faults

in the solutions. Analysis results for the 7 distinct faults are shown in Table 5.4. From the table, it can be seen that F-score is individually higher in case of ISYN than RANDOM for all the faults except `avoid_left`, thus leading to the conclusion that ISYN leads to better accuracy of classification than RANDOM for equal or lesser size of training set. It is difficult to diagnose the reason for the `avoid_left` exception because the algorithm depends on the fine tuning of many different parameters of both DTW and DBSCAN. For the fault `avoid_right`, we see that ISYN performs significantly better than RANDOM for positive examples but worse for negative examples. In this case the reason is that ISYN ends up selecting only 13 (30 being the upper bound) training examples because CLUSTER cannot find enough number of clusters even for a wide range of configuration settings.

As we noted before, training data balancing is important for the training algorithm to work well. In order to evaluate how well-balanced are the training sets obtained using the two techniques, we use *balancing ratio* i.e. the ratio of number of negative training examples and number of positive training examples. The closer this value is to 1, the better balanced are the two training sets. Table 5.4 gives a clear break down of the number of positive and negative training examples used for each fault per technique. In our evaluation, we find that this ratio was ~ 4.3 for RANDOM while it was ~ 1.2 for ISYN, ISYN leads to more balanced training sets.

Since ISYN on an average performs better than RANDOM on both accuracy measure and balancing measure, we believe that ISYN is a better choice for creating smaller yet more effective training set than random sampling.

5.6 Investigating Unknown Faults Using Clustering

CPSGrader works with a fixed pre-defined library of faults and associated test benches. This raises a natural question. How do we handle the presence of a fault that does not exist in the library yet? In other words, how do we extend this library in a data driven fashion? We attempt to answer these questions partly via semi-automated investigation of the solutions that do not meet the goals of the assignment and also do not exhibit any

| Test Bench | Training Set Size | | TPR | | TNR | | F-score | |
|---------------------------|-------------------|----------------|---------|---------|---------|---------|---------|------|
| | RANDOM | ISYN | RANDOM | ISYN | RANDOM | ISYN | RANDOM | ISYN |
| avoid_front | 23/133 + 7/74 | 15/133 + 15/74 | 133/133 | 133/133 | 67/74 | 70/74 | 0.97 | 0.99 |
| avoid_left | 23/164 + 7/45 | 15/164 + 15/45 | 164/164 | 164/164 | 42/45 | 39/45 | 0.99 | 0.98 |
| circle | 1/7 + 29/200 | 3/7 + 11/200 | 0/7 | 6/7 | 200/200 | 193/200 | 0.00 | 0.60 |
| hill_climb($\beta = B$) | 26/427 + 4/63 | 14/427 + 16/63 | 427/427 | 427/427 | 55/63 | 60/63 | 0.99 | 1.00 |
| hill_climb($\beta = M$) | 28/442 + 2/48 | 29/442 + 1/48 | 442/442 | 442/442 | 11/48 | 18/48 | 0.96 | 0.97 |
| avoid_right | 24/169 + 6/40 | 10/169 + 3/40 | 70/169 | 169/169 | 40/40 | 26/40 | 0.59 | 0.96 |

Table 5.4: Comparison of ISYN and RANDOM. Training Set Size denotes the (number of positive examples selected in the training set)/(total number of positive examples in data set) + (number of negative examples selected in the training set)/(total number of negative examples in the data set). TPR is the true positive rate of the trained classifier. TNR is the true negative rate.

faults in the existing library. We perform this analysis separately for obstacle avoidance and hill climbing objectives. For both the objectives, we first isolate the set of solutions that do not meet goals of the assignment (obstacle avoidance - avoid_right, avoid_front, avoid_left; hill climbing - hill_climb) and also do not exhibit any faults existing in the library (circle, what_hill, filter). Then we cluster the simulation traces of this set of solutions (in some simple default environment that tests the objective) using DBSCAN over pairwise DTW distances as described previously for active learning. We then do manual analysis of the clusters found by looking at similarities between the simulation traces found within a cluster and also the source code of the controllers. This leads to several interesting findings which we describe next. This analysis was carried out using the data from on-campus offering.

For the obstacle avoidance objective, we isolated a total of 114 solutions with unknown faults. DBSCAN forms 4 clusters of size 85, 17, 5, and 5 (2 points were identified as noise.) Investigation of how the minority clusters (17, 5, 5) differed from the majority one leads to two interesting findings: (1) *Symptom*: After hitting the obstacle once, the robot drives away in a direction 90 degrees rightwards of the initial orientation and rams into the wall on the right. *Possible Defect*: Presence of an extraneous unguarded transition that switches from re-orient to the drive mode; (2) *Symptom*: After hitting the obstacle first, and then hitting the wall on the right, the robot drives away in a direction 180 degrees from the initial orientation. *Possible Defect*: Improper use of the angle sensor while checking for

re-orientation success. The `angle` reads between -180 and 180 and hence absolute values should be used when comparing differences. For e.g., a guard that check for `angle < 0` will become true both when `angle` crosses from 1 to -1 and 179 to -179 degrees. Both these symptoms are easy to characterize as STL formulae.

For the hill climbing objective, we isolated a total of 174 solutions with unknown faults. DBSCAN forms 3 clusters of size 159, 4, and 4 (7 points were identified as noise.) The traces in the minority clusters are hard to distinguish from the traces in the majority cluster, hence we do not have interesting findings for this case.

5.7 Discussion

The experimental evaluation indicates that CPSGrader is both accurate and efficient. The test benches used in our evaluation capture common mistakes made by students, as observed in an on-campus offering, and even simply identifying these mistakes can be valuable feedback.

In a course survey filled by students of the edX MOOC EECS 149.1x after completion of the course, 86% of the students reported the feedback generated by the auto-grader critical in helping them debug and solve the lab exercises. The lab also featured an optional hardware track. Among the students who chose to work on hardware, more than 90% reported that their solutions that were developed on the virtual lab (equipped with CPSGrader) worked on the hardware with no or minor modifications.

The parameter synthesis requires a set of “good” and “bad” solutions. We show that a small number of labeled examples (30) is enough to get reasonable accuracy in two different scenarios. However, generation of labeled examples with good coverage of possible variations in students solutions requires an instructor to view the simulation video and label a reasonably large number of student solutions until all major variations are covered. We show that this process can be made easier with our clustering-based iterative synthesis approach, achieving better accuracy with fewer number of training examples.

Chapter 6

Conclusion and Future Work

In this thesis, we have formalized the auto-grading problem for laboratory assignments in cyber-physical systems, and presented a formal, algorithmic approach to solve it based on parameter synthesis. The approach is general and can apply beyond the particular motivating lab setting considered here. The theoretical treatment makes no assumptions about the form of the controller, environment, and simulation model. Note also that our approach can be used with any black-box simulator. We also designed and evaluated a clustering-based active learning technique for selection of labeled training examples for the synthesis algorithm. Again, this clustering-based active learning approach is general and can apply to any setting involving learning from time-series data.

There are several interesting directions for future work. One direction is to introduce cost or reward metrics into the model to quantify the quality of a student solution. Monitoring these metrics over a set of tests can help assign partial credit or extra credit to student solutions. For example, in a problem involving robot navigation to a goal location, a controller that gets closer, or takes less time, should intuitively receive more credit than one that does not. Another direction is to develop STL mining based methods for synthesizing the form of the STL formulae in test benches. More interesting would be to extend the work on identifying unknown faults by developing a general approach to synthesize test benches directly from unlabeled examples of student solutions in an unsupervised way.

As mentioned, the auto-grader has already been successfully deployed in an actual MOOC, EECS149.1x [7], and we have run user studies on its effectiveness. In a course survey filled by students of the edX MOOC EECS 149.1x after completion of the course, 86% of the students reported the feedback generated by the auto-grader critical in helping them debug and solve the lab exercises. The lab also featured an optional hardware track. Among the students who chose to work on hardware, more than 90% reported that their solutions that were developed on the virtual lab (equipped with CPSGrader) worked on the hardware with no or minor modifications.

We are exploring many avenues to use CPSGrader in other classes and labs. One interesting topic is analog and mixed signal circuits, for which *Time Frequency Logic* (TFL [35]) could be used instead of STL.

Finally, beyond the application to education, we note that our technique can be applied to debugging problems for embedded controllers where we can assume a plausible fault model and where monotonicity holds; e.g., for industrial control systems where monotonicity of PSTL has already been found widespread [26].

References

- [1] *The Year of the MOOC*, New York Times, November 2012. [Online]. Available: <http://www.nytimes.com/2012/11/04/education/edlife/massive-open-online-courses-are-multiplying-at-a-rapid-pace.html>
- [2] E. A. Lee and S. A. Seshia. EECS 149 course website. University of California, Berkeley. [Online]. Available: <http://chess.eecs.berkeley.edu/eecs149>
- [3] —, *Introduction to Embedded Systems - A Cyber-Physical Systems Approach*. Berkeley, CA: LeeSeshia.org, 2011. [Online]. Available: <http://LeeSeshia.org>
- [4] J. C. Jensen, E. A. Lee, and S. A. Seshia, *An Introductory Lab in Embedded and Cyber-Physical Systems*. Berkeley, CA: LeeSeshia.org, 2012. [Online]. Available: <http://LeeSeshia.org/lab>
- [5] Massachusetts Institute of Technology (MIT), “The iLab Project,” Last accessed: February 2014. [Online]. Available: <https://wikis.mit.edu/confluence/display/ILAB2/Home>
- [6] P. Mitros, K. Afridi, G. Sussman, C. Terman, J. White, L. Fischer, and A. Agarwal, “Teaching Electronic Circuits Online: Lessons from MITx’s 6.002x on edX,” in *International Symposium on Circuits and Systems (ISCAS)*. Beijing, China: IEEE, May 2013.
- [7] E. A. Lee, S. A. Seshia, and J. C. Jensen. Eecs149.1x cyber-physical systems. <https://www.edx.org/course/uc-berkeleyx/uc-berkeleyx-eecs149-1x-cyber-physical-1629>. UC Berkeley.
- [8] J. C. Jensen, D. H. Chang, and E. A. Lee, “A model-based design methodology for cyber-physical systems,” in *First IEEE Workshop on Design, Modeling, and Evaluation of Cyber-Physical Systems (CyPhy)*, Istanbul, Turkey, 2011. [Online]. Available: <http://chess.eecs.berkeley.edu/pubs/837.html>
- [9] J. C. Jensen, “Elements of model-based design,” Master’s thesis, University of California, Berkeley, February 2010. [Online]. Available: <http://www.eecs.berkeley.edu/Pubs/TechRpts/2010/EECS-2010-19.html>
- [10] R. Smith. Open dynamics engine. [Online]. Available: <http://ode.org>
- [11] J. C. Jensen, E. A. Lee, and S. A. Seshia, “Virtualizing cyber-physical systems: Bringing CPS to online education,” in *Proc. First Workshop on CPS Education (CPS-Ed)*, April 2013.
- [12] O. Maler and D. Nickovic, “Monitoring temporal properties of continuous signals,” in *FORMATS/FTRTFT*, 2004, pp. 152–166.
- [13] A. Donzé, “Breach: A Toolbox for Verification and Parameter Synthesis of Hybrid Systems,” in *Computer-Aided Verification*, 2010, pp. 167–170.
- [14] T. Giorgino, “Computing and visualizing dynamic time warping alignments in R: The DTW package,” *Journal of Statistical Software*, vol. 31, no. 7, pp. 1–24, 8 2009. [Online]. Available: <http://www.jstatsoft.org/v31/i07>
- [15] G. Juniwal, A. Donzé, J. C. Jensen, and S. A. Seshia, “CPSGrader: Synthesizing temporal logic testers for auto-grading an embedded systems laboratory,” in *Proceedings of the 14th International Conference on Embedded Software (EMSOFT)*, October 2014.
- [16] R. Singh, S. Gulwani, and S. Rajamani, “Automatically generating algebra problems,” in *Intl. Conf. of the Association for the Advancement of Artificial Intelligence (AAAI)*, 2012.
- [17] D. Sadigh, S. A. Seshia, and M. Gupta, “Automating exercise generation: A step towards meeting the MOOC challenge for embedded systems,” in *Workshop on Embedded Systems Education (in conjunction with ESWeek)*, Tampere, Finland, October 2012.
- [18] R. Singh, S. Gulwani, and A. Solar-Lezama, “Automated feedback generation for introductory programming assignments,” in *Programming Languages Design and Implementation (PLDI)*, 2013.
- [19] R. Alur, L. D’Antoni, S. Gulwani, D. Kini, and M. Viswanathan, “Automated grading of DFA constructions,” in *Proceedings of the 23rd International Joint Conference on Artificial Intelligence (IJCAI)*, August 2013.
- [20] A. Pnueli, “The temporal logic of programs,” in *Symposium on Foundations of Computer Science*, 1977, pp. 46–57.
- [21] R. Alur and T. A. Henzinger, “A really temporal logic,” in *Symposium on Foundations of Computer Science*, 1989, pp. 164–169.
- [22] R. Koymans, “Specifying real-time properties with metric temporal logic,” *Real-Time Syst.*, vol. 2, no. 4, pp. 255–299, 1990.
- [23] E. Asarin, A. Donzé, O. Maler, and D. Nickovic, “Parametric identification of temporal properties,” in *Runtime Verification*, 2011, pp. 147–160.

- [24] X. Wang, A. Mueen, H. Ding, G. Trajcevski, P. Scheuermann, and E. Keogh, “Experimental comparison of representation methods and distance measures for time series data,” *Data Mining and Knowledge Discovery*, vol. 26, no. 2, pp. 275–309, 2013. [Online]. Available: <http://dx.doi.org/10.1007/s10618-012-0250-5>
- [25] J. Sander, M. Ester, H.-P. Kriegel, and X. Xu, “Density-based clustering in spatial databases: The algorithm gdbscan and its applications,” *Data Mining and Knowledge Discovery*, vol. 2, no. 2, pp. 169–194, 1998. [Online]. Available: <http://dx.doi.org/10.1023/A%3A1009745219419>
- [26] X. Jin, A. Donzé, J. Deshmukh, and S. A. Seshia, “Mining requirements from closed-loop control models,” in *Hybrid Systems: Computation and Control (HSCC)*, 2013.
- [27] N. Linial and M. E. Saks, “Searching ordered structures,” *J. Algorithms*, vol. 6, no. 1, pp. 86–103, 1985.
- [28] W. Li, A. Forin, and S. A. Seshia, “Scalable specification mining for verification and diagnosis,” in *Design Automation Conference*, 2010, p. 755760.
- [29] W. Li and S. A. Seshia, “Sparse coding for specification mining and error localization,” in *Proceedings of the International Conference on Runtime Verification (RV)*, September 2012.
- [30] B. Settles, “Active learning literature survey,” *University of Wisconsin, Madison*, vol. 52, pp. 55–66, 2010.
- [31] G. E. A. P. A. Batista, R. C. Prati, and M. C. Monard, “A study of the behavior of several methods for balancing machine learning training data,” *SIGKDD Explorations*, vol. 6, no. 1, pp. 20–29, 2004.
- [32] C. A. Ratanamahatana and E. Keogh, “Making time-series classification more accurate using learned constraints.” SIAM, 2004.
- [33] T. Mitsa, *Temporal Data Mining (Chapter 3)*, 1st ed. Chapman & Hall/CRC, 2010.
- [34] A. Donzé, G. Juniwal, J. C. Jensen, and S. A. Seshia. CPSGrader website. <http://www.cpsgrader.org>. UC Berkeley.
- [35] A. Donzé, O. Maler, E. Bartocci, D. Nickovic, R. Grosu, and S. A. Smolka, “On temporal logic and signal processing,” in *ATVA*, ser. Lecture Notes in Computer Science, S. Chakraborty and M. Mukund, Eds., vol. 7561. Springer, 2012, pp. 92–106.

Appendix A

STL Semantics

The formal semantics of signal temporal logic (STL) are given as follows:

Definition 13 *The satisfaction of an STL formula relative to a signal \mathbf{x} at time t is defined inductively as*

$$\begin{array}{ll} (\mathbf{x}, t) \models \mu & \text{iff } \mathbf{x} \text{ satisfies } \mu \text{ at time } t \\ (\mathbf{x}, t) \models \neg\varphi & \text{iff } (\mathbf{x}, t) \not\models \varphi \\ (\mathbf{x}, t) \models \varphi_1 \wedge \varphi_2 & \text{iff } (\mathbf{x}, t) \models \varphi_1 \text{ and } (\mathbf{x}, t) \models \varphi_2 \\ (\mathbf{x}, t) \models \varphi_1 \mathbf{U}_{[a,b]} \varphi_2 & \text{iff } \exists t' \in [t+a, t+b] \text{ s.t.} \\ & (\mathbf{x}, t') \models \varphi_2 \text{ and} \\ & \forall t'' \in [t+a, t'), (\mathbf{x}, t'') \models \varphi_1 \end{array}$$

Extension of the above semantics to other kinds of intervals (open, open-closed, and closed-open) is straightforward. We write $\mathbf{x} \models \varphi$ as a shorthand of $(\mathbf{x}, 0) \models \varphi$.