



**NATIONAL AND KAPODISTRIAN UNIVERSITY OF ATHENS**

**SCHOOL OF SCIENCE  
DEPARTMENT OF INFORMATICS AND TELECOMMUNICATIONS**

**POSTGRADUATE STUDIES  
“COMPUTER SYSTEMS: SOFTWARE AND HARDWARE”**

**MASTER THESIS**

# **Datalog Based Symbolic Program Reasoning for Java**

**Christos V. Vrachas**

**Supervisor: Yannis Smaragdakis, Professor NKUA**

**ATHENS**

**October 2019**



**ΕΘΝΙΚΟ ΚΑΙ ΚΑΠΟΔΙΣΤΡΙΑΚΟ ΠΑΝΕΠΙΣΤΗΜΙΟ ΑΘΗΝΩΝ**

**ΣΧΟΛΗ ΘΕΤΙΚΩΝ ΕΠΙΣΤΗΜΩΝ  
ΤΜΗΜΑ ΠΛΗΡΟΦΟΡΙΚΗΣ ΚΑΙ ΤΗΛΕΠΙΚΟΙΝΩΝΙΩΝ**

**ΜΕΤΑΠΤΥΧΙΑΚΟ ΠΡΟΓΡΑΜΜΑ  
“ΥΠΟΛΟΓΙΣΤΙΚΑ ΣΥΣΤΗΜΑΤΑ: ΛΟΓΙΣΜΙΚΟ ΚΑΙ ΥΛΙΚΟ”**

**ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ**

**Συμβολική Συλλογιστική Προγραμμάτων για Java,  
βασισμένη σε Datalog**

**Χρίστος Β. Βραχάς**

**Επιβλέπων: Γιάννης Σμαραγδάκης, Καθηγητής ΕΚΠΑ**

**ΑΘΗΝΑ**

**Οκτώβριος 2019**

# **MASTER THESIS**

Datalog Based Symbolic Program Reasoning for Java

**Christos V. Vrachas      R.N.: M1608**

**SUPERVISOR:   Yannis Smaragdakis, Professor NKUA**

## **ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ**

Συμβολική Συλλογιστική Προγραμμάτων για Java, βασισμένη σε Datalog

**Χρίστος Β. Βραχάς     Α.Μ.: M1608**

**ΕΠΙΒΛΕΠΩΝ:   Γιάννης Σμαραγδάκης, Καθηγητής ΕΚΠΑ**

## ABSTRACT

Motivated by the success of theorem provers as aiding tools in symbolic execution and the convenience that declarative programming languages provide, in this thesis we attempt to introduce a strictly declarative implementation of a theorem prover in Datalog. Our approach, namely static declarative symbolic reasoning is implemented within the Doop framework for Java Pointer Analysis, and it mainly seeks to answer "*Which expressions are implied by another expression within a program*". The main motivation behind that decision was to leverage Doop's powerful infrastructure, and at the same time make it possible for Doop to utilize the reasoner in the future for any of its analyses.

**SUBJECT AREA:** Static Program Analysis, Symbolic Reasoning, Propositional Logic

**KEYWORDS:** flow-sensitivity, path-sensitivity, inference rules and logic proofs

## ΠΕΡΙΛΗΨΗ

Έχοντας ως κίνητρο την επιτυχία των προγραμμάτων απόδειξης θεωρημάτων ως υποστηρικτικά εργαλεία στην συμβολική εκτέλεση και την ευκολία που παρέχουν οι δηλωτικές γλώσσες προγραμματισμού, στην εργασία αυτή επιχειρούμε να εισάγουμε μια αυστηρώς δηλωτική υλοποίηση ενός προγράμματος απόδειξης θεωρημάτων σε Datalog. Η προσέγγισή μας, πιο συγκεκριμένα η στατικά δηλωτική συμβολική συλλογιστική, υλοποιήθηκε στα πλαίσια του εργαλείου Doop για Ανάλυση Δεικτών σε προγράμματα Java, και κυρίως επιδιώκει να δώσει απάντηση στο *"Ποιες είναι οι εκφράσεις οι οποίες συνεπάγονται από άλλες εκφράσεις εντός ενός προγράμματος"*. Το κύριο κίνητρο πίσω από αυτήν την απόφαση ήταν η αξιοποίηση των ισχυρών δομών του Doop και ταυτόχρονα η παροχή της δυνατότητας μελλοντικής χρησιμοποίησης του εργαλείου συλλογιστικής στο μέλλον.

**ΘΕΜΑΤΙΚΗ ΠΕΡΙΟΧΗ:** Στατική Ανάλυση Προγραμμάτων, Συμβολική Συλλογιστική, Προτασιακή Λογική

**ΛΕΞΕΙΣ ΚΛΕΙΔΙΑ:** "ευαισθησία"-ροής, "ευαισθησία"-μονοπατιού, κανόνες συμπερασμού και αποδείξεις προτασιακής λογικής

*To my family...*

## **ACKNOWLEDGEMENTS**

I would like to express my gratitude to my advisor prof. Yannis Smaragdakis for giving me the chance to work on such an interesting topic, within my area of interest, his valuable insights during our discussions and his constant support that kept me motivated.

I would also like to thank my labmates at PLaST lab for sharing their expertise on several occasions and providing suggestions throughout the whole process.

*October 2019*



# CONTENTS

<b>1</b>	<b>INTRODUCTION.....</b>	<b>14</b>
1.1	Thesis Structure .....	14
<b>2</b>	<b>BACKGROUND.....</b>	<b>15</b>
2.1	Static Program Analysis .....	15
2.1.1	Whole-Program vs. Modular Analyses .....	15
2.1.2	Flow-Sensitivity .....	16
2.1.3	Path-Sensitivity .....	16
2.2	Symbolic Execution .....	17
2.3	Automated Theorem Proving.....	17
2.4	Datalog and Doop .....	18
2.4.1	Datalog .....	18
2.4.2	The Doop framework .....	19
<b>3</b>	<b>STATIC DECLARATIVE SYMBOLIC REASONING.....</b>	<b>20</b>
3.1	Schema Relations and Types .....	20
3.1.1	Input Facts and Types .....	20
3.1.2	Relations.....	22
3.2	Program Expression Trees.....	23
3.3	Path Expressions.....	25
3.4	Boolean Symbolic Reasoning.....	30
3.4.1	Propositional Logic Axioms .....	30
3.4.2	Propositional Logic Inference Rules .....	32
<b>4</b>	<b>EVALUATIONS .....</b>	<b>34</b>
<b>5</b>	<b>CONCLUSIONS AND FUTURE WORK.....</b>	<b>35</b>
	<b>Abbreviations - Acronyms.....</b>	<b>36</b>
	<b>Appendices .....</b>	<b>36</b>
	<b>REFERENCES .....</b>	<b>37</b>

## LIST OF FIGURES

Figure 1: Modified Doop facts .....	21
Figure 2: New Doop facts.....	21
Figure 3: Imported facts .....	21
Figure 4: Expression Type .....	22
Figure 5: Main Symbolic Reasoning relations .....	22
Figure 6: (a) method parameter .....	23
Figure 7: (b) method invocation results .....	23
Figure 8: (c) load instance fields .....	23
Figure 9: (d) load static fields .....	23
Figure 10: (e) numeric assignments.....	24
Figure 11: (f) phi assignments.....	24
Figure 12: (a) local variable and (b) unary assignments .....	24
Figure 13: (c) binary assignments and (d) if conditions .....	24
Figure 14: Path predicate example .....	25
Figure 15: Single predecessor rules. ....	26
Figure 16: First predecessor of instruction in method .....	26
Figure 17: Last predecessor of instruction in method .....	26
Figure 18: First and last of multiple predecessors wrappers.....	27
Figure 19: Predecessors ordering rules .....	27
Figure 20: Path expressions after an instruction .....	28
Figure 21: Path expressions before an instruction - single predecessor .....	28
Figure 22: Path expression before multiple predecessors .....	29
Figure 23: Self implication.....	30
Figure 24: $A \wedge B \rightarrow A$ and $A \wedge B \rightarrow B$ .....	30
Figure 25: $A \rightarrow A \vee B$ and $B \rightarrow A \vee B$ .....	30
Figure 26: $A \wedge B \rightarrow B \wedge A$ .....	31
Figure 27: $A \vee B \rightarrow B \vee A$ .....	31
Figure 28: $A \wedge (B \vee C) \rightarrow (A \wedge B) \vee (A \wedge C)$ .....	31
Figure 29: $(A \wedge B) \vee (A \wedge C) \rightarrow A \wedge (B \vee C)$ .....	31
Figure 30: $A \vee (B \wedge C) \rightarrow (A \vee B) \wedge (A \vee C)$ .....	31
Figure 31: $(A \vee B) \wedge (A \vee C) \rightarrow A \vee (B \wedge C)$ .....	31
Figure 32: $A \rightarrow B$ and $B \rightarrow C \vdash A \rightarrow C$ .....	32

Figure 33: $A \rightarrow B$ and $A \rightarrow C \vdash A \rightarrow B \wedge C$ .....	32
Figure 34: $B \rightarrow A$ and $C \rightarrow A \vdash B \vee C \rightarrow A$ .....	32
Figure 35: $A \rightarrow B$ and $\neg B$ and $\neg A \vdash \neg B \rightarrow \neg A$ .....	32
Figure 36: $False \vdash False \rightarrow A$ .....	32
Figure 37: $True \vdash A \rightarrow True$ .....	32
Figure 38: $A \rightarrow (B \vee False) \vdash A \rightarrow B$ .....	33
Figure 39: $A \wedge True \rightarrow B \vdash A \rightarrow B$ .....	33

**LIST OF TABLES**

Table 1: Analysis evaluation ..... 34

## LIST OF SOURCE CODES

Code 1: Java source code .....	20
Code 2: Soot Jimple code.....	20

# 1. INTRODUCTION

In an era of an ever-increasing use of software technologies for almost any everyday task, the requirement for high quality software is indispensable. To that end, several tools leveraging program analysis techniques have been introduced in the literature, mainly aiming to assist in performance optimization or bug finding.

Those tools may either be highly autonomous or possibly rely on external frameworks to back their reasoning. A *symbolic execution* framework stands as an example that is highly dependent on an external constraint solver, otherwise known as theorem prover, that aids it by producing a solution to an input constraint formula which further directs the execution. These tools essentially try to yield proofs, yet they are implemented as complex software systems in a common programming language, not being easily extensible. Declarative programming languages are increasingly becoming preferable over the last few years by the Computer Science community, with applications that span from program analysis to databases and distributed systems. The main power of declarative languages is that they allow to focus on the goal, rather than how to achieve it.

In the context of this work, we attempt to introduce a strictly declarative theorem prover - *symbolic reasoner* - in Datalog. Our main motivation yields from the successful applications of such tools in the literature, while the use of Datalog as the implementation languages originates from its expressiveness and its widely recognized success, especially in the field of Program Analysis. Our reasoner is implemented as an enhancement to the Doop Framework for Java Pointer analysis, aiming to leverage its powerful constructs and the possible utilization of our approach in Doop's core in the future.

## 1.1 Thesis Structure

The rest of this thesis is organized as follows:

We provide a brief background overview of static program analysis, symbolic execution, theorem proving and Datalog in 2 chapter. In the 3 chapter

- We introduce the core types and schema relations for our reasoning in 3.1.
- We extensively describe the generation of program expression trees in 3.2
- We introduce a path-sensitive analysis to reconstruct boolean expressions and represent program paths within 3.3
- In 3.4 we build our symbolic reasoning logic, based on the constructs introduced in the previous sections.

Lastly, in 4 we provide a brief evaluation of our approach that arguably seems to scale relatively well next to a simple Doop context-insensitive analysis, and in 5 we conclude what were the most valuable key points of our approach and the next steps that would be worth exploring.

## 2. BACKGROUND

Several program analysis techniques have been proposed in the literature, in order to aid developers and users discover interesting program properties within their software. For example, one might be interested in finding out whether there are memory leaks or whether some piece of code is reachable.

Such techniques come in different flavors, either static or dynamic, and are usually based on strong mathematical concepts. In this section we provide a brief background on these concepts and frameworks tailored to the development of program analysis tools.

### 2.1 Static Program Analysis

Static program analysis is a program analysis technique that aims to reason about a program's behaviors without actually executing it. Program analysis techniques have been heavily utilized by optimizing compilers since their early stages and they have also found several other applications among the areas of software security and software correctness [16]. The question of whether a program is correct or may terminate for all possible inputs is in general undecidable. However, static analysis techniques are able to tackle undecidability by over-approximating or under-approximating the initial problem, attempting to reason over a simplified version of it.

There is a plethora of different static analysis algorithms in the literature, each one met in several domains. For instance, one may utilize analyses such as *liveness* and/or *pointer* analyses in an optimizing compiler with the intention of eliminating dead code regions or performing a constant propagation/folding optimization. Similarly, a *reachability* analysis that determines whether a specific program point is reachable could be used by a software correctness tool to make sure that an erroneous state is actually never reached.

There is a variety of design choices that may prove essential towards the *scalability* and *precision* of a static program analysis algorithm. We briefly describe some of those choices in the context of this work.

#### 2.1.1 Whole-Program vs. Modular Analyses

Whole-Program analyses should not be confused with *inter-procedural* analyses. An inter-procedural analysis considers multiple functions in order to perform its computations. In the *intra-procedural* setting instead, an analysis would restrict its reasoning within a specific function bound, overlooking the way that function calls or external dependencies affect the computations inside the function. However, someone could also perform a modular-style analysis across a relatively small set of functions, that is an *inter-procedural* analysis.

A *modular*-analysis performs its computations within modular units, disregarding external code dependencies and mainly does not attempt any reasoning across the heap. On the contrary, a *whole-program* analysis reasons about a program's properties at every program point by simultaneously constructing a whole-program heap representation of the program to be leveraged during the analysis computations.

Whether the analysis would reason about a program's properties in a whole-program or a modular manner is of great importance, because such a decision affects the way that

the analysis practitioner would tune its performance - scalability and precision. In the whole-program analysis setting the design tradeoffs between scalability and precision might prove crucial towards the overall analysis reasoning, though a modular-analysis may be able to perform more precise reasoning, due to the relatively restricted area of interest. Whole-program analyses are commonly used in the analysis of languages with complex language features. Those languages, such as the object-oriented ones need to reason about language constructs that may reside on the heap, and thus they may benefit from a whole-program analysis reasoning. Such an analysis may not take into consideration the exact ordering of the instructions of a program, sacrificing precision towards better scalability.

### 2.1.2 Flow-Sensitivity

The concept of whether an analysis is designed with respect to the instruction ordering is called *flow-sensitivity*. On the contrary, an analysis that is not designed that way is called flow-insensitive. Flow-sensitivity is tightly related with the scalability and precision of static program analyses, due to the fact that a flow-sensitive analysis keeps track of program properties at every point of it, e.g. before or after every instruction. Whole-program analyses may usually omit the control-flow constructs during any of their reasoning. However, many whole-program static analyses like the Doop framework's Pointer analysis for Java manage to add flow-sensitivity to their core analysis by slight preprocessing [17].

Several tools utilize state-of-the-art compiler technologies to convert the input source to a lower level *intermediate-representation (IR)*. For instance, there is a plethora of tools that utilize the *LLVM Compiler Infrastructure* [15] for C-like languages, or *Soot - A Java optimization framework* [18] - for JVM languages. These tools may further lower the source code in a *Static Single Assignment (SSA)* form. In SSA form, each variable is assigned exactly once, never to be re-assigned again and it is also defined before any of its uses, thus yielding a flow-sensitive representation.

### 2.1.3 Path-Sensitivity

A path-sensitive analysis is in a way an enhanced version of a flow-sensitive analysis which also considers the path taken up to a specific program point. Such an analysis introduces a path representation usually encoded as the combination of a program's neighboring branch conditional expressions, named *path predicate*. A path predicate is essentially a *Boolean formula*, that is a function of boolean variables whose assignment yields a different control-flow path.

Explicit knowledge of the path taken up to a specific program point allows the analysis to achieve higher precision, ideally eliminating the need for any approximation. However, such knowledge comes at a cost; the number of a program's paths is exponential to the number of branches within the program, thus an analysis that keeps track of all possible paths would not manage to scale. There have been suggested multiple techniques [11], [8] to address such issues in the literature. Some of those techniques manage to achieve scalability by restricting the context of the computations or by introducing loose abstractions of the initial problem. As a static analysis technique a path-sensitive analysis does not eventually execute the program, rather utilizes the path encoding in order to assert properties that hold at specific program points.



One of the adverse effects of pure Static Analysis techniques is that they result in a large amount of *false positives* by trading off precision for scalability. In the context of such techniques, a false positive is said to be an erroneous report of a static analysis tool that a property violation has been discovered within the analyzed program, though such violation does not eventually exist. There have been proposed several techniques of either (semi-) dynamic nature that execute directly the program via code instrumentation, or static solutions that attempt to simulate the execution of a program.

## 2.2 Symbolic Execution

Symbolic Execution is yet another program analysis technique that may be of either static or dynamic flavor and mainly attempts to answer whether certain properties of a program could potentially be violated by a piece of code [5]. At the same time it also results in an automated way to generate test cases for programs under testing.

In contrast to the aforementioned static analysis techniques that do not directly execute a program, Symbolic (or *Concolic*) Execution mainly attempts to simulate the execution of a program by considering symbolic (non-deterministic) values for its input. Usually, during the concrete execution of a program a single execution path is considered for any computation, whilst a Symbolic Execution engine manages to explore a plethora of paths thanks to the symbolic values assigned to its input variables. As such, a concrete execution is said to under-approximate the desired analysis. *Concolic Execution* is a mixture of symbolic and concrete execution that considers concrete input values when possible, thus making symbolic execution feasible in practice. The latter is usually described as Dynamic Symbolic Execution (DSE). Several tools have been proposed in the literature such as SAGE [12] and KLEE [7], that are considered to be the tools of choice when it comes to binary analysis and/or systems testing.

In the symbolic setting, execution is typically driven by a *symbolic execution engine* which at its core maintains some state for every explored path. The core data-structures of such an engine are the *Boolean formulas* that encode an execution path by considering the satisfied conditions taken for that exact path, and a *symbolic store* that keeps a mapping between the variables met within the path and their equivalent symbolic expressions or values. The former are constructed on each branch execution, while the latter is updated after an assignment to the corresponding variable. Established by those two structures for a path, the engine utilizes an *automated theorem prover* whose purpose is to check whether there is any violation of the property under analysis that would lead to a failed execution, besides verifying path feasibility. A path is said to be feasible in such scenario, if there exists an assignment of concrete values to its symbolic encoding that would yield the satisfiability of its boolean formula. Along with *path-explosion* due to the exponential state space exploration, non-deterministic inputs and several other vulnerabilities, symbolic execution efficiency also relies to the use of efficient external theorem provers. Ongoing research attempts to face those threats in order to provide better tools that aid developers and help discover bugs in programs.

## 2.3 Automated Theorem Proving

*Automated theorem proving*, namely the ability to automatically prove a mathematical theorem has been the insatiable desire of the scientific community for an extensive period of

time. Theorem proving's foundation is essentially mathematical logic formulas and modelling. On the other hand, *Computer Science* foundations lie within mathematics and logic. Many times software developers want to prove that a specific property is satisfied for a program and its given input values, like its termination. Computer programs can also be encoded as a logical formula and so formal theorems may be constructed in the programming setting too.

The purpose of an automated theorem prover (ATP) is ideally to answer whether a theorem can be proven or not, providing at the same time a formal proof or a counterexample. Automated theorem proving was initially considered to be part of Artificial Intelligence (AI). However, human assistance would be essential in many cases, leading at the same time to the recognition of unable to be proven theorems (such as termination).

Automated theorem proving has led to the creation of several tools following different approaches to solve the same problem: proof construction. Of particular interest have been satisfiability modulo theories (SMT) provers (or solvers) that are the main tools employed in symbolic execution. SMT is a generalization of the traditional boolean satisfiability problem (SAT) that utilizes underlying decidable mathematical theories in order to encode and give semantical meaning to function predicates instead of boolean (binary) variables. Such theories are the theory of linear arithmetic, bitvectors and arrays and provide a powerful expressiveness to the related solvers which could also be considered as constraint satisfaction engines. Provided a predicate formula encoding, a prover aims to check whether the formula is satisfiable; there is an assignment of values to the compounded predicates. The state of the art tools for theorem proving have been the Z3 Theorem Prover [9], and the Coq Proof Assistant [10], an interactive theorem prover

## 2.4 Datalog and Doop

### 2.4.1 Datalog

*Datalog* is a declarative programming language, essentially a subset of the *Prolog* programming language due to the lack of function symbols, the cut operator among other features. It is also usually considered a query language and has been described as pure SQL enhanced with recursion.

Clause ordering does not matter in Datalog, and thus any order of computation is guaranteed to yield the same results, unlike Prolog. Computation in Datalog is established by declaring *input facts* along with *reasoning rules*. Based on those, anything that can be inferred is computed by the underlying engine reaching a convergence state due to the fact that only new facts, that is asserted truth, can be produced by the inference rules. Such rules are called *monotonic* and the associated computation is the *fixed-point*: rules produce new facts until a convergence state is reached.

The main statements of Datalog are the inference rules, usually consisting of the *head* and the *body*. Both the head and the body may consist of conjunctions (expressed with the comma operator ,) or disjunctions (expressed with the semicolon operator;) of predicates, equivalently called relations. Sharing many similarities with databases, a relation can be thought of as a database table, whilst input facts and inferred knowledge can be thought of as the database rows. The comma operator along with the variable constraints that may form along the operand relations describes what is essentially known as the join operator in the database community.

Datalog has been extensively used both in the academic world and the industry. Several prominent knowledge databases have been designed with Datalog being the underlying engine (i.e. LogicBlox). Datalog also finds use cases in the construction of network specification files, distributed systems, security analysis of smart contracts [13] and static program analysis. One of the static analysis tools that have successfully embodied Datalog to express its analyses is the Doop framework.

### 2.4.2 The Doop framework

The Doop framework is a static program analysis tool particularly performing pointer analysis for the languages of the Java Virtual Machine (JVM) ecosystem. Late work has also included analysis of Python, and especially TensorFlow programs.

Doop analyses are expressed in the Datalog language in a purely declarative setting and the framework itself is known to be the first to introduce full, end-to-end context-sensitive analyses for JVM based languages. The use of Datalog has lead to the generation of new algorithms for pointer analysis that manage to scale, yet provide accurate results.

Doop at first utilizes the Soot framework to convert the initial program, usually compressed as a Java ARchive (JAR) into Java bytecode. Based on the bytecode version of the program, Doop generates the input facts for the analysis, introducing several relations that can be pre-computed. The framework then includes the Datalog based Souffle tool for static analysis [14] in order to perform any reasoning starting from the provided input facts. That is, Souffle is the tool in which any reasoning happens. Doop analyses have been expressed by a bunch of Datalog rules in Souffle Datalog. Even though several context-sensitive analyses have been expressed, there is not any path-sensitive analysis within the framework, thus it relies on the flow-sensitivity introduced by Soot.

### 3. STATIC DECLARATIVE SYMBOLIC REASONING

In this chapter we present our approach towards static declarative symbolic reasoning. The core of the reasoning is encoded as a set of Souffle Datalog rules, along with the related input facts. Doop includes several preprocessing steps, expressive analyses and postprocessing parts that can be leveraged to enhance the core of our reasoning at relatively swift speed. The motivating Java example along with the associated Jimple code (Soot's IR) are presented in 1 and 2 respectively.

<pre> boolean and(boolean t, boolean f) {     boolean x = t &amp;&amp; f;     return x; } </pre>	<pre> boolean and(boolean t, boolean f) {     boolean y, x, z;      if t == 0 goto label1;     if f == 0 goto label1;     x = 1;     goto label2; label1:     y = 0; label2:     z = phi(x, y);     return z; } </pre>
--	--

Code 1: Java source code

Code 2: Soot Jimple code

## 3.1 Schema Relations and Types

### 3.1.1 Input Facts and Types

Doop includes a set of fact generators that provide us with an extensive set of input facts that may be leveraged for our reasoning. Those generators either target different input languages or different intermediate representations (IR). As an example, there are generators that target the Java language, Android technology or Python. These generators may further utilize different tools that would take care of the translation of the input program into the corresponding IR; the current tools being used by Doop are namely Soot and WALA. In the case of Java, the generators ought to produce the same set of input facts that Doop consumes for its main analyses, and thus there is a set of files that define the common predicates-facts that should be produced by the generators. The existing set of generated facts is quite extensive and it is constantly further enriched, yet it does not provide us with a set of facts crucial for the reasoning we introduce. Hosting a rich whole of various pointer analyses, Doop mainly needs knowledge that is related to pointers, or, in the Java setting, reference types. As such, Doop lacks presence of multiple facts in regard to primitive types or values and in case that they exist they are relatively inadequate. Towards this end we decided to modify the common fact generator in a way that enhances Doop with facts to be consumed by our symbolic reasoning approach. At the same time these facts may also come handy for Doop itself in the future. In order to be able to construct a fine-grained analysis we modified a certain amount of facts as follows in figure 1.

```

__AssignOperFrom(insn, position, from).
__IfVar(insn, position, var).

```

**Figure 1: Modified Doop facts**

`__AssignOperFrom` introduces a new position variable that keeps track of each subsequent operand of an instruction's right-hand side (rhs), whilst the same also holds for `__IfVar` too. Next, we provide a glimpse of all the newly introduced generated facts in figure 2.

<code>__OperatorAt(insn, operator).</code>	<code>AssignOper__From(insn, position, from).</code>
<code>__DummyIfVar(insn, var).</code>	<code>AssignOper__FromConstant(insn, position, from)</code>
<code>__AssignOperFromConstant(insn, position, from).</code>	<code>If__Var(insn, position, var).</code>
<code>__IfConstant(insn, position, constant).</code>	<code>If__Constant(insn, position, constant).</code>
	<code>DummyIf__Var(insn, var).</code>
	<code>Operator__At(insn, operator).</code>

**Figure 2: New Doop facts****Figure 3: Imported facts**

These input facts share many similarities in regards to the associated variables - relation columns - though each one encodes crucial knowledge. The `__OperatorAt` relation holds facts that encode the operation being present at a given instruction. `__AssignOperFromConstant` is an alternate rule of `__AssignOperFrom`, while `__IfConstant` is analogous to `__IfVar`. A constant type is different to a variable type in the context of the framework, thus the need for a more detailed differentiation between the two. Such differentiations have extensively been made during the development of our approach, but in many cases will be omitted in the context of this text. These relations essentially encode the operands met in a unary/binary or if-condition expression accordingly. Lastly, the `__DummyIfVar` relation has been introduced in order to encode the conditional expressions at a given branch instruction, that is quite necessary to our implementation and its need is evoked due to the IR provided by Soot. It is also worth mentioning that the facts and the associated relations as described above are provided by the generators written in Java. Mere preprocessing takes place before importing the final facts into Doop's knowledge base in order to follow its conventions. That preprocessing leads to the following massaged relations as described in figure 3.

The relations described thus far are only a small set of the input to our reasoning, as we have made use of the various relations and rules that already existed in Doop. Here we have only presented the brand new input relations or the ones that we modified. However we cannot provide an extensive presentation of all, thus we will briefly describe them in the following sections, in the context of the rules that they are used.

Our symbolic reasoning approach reasons about a program's possible expressions. That said, there is a need for the representation of a program's expression, whether primitive or not. Doop does not include such an encoding, due to the fact that its analyses mainly reason about pointers. In order to address this and also lay the foundations of our approach, we have introduced the appropriate types as described in figure 4.

The types that are defined using a (`|`) symbol are said to inherit from the associated types, while the latter definition introduces a new record type. The expression type manages to represent both unary, binary and constant expressions at the same time, by distinguishing each case with the help of the *base* field. If *base* is of type *SymbolicInput* then the latter *Expr* fields are nil, while in the case of *base* being an *Operator* at least the first of the *Expr* fields is indeed an expression. For each distinct expression we also preserve its

```

.type SymbolicInput = Var | MethodInvocation | NumConstant
.type Operator = symbol
.type ExpressionType = PrimitiveType | ReferenceType
.type Base = SymbolicInput | Operator
.type Expr = [
  base: Base ,
  type: ExpressionType ,
  left: Expr ,
  right: Expr
]

```

Figure 4: Expression Type

type, which can either be *PrimitiveType* or *ReferenceType*. These types are the building blocks for the construction of the program's expressions. It may come natural to the reader that the encoding leads to complex *expression trees*. This expression type is essentially what powers our declarative symbolic reasoning approach, encoding both the primitive expressions of the input program and its control flow constructs at the same time.

### 3.1.2 Relations

We provide here a brief presentation of the *main* declared relations that make our approach work in figure 5.

```

.decl ResolveExpr(meth: Method, var: symbol, expr: Expr)
.decl isExpr(expr: Expr)
.decl isArithmeticExpr(expr: Expr)
.decl isReferenceExpr(expr: Expr)
.decl isBooleanExpr(expr: Expr)
.decl isBooleanExprLeftRight(exprOther: Expr, exprX: Expr, exprY: Expr,
                             op: Operator)
.decl isBooleanExprLeftRightInMethod(exprOther: Expr, exprX: Expr, exprY: Expr,
                                     op: Operator, meth: Method)
.decl BuildPathExprBefore(meth: Method, prev: Instruction, exprBase: Expr,
                          insn: Instruction)
.decl PathExpressionBefore(meth: Method, insn: Instruction, pathExpr: Expr)
.decl PathExpressionAfterTrue(meth: Method, insn: Instruction, pathExpr: Expr)
.decl PathExpressionAfterFalse(meth: Method, insn: Instruction, pathExpr: Expr)
.decl ExprImpliesOther(expr: Expr, exprOther: Expr)
.decl ExprIsAlwaysTrue(expr: Expr)
.decl ExprIsAlwaysFalse(expr: Expr)
.decl ExprIsNegationOfOther(expr: Expr, exprOther: Expr)

```

Figure 5: Main Symbolic Reasoning relations

These relations and their associated rules are going to be thoroughly described in the following section where we will be introducing the main parts of *declarative symbolic reasoning*. *ResolveExpr* is the main relation that constructs the program's expressions derived from the input statements and expressions, while the *PathExpression\** relations are responsible for the construction of the control flow expressions, essentially encoding the

branches taken up to a specific program point. Lastly, our declarative symbolic reasoner is powered by the rules of *ExprImpliesOther*, *ExprIsAlwaysTrue*, *ExprIsAlwaysFalse* and *ExprIsNegationOfOther* relations, attempting to prove *implications* between the expressions and assert those that hold *true* or *false* equivalently, if possible.

### 3.2 Program Expression Trees

Introducing a Souffle type that represents a program's expressions is only one step towards symbolic reasoning over a program. For our reasoner to be able to provide us with meaningful results, we have to generate its world of expressions. To that end, we have introduced several rules that build expression trees from the given program's input relations, which essentially encode the expressions of a program in a symbolic manner. The first set of these rules originate from the *ResolveExpr* relation that was declared in figure 5 of the previous section.

The *ResolveExpr* rules are divided into two sets of rules. The ones that encode the base, *symbolic input* expressions and those that lead to the construction of the complex composite expressions. Symbolic input expressions represent the symbolic inputs of a method. We consider (a) method parameters, (b) method invocation results, (c) instance and static fields loading, (d) numeric constant assignments and (e) phi assignments as symbolic inputs. These input cases form the basis of our expressions, and it is also worth mentioning that our reasoner mainly asserts any new knowledge in an intra-procedural fashion. The rules for the enumerated symbolic input expressions are presented in the following figures.

```
ResolveExpr(?meth, ?var, ?expr) :-
  isAssignLocal_Insn(?insn),
  AssignLocal_From(?insn, ?param),
  FormalParam(_, ?meth, ?param),
  Var_Type(?param, ?type),
  (isPrimitiveType(?type) ;
  isReferenceType(?type)),
  AssignInstruction_To(?insn, ?var),
  Instruction_Method(?insn, ?meth),
  ?expr = [?param, ?type, nil, nil] .
```

Figure 6: (a) method parameter

```
ResolveExpr(?meth, ?var, ?expr) :-
  isMethodInvocation(?insn),
  MethodInvocation_Method(?insn, ?sig),
  Method_ReturnType(?sig, _, ?type),
  (isPrimitiveType(?type) ;
  isReferenceType(?type)),
  Instruction_Method(?insn, ?meth),
  AssignReturnValue(?insn, ?var),
  ?expr = [?param, ?type, nil, nil] .
```

Figure 7: (b) method invocation results

```
ResolveExpr(?meth, ?var, ?expr) :-
  isLoadInstanceField_Insn(?insn),
  LoadInstanceField_To(?insn, ?var),
  Var_Type(?var, ?type),
  (isPrimitiveType(?type) ;
  isReferenceType(?type)),
  Instruction_Method(?insn, ?meth),
  ?expr = [?var, ?type, nil, nil] .
```

Figure 8: (c) load instance fields

```
ResolveExpr(?meth, ?var, ?expr) :-
  isLoadInstanceField_Insn(?insn),
  LoadInstanceField_To(?insn, ?var),
  Var_Type(?var, ?type),
  (isPrimitiveType(?type) ;
  isReferenceType(?type)),
  Instruction_Method(?insn, ?meth),
  ?expr = [?var, ?type, nil, nil] .
```

Figure 9: (d) load static fields

The various relations used in the rules bodies are quite self-explanatory which proves the expressiveness and high quality of Doop's code. We present the recursive rules that



```
ResolveExpr(?meth, ?var, ?expr) :-
    isAssignNumConstant_Insn(?insn),
    AssignInstruction_To(?insn, ?var),
    Var_Type(?var, ?type),
    (isPrimitiveType(?type) ;
    isReferenceType(?type)),
    Instruction_Method(?insn, ?meth),
    AssignNumConstant_Id(?insn, ?const),
    ?expr = [?const, ?type, nil, nil] .
```

Figure 10: (e) numeric assignments

```
ResolveExpr(?meth, ?var, ?expr) :-
    Instruction_Method(?insn, ?meth)
    isAssignLocal_Insn(?insn),
    AssignInstruction_To(?insn, ?var),
    Var_Type(?var, ?type),
    (isPrimitiveType(?type) ;
    isReferenceType(?type)),
    c = count: AssignInstruction_To(_, ?var),
    c > 1, ?expr = [?var, ?type, nil, nil] .
```

Figure 11: (f) phi assignments

apply for the cases of (a) local variable assignments, (b) unary assignments, (c) binary expression assignments and (d) if conditions in the associated figures.

```
ResolveExpr(?meth, ?var, ?expr) :-
    Instruction_Method(?insn, ?meth)
    isAssignLocal_Insn(?insn),
    AssignInstruction_To(?insn, ?var),
    Var_Type(?var, ?type),
    (isPrimitiveType(?type) ;
    isReferenceType(?type)),
    c = count: AssignInstruction_To(_, ?var),
    c = 1, AssignLocal_From(?insn, ?from),
    ResolveExpr(?meth, ?from, ?expr).
```

Figure 12: (a) local variable and (b) unary assignments

```
ResolveExpr(?meth, ?var, ?expr) :-
    isAssignUnop_Insn(?insn),
    AssignInstruction_To(?insn, ?var),
    Var_Type(?var, ?type),
    isPrimitiveType(?type) ;
    isReferenceType(?type)),
    Instruction_Method(?insn, ?meth),
    AssignOper_From(?insn, _, ?right),
    _OperatorAt(?insn, ?op),
    ResolveExpr(?meth, ?right, ?rExpr),
    ?expr = [?op, ?type, ?rExpr, nil] .
```

```
ResolveExpr(?meth, ?var, ?expr) :-
    FUNCTION(?var, ?right1, ?right2, ?op),
    Var_Type(?var, ?type),
    (isPrimitiveType(?type) ;
    isReferenceType(?type)),
    ResolveExpr(?meth, ?right1, ?lExpr),
    ResolveExpr(?meth, ?right2, ?rExpr),
    ?expr = [?op, ?type, ?lExpr, ?rExpr] .
```

```
ResolveExpr(?meth, ?var, ?expr) :-
    FUNCTION(?var, ?right1, ?right2, ?op),
    ResolveExpr(?meth, ?right1, ?lExpr),
    ResolveExpr(?meth, ?right2, ?rExpr),
    ?expr = [?op, "boolean", ?lExpr, ?rExpr] .
```

Figure 13: (c) binary assignments and (d) if conditions

For simplicity, several details have been omitted from the rules as defined in figure 13 and they have been hidden under the *FUNCTION* relation. This function relation is for demonstration purpose only and it is responsible for emitting the appropriate variables and generating their corresponding expressions in either case of binary expression assignment or if conditions for the rules to work. The rules behind the *FUNCTION* relation essentially attempt to generate the appropriate combination of operands and expressions based on whether an operand is a constant or a variable.

As previously mentioned, the rules presented in this section constitute one of the most fundamental parts of our reasoner. They construct the most primitive, yet essential expressions of the analyzed program. However, the expressions are guarded in the sense that, if unrolled, they would all lead to complex expressions between symbolic expressions. We shall mention here that any reasoning is purely symbolic, and thus there is not



any concrete evaluation of the expressions in the sense of constant folding. One could also notice the lack of boolean expressions up to this point. This is due to Soot's IR which converts complex boolean expressions into different control flow paths. In the following section we provide an overview of the algorithm responsible for generating these boolean expressions that encode the control flow constructs of the methods of the input program.

### 3.3 Path Expressions

Any Doop reasoning happens on facts produced at a lower level IR called Jimple, Soot's IR. Jimple omits several complex constructs, such as boolean expressions. In our case though, we need to have an exact representation of these high level control flow constructs in order to be able to reason about a program and answer several useful questions about it. In order to further facilitate and enhance our core symbolic reasoning technique we have introduced several algorithms that are responsible for the reconstruction of these expressions representing the program's branch conditions or complete path predicates.

Even though the analyzed program is quite simplified and the various complex control flow structures have been lowered to way simpler forms, the reconstruction of the desired predicates is not trivial at all. Complex expressions between operands that contain operators such as `&&`, `||` and `!` have been simplified to three-address code (TAC) that completely trims these operators in favor of if statements with simplified boolean conditions. Even though we could modify the code transformer that lowers the initial program in a way that would avoid these transformations, we decided to preserve any external framework functionality as it is, implementing any predicate reformation in pure Datalog. Following we introduce the algorithm that implements the desired functionality in two parts. For the sake of pure logic formality we consider  $\wedge$ ,  $\vee$  and  $\neg$  the equivalents of *and*, *or* and *negation* operators.

At first we provide a high level overview of the algorithms responsible for the restoration of the boolean predicates. The basis of each predicate is the first boolean expression met in the program, that is usually the first branch instruction. The algorithm keeps track of the path predicates at any program point before and after each instruction. Each branch instruction splits two distinct path predicates, whilst any other instruction does not produce any new path predicate but rather emits the predicate that holds at the program point before it. However, an instruction may have multiple predecessors, that lead the execution path up to the program point before it. The decision to consider only a single predecessor would lead to inaccurate results as several program paths would have been lost. In order for our algorithm to yield sound path predicates we join the expressions that hold after each predecessor at the program point before such an instruction. Recall that any reasoning is symbolic and there is not any concrete evaluation, thus we cannot be sure of what path predicates should be discarded. This brief introduction reveals that our algorithm needs to implement the following two functionalities: (a) the enumeration of all possible predecessors of an instruction and (b) the making of the path predicates for every program point. A path-predicate for label1 in 2 should be thought as follows:

$$(f == 0 \wedge t \neq 0) \vee t == 0 \quad (3.1)$$

**Figure 14: Path predicate example**

In the following figures we are going to introduce the algorithmic parts responsible for the

instruction predecessors enumeration. The whole logic is implemented in pure Datalog rules like the rest of our work. Doop contains a set of rules populating relations that represent various control flow constructs, such as the set of basic blocks in each function which we made use of. The simplest case of this enumeration is the single predecessor case. We present the relations and the associated rules for this one, along with the *helper* rules that populate the relations which keep track of the first and last predecessor of an instruction in the following figures.

```
SinglePredOfInsnInMethod(?pred, ?insn, ?meth) :-
    Instruction_Method(?insn, ?meth),
    FirstPredOfInsnInMethod(?meth, ?pred, ?insn),
    LastPredOfInsnInMethod(?meth, ?pred, ?insn).
```

**Figure 15: Single predecessor rules.**

```
FirstPredOfInsnInMethodOrdinal(?meth, ?firstOrd, ?insn) :-
    Instruction_Method(?insn, ?meth),
    ?firstOrd = min ord(?prev): MayPredecessorModuloThrow(?prev, ?insn).
```

```
FirstPredOfInsnInMethod(?meth, ?first, ?insn) :-
    Instruction_Method(?insn, ?meth),
    MayPredecessorModuloThrow(?first, ?insn),
    ?min_ord = ord(?first),
    FirstPredOfInsnInMethodOrdinal(?meth, ?min_ord, ?insn).
```

**Figure 16: First predecessor of instruction in method**

```
LastPredOfInsnInMethodOrdinal(?meth, ?lastOrd, ?insn) :-
    Instruction_Method(?insn, ?meth),
    ?lastOrd = max ord(?prev): MayPredecessorModuloThrow(?prev, ?insn).
```

```
LastPredOfInsnInMethod(?meth, ?last, ?insn) :-
    Instruction_Method(?insn, ?meth),
    MayPredecessorModuloThrow(?last, ?insn),
    ?max_ord = ord(?last),
    LastPredOfInsnInMethodOrdinal(?meth, ?max_ord, ?insn).
```

**Figure 17: Last predecessor of instruction in method**

The above rules introduce the handling of the least complicated cases. With just a few rules we manage to identify the one and only predecessor of an instruction. In figure 15 we present the rule that populates the relation keeping track of the single predecessors of an instruction which is powered by the rules of FirstPredOfInsnInMethod and LastPredOfInsnInMethod relations. In these two presented in figures 16 and 17 we leverage Doop's control-flow graph (CFG) analysis by using the MayPredecessorModuloThrow relation that provides us with the predecessors of an instruction. We further aggregate with the help of min and max aggregation functions of Souffle in order to distinguish the first and last predecessors of an instruction accordingly. These rules are also utilized during the enumeration of the multiple predecessors of a single instruction, but they are wrapped into complex rules that separate the multiple from the single predecessor case. We present these wrappers in figure 18.

The negation of LastPredOfInsnInMethod in the first rule is required so that the FirstOfMultiplePredsOfInsnInMethod relation considers instructions with multiple predecessors. Otherwise it would be equivalent to the SinglePredOfInsnInMethod relation, while the

```

FirstOfMultiplePredsOfInsnInMethod(?pred, ?insn, ?meth) :-
    Instruction_Method(?insn, ?meth),
    FirstPredOfInsnInMethod(?meth, ?pred, ?insn),
    !LastPredOfInsnInMethod(?meth, ?pred, ?insn).
LastOfMultiplePredsOfInsnInMethod(?pred, ?insn, ?meth) :-
    Instruction_Method(?insn, ?meth),
    !FirstPredOfInsnInMethod(?meth, ?pred, ?insn),
    LastPredOfInsnInMethod(?meth, ?pred, ?insn).

```

**Figure 18: First and last of multiple predecessors wrappers**

same also holds for the LastOfMultiplePredsOfInsnInMethod. Having introduced the first and last predecessors rules for either case, there is the need of introducing the rules of the relations associated with yielding an ordering of each distinct predecessor. Datalog does not come with any explicit iteration syntax, thus we need to come with a way to enumerate those possible predecessors in a recursive style. We present the steps towards the multiple predecessors enumeration in figure 19.

```

NotNextPredOfInsnInMethod(?meth, ?prev, ?next, ?insn) :-
    Instruction_Method(?insn, ?meth),
    MayPredecessorModuloThrow(?prev, ?insn),
    MayPredecessorModuloThrow(?next, ?insn),
    MayPredecessorModuloThrow(?nextPossible, ?insn),
    ord(?prev) < ord(?next),
    ord(?prev) < ord(?nextPossible),
    ord(?nextPossible) < ord(?next).

NextPredOfInsnInMethod(?meth, ?prev, ?next, ?insn) :-
    Instruction_Method(?insn, ?meth),
    MayPredecessorModuloThrow(?prev, ?insn),
    MayPredecessorModuloThrow(?next, ?insn),
    ord(?prev) < ord(?next),
    !NotNextPredOfInsnInMethod(?meth, ?prev, ?next, ?insn).

NextOfMultiplePredsOfInsnInMethod(?next, ?prev, ?insn, ?meth) :-
    Instruction_Method(?insn, ?meth),
    !LastPredOfInsnInMethod(?meth, ?prev, ?insn),
    NextPredOfInsnInMethod(?meth, ?prev, ?next, ?insn).

```

**Figure 19: Predecessors ordering rules**

The last rule in figure 19 is yet again a wrapper which is used for the predecessors enumeration. It populates the associated relation excluding the last predecessor. Even though there are many different ways to tackle this problem, we decided to provide an intermediate relation which when negated yields the next predecessor in the enumeration. The rule feeding that relation, named NotNextPredOfInsnInMethod, considers the previous and next possible predecessor of an instruction and asserts that in case there is another one between them the rule shall fail. Thus, the negation of this relation yields the actual next possible predecessors which are then used to proceed with the enumeration.

After introducing the part responsible for the predecessors enumeration, we may proceed to the main algorithm that creates the path predicates at every single program point. A high

level overview of the algorithm has been provided earlier, and we now present the core rules of this part in the following figures. We have already introduced the notion of *before* and *after* program points in the context of a distinct instruction and at first we describe the rules responsible for the construction of the path predicates after an instruction via the `PathExpressionAfterTrue` and `PathExpressionAfterFalse` relations.

```

PathExpressionAfterTrue(?meth, ?insn, ?cond),
PathExpressionAfterFalse(?meth, ?insn, ?negatedCond) :-
    FirstIfInsnInMethod(?insn, ?meth),
    NegationOfConditionAtIfInsn(?cond, ?negatedCond, ?insn).
PathExpressionAfterTrue(?meth, ?insn, ?pathExpr1),
PathExpressionAfterFalse(?meth, ?insn, ?pathExpr2) :-
    PathExpressionBefore(?meth, ?insn, ?pathExpr),
    NegationOfConditionAtIfInsn(?cond, ?negatedCond, ?insn),
    ?cond = [?op, ?type, ?exprLeft, ?exprRight] ,
    ?pathExpr1 = ["&&", ?type, ?cond, ?pathExpr] ,
    ?pathExpr2 = ["&&", ?type, ?negatedCond, ?pathExpr] .
PathExpressionAfterTrue(?meth, ?insn, ?pathExpr) :-
    PathExpressionBefore(?meth, ?insn, ?pathExpr),
    !isIf_Insn(?insn).

```

**Figure 20: Path expressions after an instruction**

```

PathExpressionBefore(?meth, ?insn, ?expr) :-
    PathExpressionAfterTrue(?meth, ?pred, ?expr),
    SinglePredOfInsnInMethod(?pred, ?insn, ?meth),
    isIf_Insn(?pred), isJumpTarget(?insn).
PathExpressionBefore(?meth, ?insn, ?expr) :-
    PathExpressionAfterFalse(?meth, ?pred, ?expr),
    SinglePredOfInsnInMethod(?pred, ?insn, ?meth),
    isIf_Insn(?pred), !isJumpTarget(?insn).
PathExpressionBefore(?meth, ?insn, ?expr) :-
    PathExpressionAfterTrue(?meth, ?pred, ?expr),
    SinglePredOfInsnInMethod(?pred, ?insn, ?meth),
    !isIf_Insn(?pred).

```

**Figure 21: Path expressions before an instruction - single predecessor**

The path expression predicates are split into two disjoint sets of *true* and *false* conjuncts due to the separation of a program path after a branch instruction. A branch instruction leads to different execution paths based on whether the condition succeeds or not. In order to differentiate between the two, we have introduced the `NegationOfConditionAtIfInsn` relation that yields the negated condition of a branch condition that would be used to generate the complement of a path predicate. The first rule in figure 20 represents the basis of any path predicate formation, while the second one further splits each path predicate into two subsequent expressions. The last rule represents the path predicate after any statement that is not a branch instruction. However, these rules heavily depend on the existence of path predicates on the program point right before each instruction. In figure 21 we present the rules of the relations that yield the path predicates at the program point before each instruction for the single predecessor case.

The complicated case of the construction of the before-path predicates is essentially when an instruction has multiple predecessors that flow up to it. In order to facilitate the rules of the relation `PathExpressionBefore`, we introduced the `BuildPathExprBefore` relation, responsible for constructing a disjunction of all the predecessor path expressions that meet at the program point right before the instruction. The former relation rules utilize the predecessor enumeration relations introduced before, by implicitly enumerating all over the predecessors. Lastly, we present the logic of the disjunct formation in figure 22.

```

PathExpressionBefore(?meth, ?insn, ?pathExpr) :-
    LastOfMultiplePredsOfInsnInMethod(?pred, ?insn, ?meth),
    BuildPathExprBefore(?meth, ?pred, ?pathExpr, ?insn).
BuildPathExprBefore(?meth, ?pred, ?expr, ?insn) :-
    FirstOfMultiplePredsOfInsnInMethod(?pred, ?insn, ?meth),
    isIf_Insn(?pred),
    IsJumpTarget(?insn),
    PathExpressionAfterTrue(?meth, ?pred, ?expr).
BuildPathExprBefore(?meth, ?pred, ?expr, ?insn) :-
    FirstOfMultiplePredsOfInsnInMethod(?pred, ?insn, ?meth),
    isIf_Insn(?pred),
    !IsJumpTarget(?insn),
    PathExpressionAfterFalse(?meth, ?pred, ?expr).
BuildPathExprBefore(?meth, ?pred, ?expr, ?insn) :-
    FirstOfMultiplePredsOfInsnInMethod(?pred, ?insn, ?meth),
    !isIf_Insn(?pred),
    PathExpressionAfterTrue(?meth, ?pred, ?expr).
BuildPathExprBefore(?meth, ?next, ?pathExpr, ?insn) :-
    BuildPathExprBefore(?meth, ?prev, ?expr, ?insn),
    NextOfMultiplePredsOfInsnInMethod(?next, ?prev, ?insn, ?meth),
    !isIf_Insn(?next),
    PathExpressionAfterTrue(?meth, ?next, ?exprPrev),
    ?pathExpr = ["||", "boolean", ?expr, ?exprPrev] .
BuildPathExprBefore(?meth, ?next, ?pathExpr, ?insn) :-
    BuildPathExprBefore(?meth, ?prev, ?expr, ?insn),
    NextOfMultiplePredsOfInsnInMethod(?next, ?prev, ?insn, ?meth),
    isIf_Insn(?next),
    isJumpTarget(?insn),
    PathExpressionAfterTrue(?meth, ?next, ?exprPrev),
    ?pathExpr = ["||", "boolean", ?expr, ?exprPrev] .
BuildPathExprBefore(?meth, ?next, ?pathExpr, ?insn) :-
    NextOfMultiplePredsOfInsnInMethod(?next, ?prev, ?insn, ?meth),
    BuildPathExprBefore(?meth, ?prev, ?expr, ?insn),
    isIf_Insn(?next),
    !isJumpTarget(?insn),
    PathExpressionAfterTrue(?meth, ?next, ?exprPrev),
    ?pathExpr = ["||", "boolean", ?expr, ?exprPrev] .

```

**Figure 22: Path expression before multiple predecessors**

All of the rules introduced in this section generate the boolean logic formulas, disjuncts and conjuncts. These formulas are modeled with the help of the expression type introduced during the first sections of this thesis. We may now proceed to the description of the core algorithm of our analysis, static declarative symbolic reasoning.

### 3.4 Boolean Symbolic Reasoning

A symbolic reasoner mainly deduces knowledge about the expressions of a program. This knowledge may get as complex as one wants to. The symbolic reasoner tool is essentially a theorem prover that given a friendly input representation and a set of guiding facts tries to reason about a program by constructing formal proofs with the help of a set of inference rules. In the context of this work we have introduced a fundamental symbolic reasoner that mainly deduces knowledge about the boolean expressions of a program in a purely symbolic manner. That is, our reasoner contains pure propositional logic formulas and it does not concretely evaluate any expression on the fly, but rather constantly emits new knowledge until a fixpoint is reached: there can be no more knowledge produced by the set of inference rules. In this chapter we exhibit the core implementation of our reasoner by describing the set of propositional logic axioms and inference rules that bring into life our technique. The whole of those axioms and rules are implemented as pure Datalog rules, without the help of any external tool. The main work of the reasoner is to ask the question of “*What are the program expressions that may be implied by another program expression?*”.

The main relation of the core reasoning algorithm is the ExprImpliesOther relation. Formally, this relation represents the following logical statement:

$$P \rightarrow Q \quad (3.2)$$

where  $P$  and  $Q$  represent any logical (boolean) formula.

#### 3.4.1 Propositional Logic Axioms

We begin with the enumeration of the axioms that constitute the basis of our reasoner. In the following figures we present the Datalog rules describing those axioms.

```
ExprImpliesOther(?expr, ?expr) :-isBooleanExpr(?expr).
```

**Figure 23: Self implication**

```
ExprImpliesOther(?expr, ?exprLeft),
ExprImpliesOther(?expr, ?exprRight) :-
  isBooleanExprLeftRightInMethod(?expr, ?exprLeft, ?exprRight, "&&", ?meth).
```

**Figure 24:  $A \wedge B \rightarrow A$  and  $A \wedge B \rightarrow B$**

```
ExprImpliesOther(?expr, ?exprLeft),
ExprImpliesOther(?expr, ?exprRight) :-
  isBooleanExprLeftRightInMethod(?expr, ?exprLeft, ?exprRight, "&&", ?meth).
```

**Figure 25:  $A \rightarrow A \vee B$  and  $B \rightarrow A \vee B$**

```
ExprImpliesOther(?expr, ?exprOther) :-
    isBooleanExprLeftRightInMethod(?expr, ?exprLeft, ?exprRight, "&&", ?meth),
    isBooleanExprLeftRightInMethod(?exprOther, ?exprRight, ?exprLeft, "&&", ?meth).
```

**Figure 26:**  $A \wedge B \rightarrow B \wedge A$

```
ExprImpliesOther(?expr, ?exprOther) :-
    isBooleanExprLeftRightInMethod(?expr, ?exprLeft, ?exprRight, "||", ?meth),
    isBooleanExprLeftRightInMethod(?exprOther, ?exprRight, ?exprLeft, "||", ?meth).
```

**Figure 27:**  $A \vee B \rightarrow B \vee A$

```
ExprImpliesOther(?expr, ?exprOther) :-
    isBooleanExprLeftRightInMethod(?expr, ?exprLeft, ?exprRight, "&&", ?meth),
    ?exprOr = ["||", "boolean", ?exprB, ?exprC] ,
    ?exprLeft = ["&&", "boolean", ?exprA, ?exprB] ,
    ?exprRight = ["&&", "boolean", ?exprA, ?exprC] ,
    ?exprOther = ["||", "boolean", ?exprLeft, ?exprRight] .
```

**Figure 28:**  $A \wedge (B \vee C) \rightarrow (A \wedge B) \vee (A \wedge C)$

```
ExprImpliesOther(?expr, ?exprOther) :-
    isBooleanExprLeftRightInMethod(?expr, ?exprLeft, ?exprRight, "&&", ?meth),
    ?exprLeft = ["&&", "boolean", ?exprA, ?exprB] ,
    ?exprRight = ["&&", "boolean", ?exprA, ?exprC] ,
    ?exprNRight = ["||", "boolean", ?exprB, ?exprC] ,
    ?exprOther = ["&&", "boolean", ?exprA, ?exprNRight] .
```

**Figure 29:**  $(A \wedge B) \vee (A \wedge C) \rightarrow A \wedge (B \vee C)$

```
ExprImpliesOther(?expr, ?exprOther) :-
    ?exprBC = ["&&", "boolean", ?exprB, ?exprC] ,
    isBooleanExprLeftRightInMethod(?expr, ?exprA, ?exprAB, "||", ?meth),
    ?exprLeft = ["||", "boolean", ?exprA, ?exprB] ,
    ?exprRight = ["||", "boolean", ?exprA, ?exprC] ,
    isBooleanExprLeftRightInMethod(?exprOther, ?exprLeft, ?exprRight, "&&", ?meth).
```

**Figure 30:**  $A \vee (B \wedge C) \rightarrow (A \vee B) \wedge (A \vee C)$

```
ExprImpliesOther(?expr, ?exprOther) :-
    ?exprAB = ["||", "boolean", ?exprA, ?exprB] ,
    ?exprBC = ["||", "boolean", ?exprB, ?exprC] ,
    isBooleanExprLeftRightInMethod(?expr, ?exprAB, ?exprBC, "&&", ?meth),
    ?exprRight = ["&&", "boolean", ?exprB, ?exprC] ,
    isBooleanExprLeftRightInMethod(?exprOther, ?exprA, ?exprRight, "&&", ?meth).
```

**Figure 31:**  $(A \vee B) \wedge (A \vee C) \rightarrow A \vee (B \wedge C)$



### 3.4.2 Propositional Logic Inference Rules

The axioms described before form the ground truth for the reasoner. However, for the reasoner to infer new knowledge we have to define a set of inference rules. An inference rule is the actual means of proving. Based on any number of logical premises an inference rule infers one or maybe even more conclusions.

The formal definition of a logical inference rule that takes only two premises is the following:

$$P, P \rightarrow Q \vdash Q \quad (3.3)$$

where P and Q are metavariables which are essentially logical sentences that may consist of multiple logical statements. Based on the definitions above we introduce a set of inference rules that are the ones responsible for powering our analysis.

```
ExprImpliesOther(?expr, ?exprOther) :-
    ExprImpliesOther(?expr, ?exprInter),
    ExprImpliesOther(?exprInter, ?exprOther).
```

**Figure 32:**  $A \rightarrow B$  and  $B \rightarrow C \vdash A \rightarrow C$

```
ExprImpliesOther(?exprA, ?exprOther) :-
    isBooleanExprLeftRightInMethod(?exprOther, ?exprB, ?exprC, "&&", ?meth),
    ExprImpliesOther(?exprA, ?exprB),
    ExprImpliesOther(?exprA, ?exprC).
```

**Figure 33:**  $A \rightarrow B$  and  $A \rightarrow C \vdash A \rightarrow B \wedge C$

```
ExprImpliesOther(?exprLeft, ?exprA) :-
    isBooleanExprLeftRightInMethod(?exprLeft, ?exprB, ?exprC, "||", ?meth),
    ExprImpliesOther(?exprB, ?exprA),
    ExprImpliesOther(?exprC, ?exprA).
```

**Figure 34:**  $B \rightarrow A$  and  $C \rightarrow A \vdash B \vee C \rightarrow A$

```
ExprImpliesOther(?exprBNeg, ?exprANeg) :-
    ExprIsNegationOfOther(?exprA, ?exprANeg),
    ExprIsNegationOfOther(?exprB, ?exprBNeg).
```

**Figure 35:**  $A \rightarrow B$  and  $\neg B$  and  $\neg A \vdash \neg B \rightarrow \neg A$

```
ExprImpliesOther(?expr, ?exprOther) :-
    ExprIsAlwaysFalse(?expr),
    isBooleanExprLeftRightInMethod(?expr, __, __, __, ?meth),
    isBooleanExprLeftRightInMethod(?exprOther, __, __, __, ?meth).
```

**Figure 36:**  $False \vdash False \rightarrow A$

```
ExprImpliesOther(?exprOther, ?expr) :-
    ExprIsAlwaysTrue(?expr),
    isBooleanExprLeftRightInMethod(?expr, __, __, __, ?meth),
    isBooleanExprLeftRightInMethod(?exprOther, __, __, __, ?meth).
```

**Figure 37:**  $True \vdash A \rightarrow True$



```
ExprImpliesOther(?expr, ?exprOther) :-
  ExprImpliesOther(?expr, ?exprCompOther),
  ?exprCompOther = ["||", ?type, ?exprOther, ?exprFalse] ,
  ExprIsAlwaysFalse(?exprFalse).
```

**Figure 38:**  $A \rightarrow (B \vee False) \vdash A \rightarrow B$

```
ExprImpliesOther(?expr, ?exprOther) :-
  ExprImpliesOther(?exprCompOther, ?exprOther),
  ?exprCompOther = ["&&", ?type, ?expr, ?exprTrue] ,
  ExprIsAlwaysTrue(?exprTrue).
```

**Figure 39:**  $A \wedge True \rightarrow B \vdash A \rightarrow B$

As seen above, we have encoded our tool's core reasoning within a set of few Datalog rules. Even though our approach builds on top of propositional logic theory it manages to express several logical rules of inference and forms the basis for introducing more powerful theories within the scope of our reasoner, such as the theory of arithmetic. It is also worth mentioning that in the context of this text we have not included several implementation details. The whole implementation may be found in Doop's public Github repository [2] under the *symbolic-reasoning* addon logic.

## 4. EVALUATIONS

In this chapter we provide a short evaluation of our approach. The main purpose of this work was to investigate how easily could a theorem prover be implemented in the Datalog programming language. Our main motivation for this work has been the successful application of theorem provers in techniques such as symbolic execution. To this end, we wanted our approach to lay the foundations for the integration of a theorem prover in Doop’s reasoning. We evaluated our approach for three different input JARs based on the DaCapo benchmarks [6]. These JARs are associated to ANTLR [1], HSQLDB [3] and Jython [4]. In the following table we provide some measurements regarding the time needed for the analyses to run. We have evaluated our approach on Doop’s context-insensitive analysis with symbolic reasoning either turned on or off. For the symbolic reasoning runs, we also wanted to evaluate how many program expressions are identified along with how many expression implications are inferred.

**Table 1: Analysis evaluation**

Input Jar	Fact Generation	CI + Symbolic	CI	Expressions	ExprImpliesOther
ANTLR	49 sec	48 sec	33 sec	608.547	863.864
HSQLDB	51 sec	43 sec	25 sec	653.207	916.722
Jython	33 sec	79 sec	77 sec	418.389	590.549

Doop’s context-insensitive analysis performs a pointer analysis without any context consideration for its reasoning. Our symbolic reasoning approach is implemented as complement to any of Doop main analyses, and thus it further burdens the analysis execution times, as observed in table 1. It is also worth mentioning that Doop does not benefit in any way from our approach at the moment. However we may easily notice that our reasoner, even though quite simple, manages to identify a set of expressions, also providing a relatively satisfying number of expression implications. These numbers are quite encouraging, as further restricting the reasoning to discard certain sets of expressions and introducing more logic theories as part of our reasoner would certainly provide more accurate results.

## 5. CONCLUSIONS AND FUTURE WORK

In the previous chapter we provided a brief evaluation of our approach. By introducing our reasoner as part of Doop, we managed to demonstrate several encouraging results. We have also managed to demonstrate the simplicity and expressiveness of a declarative language such as Datalog. The implementation of both core reasoner and its preliminaries came quite naturally to us, proving Datalog to be a tool to be considered for such applications. At the same time we also enhanced Doop with several constructs such as the expression type which could prove valuable to Doop for further usage.

Our work opens a plethora of research directions to investigate in the context of Doop and declarative static analysis in general. For example, it would make sense to explore the case of completely integrating the reasoner in a pointer analysis performed by Doop. Such an investigation would require the utilization of the introduced path-predicates and inference rules. The path-predicates are essentially expressions that encode the control-flow constructs of a program, and thus inferencing over them would lead to the identification of unreachable program locations not to be included during a pointer analysis reasoning. As of the moment Doop completely lacks path-sensitivity, thus we believe that it would probably be of value to research towards this direction.

## ABBREVIATIONS - ACRONYMS

IR	Intermediate Representation
JVM	Java Virtual Machine
CI	Context-Insensitive
CFG	Control-Flow Graph
SSA	Static Single Assignment
DSE	Dynamic Symbolic Execution
SMT	Satisfiability Modulo Theories
SAT	Satisfiability
ATP	Automated Theorem Prover
AI	Artificial Intelligence
JAR	Java ARchive

## REFERENCES

- [1] ANTLR. <https://www.antlr.org>.
- [2] DOOP. <https://github.com/plast-lab/doop-mirror>.
- [3] HSQLDB. <http://hsqldb.org/>.
- [4] Jython. <https://www.jython.org>.
- [5] Roberto Baldoni, Emilio Coppa, Daniele Cono D’elia, Camil Demetrescu, and Irene Finocchi. A survey of symbolic execution techniques. *ACM Computing Surveys (CSUR)*, 51(3):50, 2018.
- [6] S. M. Blackburn, R. Garner, C. Hoffman, A. M. Khan, K. S. McKinley, R. Bentzur, A. Diwan, D. Feinberg, D. Frampton, S. Z. Guyer, M. Hirzel, A. Hosking, M. Jump, H. Lee, J. E. B. Moss, A. Phansalkar, D. Stefanović, T. VanDrunen, D. von Dincklage, and B. Wiedermann. The DaCapo benchmarks: Java benchmarking development and analysis. In *OOPSLA ’06: Proceedings of the 21st annual ACM SIGPLAN conference on Object-Oriented Programming, Systems, Languages, and Applications*, pages 169–190, New York, NY, USA, October 2006. ACM Press.
- [7] Cristian Cadar, Daniel Dunbar, Dawson R Engler, et al. Klee: Unassisted and automatic generation of high-coverage tests for complex systems programs. In *OSDI*, volume 8, pages 209–224, 2008.
- [8] Manuvir Das, Sorin Lerner, and Mark Seigle. Esp: Path-sensitive program verification in polynomial time. In *ACM Sigplan Notices*, volume 37, pages 57–68. ACM, 2002.
- [9] Leonardo De Moura and Nikolaj Bjørner. Z3: An efficient smt solver. In *International conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 337–340. Springer, 2008.
- [10] The Coq development team. *The Coq proof assistant reference manual*. LogiCal Project, 2004. Version 8.0.
- [11] Isil Dillig, Thomas Dillig, and Alex Aiken. Sound, complete and scalable path-sensitive analysis. In *ACM SIGPLAN Notices*, volume 43, pages 270–280. ACM, 2008.
- [12] Patrice Godefroid, Michael Y Levin, and David Molnar. Sage: whitebox fuzzing for security testing. *Communications of the ACM*, 55(3):40–44, 2012.
- [13] Neville Grech, Michael Kong, Anton Jurisevic, Lexi Brent, Bernhard Scholz, and Yannis Smaragdakis. Madmax: Surviving out-of-gas conditions in ethereum smart contracts. *Proceedings of the ACM on Programming Languages*, 2(OOPSLA):116, 2018.
- [14] Herbert Jordan, Bernhard Scholz, and Pavle Subotić. Soufflé: On synthesis of program analyzers. In *International Conference on Computer Aided Verification*, pages 422–430. Springer, 2016.
- [15] Chris Lattner and Vikram Adve. Llvm: A compilation framework for lifelong program analysis & transformation. In *Proceedings of the International Symposium on Code Generation and Optimization: Feedback-directed and Runtime Optimization*, CGO ’04, pages 75–, Washington, DC, USA, 2004. IEEE Computer Society.
- [16] Anders Møller and Michael I. Schwartzbach. Static program analysis, October 2018. Department of Computer Science, Aarhus University, <http://cs.au.dk/~amoeller/spa/>.
- [17] Yannis Smaragdakis and George Balatsouras. Pointer analysis. *Found. Trends Program. Lang.*, 2(1):1–69, April 2015.
- [18] Raja Vallée-Rai, Phong Co, Etienne Gagnon, Laurie Hendren, Patrick Lam, and Vijay Sundaresan. Soot - a java bytecode optimization framework. In *Proceedings of the 1999 Conference of the Centre for Advanced Studies on Collaborative Research*, CASCON ’99, pages 13–. IBM Press, 1999.