# SWEETPEA: A LANGUAGE FOR EXPERIMENTAL DESIGN

by

Anastasia Cherkaev

A thesis submitted to the faculty of The University of Utah in partial fulfillment of the requirements for the degree of

Master of Science

in

Computer Science

Department of Computer Science

The University of Utah

October 2018

Copyright © Anastasia Cherkaev 2018 All Rights Reserved

#### The University of Utah Graduate School

#### STATEMENT OF DISSERTATION APPROVAL

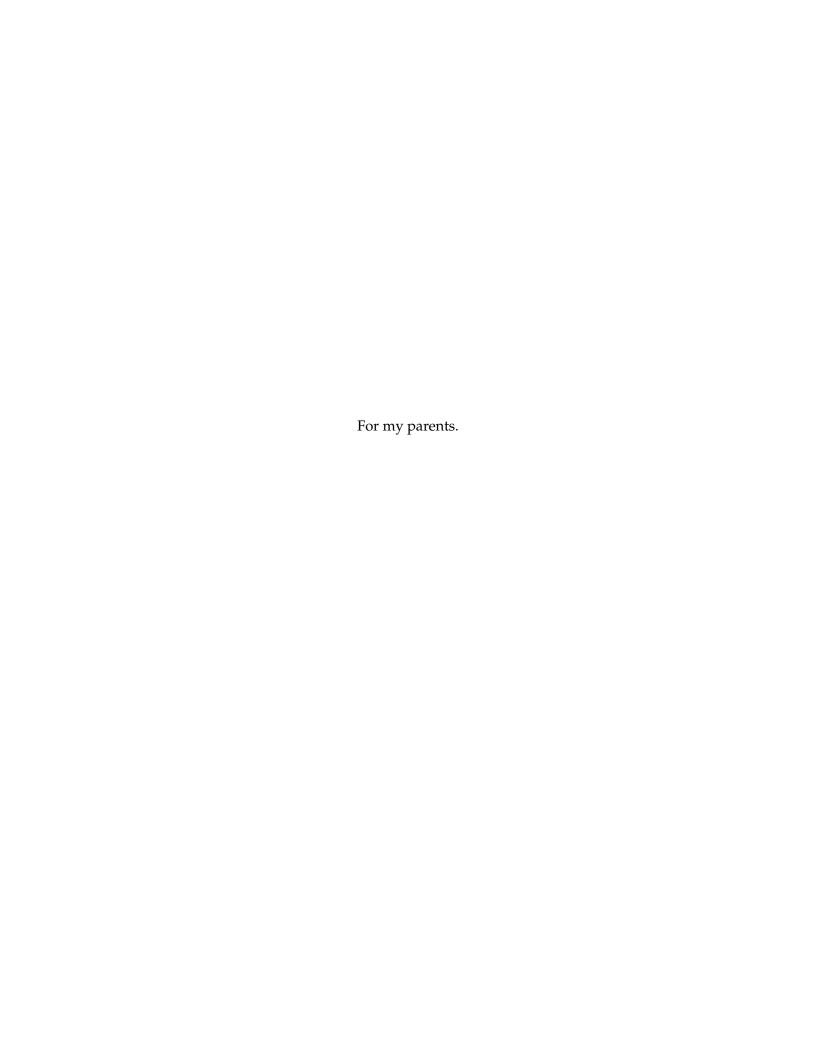
The dissertation of <u>Anastasia Cherkaev</u>
has been approved by the following supervisory committee members:

Vivek Srikumar, Chair(s) 22 Oct 2018 Date Approved Matthew Flatt, Member 22 Oct 2018 Date Approved Jonathan Cohen , Member 22 Oct 2018 Date Approved none , Member 22 Oct 2018 Date Approved 22 Oct 2018 Member none, Date Approved

by <u>Ross Whitaker</u>, Chair/Dean of the Department/College/School of <u>Computer Science</u> and by <u>David Kieda</u>, Dean of The Graduate School.

#### **ABSTRACT**

The replicability crisis in experimental science is fueled by a lack of transparent and explicit discussion of experimental design in published work. An experimental design is a description of experimental factors and how to map those factors onto a sequence of trials such that researchers can draw statistically valid conclusions. This thesis introduces SweetPea, a SAT-sampler aided language that facilitates creating understandable, reproducible and statistically robust experimental designs. SweetPea consists of (1) a high-level language to declaratively describe an experimental design, and (2) a low-level runtime to generate unbiased sequences of trials given satisfiable constraints. The high-level language provides primitives that closely match natural descriptions of experimental designs. To ensure statistically significant results, every possible sequence of trials that satisfies the design must have an equal likelihood of being chosen for the experiment. The low-level runtime samples sequences of trials by compiling experimental designs into Boolean logic, which are then passed to a SAT-sampler. The SAT-sampler provides guarantees that the solutions it finds are statistically robust.



## **CONTENTS**

AB	STRACT	ii
LIS	ST OF FIGURES v	ii
LIS	ST OF TABLESvi	ii
NO	OTATION AND SYMBOLS	ix
	APTERS	
1.	MOTIVATION	1
	<ul><li>1.1 Reliable Experimental Design and Reproducibility Crisis</li><li>1.2 Declarative Programming: Science without the Engineering Burden</li></ul>	1 1 2
2.	SWEETPEA OVERVIEW	3
	2.2 A Runtime for Uniform Sampling	3 4
3.	RELATED WORK	5
	3.1.1 Reproducibility Crisis  3.2 Domain Specific Languages  3.2.1 Solver Aided Languages  3.3 Combinatorial Search Spaces  3.3.1 Sampling Methods  3.3.2 Boolean Satisfiability	5 5 5 6 6
4.	SWEETPEA LANGUAGE	7
	4.1.1 Descriptions of Stimuli 4.1.2 Ordering Constraints 4.1.3 Experimental Design and Balancing 4.1.4 Experimental Structure 4.2 SweetPea Primitives 4.2.1 Factors and Levels 4.2.2 Derived Levels	7 7 8 8 8 8 9 9
5.	SWEETPEA RUNTIME	0

	<ul> <li>5.1 Conjunctive Normal Form (CNF)</li> <li>5.2 Representing SweetPea Primitives in CNF</li> <li>5.2.1 Representing Levels and Factors</li> <li>5.2.2 Representing Derived Levels and Derivation Functions</li> <li>5.3 Communicating with the SAT-Sampler</li> <li>5.4 Correctness Guarantees</li> </ul>	10 10 11 11
6.	IMPLEMENTING SWEETPEA	12
	<ul> <li>6.1 Language Implementation: Decisions and Alternatives <ul> <li>6.1.1 Embedding Language</li> <li>6.1.2 Internal Representations</li> </ul> </li> <li>6.2 Runtime Implementation: Tsietin Transform <ul> <li>6.2.1 Binding Variables: Iff</li> <li>6.2.2 Adders</li> <li>6.2.3 Ripple Carry Adders</li> <li>6.2.4 Pop Count Circuit</li> <li>6.2.5 Exhaustive Testing</li> </ul> </li> </ul>	12 12 12 12 13 13
7.	COMPILING STROOP IN EXCRUTIATING DETAIL	14
	<ul> <li>7.1 Stroop in SweetPea</li> <li>7.2 Compiling Stroop to CNF</li> <li>7.3 Synthesizing an Experimental Sequence</li> <li>7.4 Verifying the Uniformity Guarantee</li> </ul>	14 14
8.	FUTURE WORK	15
	8.1 Beyond Psychology 8.2 Future Language 8.2.1 Weighted Crossings 8.2.2 Sampling Continuous Factors 8.2.3 Automated Experimental Design 8.2.4 Syntactic Sugar 8.3 Future Runtime 8.3.1 Verified Core 8.3.2 Debugging unSAT experiments 8.3.3 Iterative experimental design and Partial Satisfiability 8.3.4 Optimizations	15 15 15 15 15 15
9.	CONCLUSION	17
RE	FERENCES	18

## LIST OF FIGURES

## LIST OF TABLES

## **NOTATION AND SYMBOLS**

α	fine-structure (dimensionless) constant, approximately 1/137
α	radiation of doubly-ionized helium ions, He++
β	radiation of electrons
$\gamma$	radiation of very high frequency, beyond that of X rays
$\gamma$	Euler's constant, approximately 0.577 215
$\delta$	stepsize in numerical integration
$\delta(x)$	Dirac's famous function
$\epsilon$	a tiny number, usually in the context of a limit to zero
$\zeta(x)$	the famous Riemann zeta function
	•••
$\psi(x)$	logarithmic derivative of the gamma function
$\omega$	frequency

#### **MOTIVATION**

If a scientific experiment is incorrectly designed or biased, the results of that experiment are invalid and meaningless. Creating correct, statistically unbiased, reproducible experimental designs is paramount to performing meaningful experiments. todo

## 1.1 Reliable Experimental Design and Reproducibility Crisis

The replicability crisis in experimental science is fueled by a lack of transparent and explicit discussion of experimental design in published work. While there are many software tools for modeling and running experiments [1] [2] [4], there are, to the best of our knowledge, none for designing them. Currently, scientists design experiments by writing complex scripts which manually balance the experimental factors of interests. There are two major issues with this approach: the first is that it may not (and often does not) produce unbiased sequence of trials. In practice, researchers construct these sequences without any statistical guarantees because the brute force solution for constructing unbiased sequences by enumerating all options is intractable; for a typical experiment the size of the search space is 10<sup>100</sup>. The second issue is that this approach is brittle. It is easy to introduce bugs that go unnoticed but which may have large consequences, and it is difficult to verify and reproduce another researcher's implementation.

## 1.2 Declarative Programming: Science without the Engineering Burden

Virtually all fields of scientific research increasingly rely on computational tools. Computational tools open the door to analyses that are otherwise intractable, time consuming and error-prone. Moreover, writing programs contributes to making research reproducible because it creates digital artifacts like files, which allow one scientist to share, analyze and run a program written by another. The drawback, however, is that writing correct,

maintainable, complex programs takes significant engineering effort. Requiring scientists to be engineers in addition to being highly trained domain specialists needlessly impedes the progress of science. Declarative languages allow their users to describe the result they want, as contrasted with imperative languages which require their users to describe the how to construct the result.

## 1.3 Requirements Statement

There is a need for a software system that allows domain scientist to design unbiased, replicable experiments. Moreover, this system needs to provide an easy-to-use, declarative interface so that scientists who are not necessarily trained as software engineers can create and reason about complicated experimental designs, and transparently share their experimental setups and design choices. SweetPea is just such a system; it is a language which provides semantics for describing experiments, a runtime for synthesizing experimental sequences from specifications, and a set of tools for debugging over-constrained designs. While the need for a system to automate experimental design is general to many types of science, we have built a prototype targeted for psychology and neuroscience, where issues of reproducibility and complexity of design have become a focus of attention.

## **SWEETPEA OVERVIEW**

- there's a language and a runtime
- language has to do with the domain specific representation (plus thinking about how choice of representation influences whether its amenable to SAT)
- runtime has to do with executing a representation of the experiment with the SAT sampler; decisions about the SAT representation have to do with the runtime.

#### 2.1 A Language for Experimental Design

The independent and control variables in an experiment are called *factors*, and a trial is specified by a combination of levels of different factors. As a running example, consider an experiment where subjects are shown shapes of different colors; the factors are "color" and "shape", and each colored shape is a trial. Many experiments have additional constraints on the trials, such as "no more than 4 red shapes in a row".

- briefly mention the other parts of experiments: design, crossings, blocks
- the purpose of the language is make it easy to represent experimental designs, while also being amenable to being translated into SAT

## 2.2 A Runtime for Uniform Sampling

SweetPea can be viewed as a domain-specific interface to SAT-sampling, and while there are other languages that rely on SAT-solvers [5], none that we know of leverage the guarantees provided by SAT-samplers. To ensure statistically significant results, every possible trial sequence that satisfies the constraints must have an equal likelihood of being chosen for the experiment. This guarantees that the method for generating trial sequences is not introducing bias. In practice, however, researchers construct these trial sequences without statistical guarantees. The number of valid sequences is both intractably large and sparse in the space of all sequences, so it is not possible to find a valid sequence by

randomly sampling all sequences or by enumerating all valid sequences.

The runtime to generates unbiased sequences of trials given satisfiable constraints. At the heart of the bias problem is the need to sample from constrained combinatorial spaces with statistical guarantees; SweetPea samples sequences of trials by compiling experimental designs into Boolean logic, which are then passed to a SAT-sampler. The SAT-sampler Unigen provides statistical guarantees that the solutions it finds are approximately uniformly probable in the space of all valid solutions. This means that while producing sequences of trials that are perfectly unbiased is intractable, we do the next best thing– produce sequences that are *approximately* unbiased.

- the runtime provides the uniformity guarantee because the sampler provides the guarantee
- because we're compiling to SAT there are a lot of low level "representation in SAT" decision to be made; need to ensure an \*efficient\* representation
  - provide an estimate for no. of vars and no. of clauses

## 2.3 Running Example: the Stroop Experiment

- outline the simplest stroop (2x2 experiment)
- explain all the parts
- explain the desired analysis
- show a sample sequence
- state that all 4 sequences should be equally likely
- make this into a block diagram perhaps

#### **RELATED WORK**

- there are three related areas: work related to the pyscology computation and experimental design, work related to domain specific and solver-aided languages, and work related to sampling combinatorial spaces.

## 3.1 Psychology Toolboxes

- psyScope [1]
- psychoPy [2]
- OpenSesame [4]

#### 3.1.1 Reproducibility Crisis

- TODO: surely something goes here

## 3.2 Domain Specific Languages

- probably don't need to cite anything here?

#### 3.2.1 Solver Aided Languages

- rosette [5]
- sketch: "Domain-Specific Symbolic Compilation"
- dafny maybe
- hyperkernel: co-designing a language and the verification

## 3.3 Combinatorial Search Spaces

- finding solutions in a large search space- no really, very large
- how large?
- so large
- what is the nature of our search constraints? things like 60 red words, 60 blue words

(in Stroop, see chapter 2).

#### 3.3.1 Sampling Methods

- sampling is the problem of finding solutions
- could try solutions at random: turns out they are sparse (most examples don't have 60 red, 60 blue)
- could try to generate all solutions: turns out there are too many (lots of possible arrangements)
  - could use MCMC, but doesn't provide guarantees
- this project is \*really\* about providing this guarantee that we're not introducing bias because this is a huge deal

#### 3.3.2 Boolean Satisfiability

- SAT is a classic problem, NP-complete
- SAT solvers are really efficient solvers
- To specify something in SAT, you use variables and specify invariants
- the SAT solver finds an assignment that satisfies the invariants
- we use a SAT sampler which finds multiple assignments, and with guarnatees
- an alternative is using SMT constraints
- we compile to SAT because unigen is available; to use a different tool we could pretty easily swap out the backend

## 3.3.3 Uniform Sampling

[3]

- what problem is it solving? want uniform for coverage
- who else cares about this problem
- how is it solved: universal hash functions
- what alternatives exist

#### **SWEETPEA LANGUAGE**

The goal of the SweetPea language is to have semantics that match the terms researchers use to describe their experiments, while also being amenable to being translated efficiently to SAT.

- summary here

TODO: everytime I saw "sometimes" come up with an example.

## 4.1 Components of an Experiment

To motivate our choice of primitives, we will first look at the components of an experiment.

#### 4.1.1 Descriptions of Stimuli

- Usually factor with discrete levels
- experiments typically have 2-7 factors with 2-4 levels each: important because it means a large search space
  - sometimes they might be nested (light color, dark color)
  - sometimes want to sample a continous distribution
- sometimes don't fully cross all of them because can't keep the person there that long, but really want to try all combinations

#### 4.1.2 Ordering Constraints

- need to specify the ordering to run Experiments, ie if you're testing the effect of the presence of A, you better be able to run it with and without A etc
- really these are about relationships (presence or absense) of levels or factors, or arbitrary nestings. (1) congruence and (2) transitions as examples of derived levels
  - sometimes about relationships of other relationships! like, you can imagine balancing

#### transitions

- usually the complexity is limited by experimental limits, ie people actually have to complete these

#### 4.1.3 Experimental Design and Balancing

- often want a full-crossing
- sometimes experiment is too large for full crossing, so while experiment has over attributes (factors) they might not appear in all examples
- sometimes its impossible to fully-counterweight and it'd be awesome if the tool could tell you if you were trying to do something impossible
  - sometimes can fix impossible by weighted crossing
  - sometimes can fix by "near balancing" or toss-away block, ie balancing transitions
- ultimate declarative would be just say "these are what I want to analyze, balance for me"

#### 4.1.4 Experimental Structure

- Experimental structure (ie multiple blocks)
- sometimes want prologue / epilogue blocks
- sometimes need these to balance transitions
- sometimes have an experiment that consists of multiple experiments
- sometimes want to reason between subjects because space is so large

#### 4.2 SweetPea Primitives

- for all: how does it map onto the pysch and why is it amenable to SAT
- details about SAT encodings in the next section

#### 4.2.1 Factors and Levels

- possibly nested lists; that is to say trees
- an experiment has one option for a level on at a time
- we can represent whether a level is on or not as a boolean variable, then have a boolean constraint that says that only one can be on at a time
  - a nesting isn't represented in the boolean logic, it's just a reference to arbitraty group-

ings of levels that can be used in the relationships; can build a factor directly or from these groups. Q: what if overlapping level in multiple groupings? probably need to define it as a level then construct multiple references.

- don't currently support sampling factors that are cont distributions because it's challenging to translate that to discrete boolean SAT

#### 4.2.2 Derived Levels

- how do we represent these relationships?
- because they all have to do with ordering, let's consider a "window". Windows have a width and a stride.
  - example: transitions
  - example: congruence
  - windows implicitly "select" and group levels.
- we can then define functions. they are allowed to depend on the state (on or off) of levels, and use this state to say whether or not \*they\* are "on". Example, congruence. In this way, they're defining new levels (since a level is just a thing that knows whether it's on or off) which is why they're derived factors.
  - these should allow us to define things that skip elements, like the "bait" example
- why is this amenable to SAT? we've got discrete boolean elements whose state depends only upon the state of other boolean elements.
- an advantage is that Levels can have any [printable] type (can be strings, numbers, etc) and then you can build your derivation functions however you like based on those properties. maybe list an example with numbers.
- two options: you can either define a factor with a partition (levels are "where the func is true" and "where the func is false") or by defining and combining levels

#### 4.2.3 Experimental Design, Balacing and Experimental Structure

- really straightforward. An experimental sequence is a list of experimental blocks. An experimental block is at the high level these derivations and at the low level linearized into a list of SAT that the runtime can execute. A design is literally which factors are visible and for balancing we currently only support full-crossings.
  - we'll discuss how these are represented in SAT in the next section

#### **SWEETPEA RUNTIME**

- chapter summary here

## 5.1 Conjunctive Normal Form (CNF)

- variables are index based (important for examples)
- keep track of : 1) variables which need assignemnts, 2) boolean formulas using those vars
  - CNF is standard because it makes search easy; what is it
  - can efficently (how efficeintly) translate into CNF

## 5.2 Representing SweetPea Primitives in CNF

- section summary here:
- mention counting constraints in two flavors

#### 5.2.1 Representing Levels and Factors

- as mentioned in the last chapter, levels correspond naturally to boolean values
- need additional constraint; one hot encoding, exactly one true at a time. literally represent as (a and not b) or (not a and b) this would be a bad choice if there were many levels but given the experimental size this is the best decision.
  - walk through an example
- also need "sum" constraints. These also count up, but we expect these to be much bigger numbers, ie in the experiment w/7 factors w/2 levels, there are 64 of each level. So these we compile to "counting constraints". Basically these are linear inequality. We compile them to SAT efficiently through the Tsietin transform (see next chapter for details).
  - walk through an example

#### 5.2.2 Representing Derived Levels and Derivation Functions

Internally, we pick out the levels, create new levels whose value depends on a boolean relationship of those levels. The relationship is then defined by a literal truth table: - generate a truth table for the defined function by: - each input is a factor - then take all boolean combos and literally run it - that generates a truth table - encode that truth table in the logic by translating to CNF

- Work an example, for instance congruence.

## 5.3 Communicating with the SAT-Sampler

- block diagram of all runtime
- variables marked as "important" vs aux
- translate output to be human readable; easy to integrate with other environments

#### 5.4 Correctness Guarantees

- cite guarantee from unigen
- postulate why this is preserved: TODO

## **IMPLEMENTING SWEETPEA**

- chapter summary

## 6.1 Language Implementation: Decisions and Alternatives

- section summary

#### 6.1.1 Embedding Language

- was in haskell
- now is in python
- in the future, perhaps in racket for macro options
- why is this important / consequences to these choices

#### 6.1.2 Internal Representations

- pros / cons of using one-hot vs binary : trade-off between more variables and more clauses
  - todo: more examples

## 6.2 Runtime Implementation: Tsietin Transform

- necessary to efficiently encode on the scale of 60 of these 120 boolean vars need to be true
  - how efficient is it?
  - works by mimicing popcount circuitry
  - generates a bunch of "junk variables"

#### 6.2.1 Binding Variables: Iff

- new variable "equals" if their values are in lockstep

#### 6.2.2 Adders

- half adder example
- full adder example

## 6.2.3 Ripple Carry Adders

- multiple options here, the one I went with
- trade-offs of options

## 6.2.4 Pop Count Circuit

- pop count circuit example

## 6.2.5 Exhaustive Testing

- example of input / output: see how it's prone to error
- tested for a given size all assignments

## COMPILING STROOP IN EXCRUTIATING DETAIL

- chapter summary
- maybe delete this chapter?
- maybe walk through small Stroop top-to-bottom with all the details
- and then look at the distribution

## 7.1 Stroop in SweetPea

- example code listing
- it'd be nice to have a simple derivation function
- show how these get translated on the high level (ie, a list: color one hot constraints, text one hot constraints, color sum constraints)

## 7.2 Compiling Stroop to CNF

- show how the high level constraints get translated to low level constraints
- explain one of them
- show full DIMACS file

## 7.3 Synthesizing an Experimental Sequence

- show input to unigen and output
- show how that output is translated back to human-readable

## 7.4 Verifying the Uniformity Guarantee

- run it a bunch and histogram results

## **FUTURE WORK**

- chapter summary

## 8.1 Beyond Psychology

- other science domains and what would need to change

#### 8.2 Future Language

- section summary: discussed wishes in chapter 2 section 1, these are the ones that are currently unimplemented

#### 8.2.1 Weighted Crossings

- discussed weighted crossings as a necessary component of experiments in chapter 2; not yet implemented

#### 8.2.2 Sampling Continuous Factors

- this is challenging because how does this get translated to SAT?

#### 8.2.3 Automated Experimental Design

- ANOVA experimental design
- need to make sure that this is correct in the domain in many different flavors of experiment

#### 8.2.4 Syntactic Sugar

- don't have to write the name of the factor when the level name is unique

#### 8.3 Future Runtime

- section summary

\_

#### 8.3.1 Verified Core

- the motivation for SweetPea is that it will make it easy to write correct experiments; that only holds if sweetpea doesn't itself have bugs.
- the tsietin transform is ripe for bugs because it's doing a non-human-readable transformation
  - would be cool to formally verify that the transformations are correct

#### 8.3.2 Debugging unSAT experiments

- why is this a problem: usability: if it's unSAT it just say unSAT but not why. Not very useful.
- can get "minimally unsat core" from SAT solver, maybe can translate that back to user-defined levels to guess what the problem is

#### 8.3.3 Iterative experimental design and Partial Satisfiability

- really want to know if experiment is over-constrained
- maybe could try specify subsets to try to iterately find the most constraints that can be simultaneously satisfied – solvers let you push / pop clauses

#### 8.3.4 Optimizations

- xor constraints
- truth table simplification (QuineMcCluskey)
- choice of SAT encodings and variable v. clauses

## **CONCLUSION**

- Problem we tried to solve
- Why it's important
- What we did
- What are the consequences + forward looking

#### **REFERENCES**

- [1] J. Cohen, B. MacWhinney, M. Flatt, and J. Provost, *Psyscope: An interactive graphic system for designing and controlling experiments in the psychology laboratory using macintosh computers*, Behavior Research Methods, Instruments, & Computers, 25 (1993), pp. 257–271.
- [2] S. Mathôt, D. Schreij, and J. Theeuwes, *Opensesame: An open-source, graphical experiment builder for the social sciences*, Behavior research methods, 44 (2012), pp. 314–324.
- [3] K. S. MEEL, M. Y. VARDI, S. CHAKRABORTY, D. J. FREMONT, S. A. SESHIA, D. FRIED, A. IVRII, AND S. MALIK, Constrained sampling and counting: Universal hashing meets sat solving., 2016.
- [4] J. W. Peirce, *Generating stimuli for neuroscience using psychopy*, Frontiers in neuroinformatics, 2 (2009), p. 10.
- [5] E. TORLAK AND R. BODIK, A lightweight symbolic virtual machine for solver-aided host languages, in ACM SIGPLAN Notices, vol. 49, ACM, 2014, pp. 530–541.