



BRNO UNIVERSITY OF TECHNOLOGY

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

FACULTY OF INFORMATION TECHNOLOGY

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ

DEPARTMENT OF INTELLIGENT SYSTEMS

ÚSTAV INTELIGENTNÍCH SYSTÉMŮ

**ADVANCED STATIC ANALYSIS OF ATOMICITY
IN CONCURRENT PROGRAMS THROUGH
FACEBOOK INFER**

POKROČILÁ STATICKÁ ANALÝZA ATOMIČNOSTI V PARALELNÍCH PROGRAMECH
V PROSTŘEDÍ FACEBOOK INFER

MASTER'S THESIS

DIPLOMOVÁ PRÁCE

AUTHOR

AUTOR PRÁCE

Bc. DOMINIK HARMIM

SUPERVISOR

VEDOUČÍ PRÁCE

prof. Ing. TOMÁŠ VOJNAR, Ph.D.

BRNO 2021

Master's Thesis Specification



Student: **Harmim Dominik, Bc.**

Programme: Information Technology and Artificial Intelligence

Specialization: Software Verification and Testing

n:

Title: **Advanced Static Analysis of Atomicity in Concurrent Programs through Facebook Infer**

Category: Software analysis and testing

Assignment:

1. Study limitations of the atomicity analyser Atomer developed in your bachelor thesis as well as the latest developments concerning the Facebook Infer framework.
2. Propose ways of significantly improving precision and/or scalability of the analysis even if for the price of the user providing more input and/or combining it with dynamic analysis.
3. Implement a new version of Atomer including the proposed improvements and supporting analysis of programs written in more programming languages than just C supported by the first version of Atomer.
4. Evaluate the new version of Atomer on suitable benchmarks, including at least real-life code in which some atomicity problems were previously detected.
5. Describe and discuss the achieved results and their further possible improvements.

Recommended literature:

1. Rival, X., Yi, K.: Introduction to Static Analysis: An Abstract Interpretation Perspective. MIT Press, 2020.
2. Blackshear, S., Gorogiannis, N., O'Hearn, P. W., Sergey, I.: RacerD: Compositional Static Race Detection. In: Proc. of OOPSLA'18, PACMPL 2(OOPSLA):144:1-144:28, 2018.
3. Gorogiannis, N., O'Hearn, P.W., Sergey, I.: A True Positives Theorem for a Static Race Detector. In: Proc. of POPL'19, PACMPL 3(POPL):57:1-57:29, 2019.
4. Dias, R.J., Ferreira, C., Fiedor, J., Lourenço, J.M., Smrčka, A., Sousa, D.G., Vojnar, T.: Verifying Concurrent Programs Using Contracts, In: Proc. of ICST'17, IEEE, 2017.
5. Harmim, D.: Static Analysis Using Facebook Infer to Find Atomicity Violations. Bachelor thesis, Brno University of Technology, 2019.
6. Marcin, V.: Static Analysis Using Facebook Infer Focused on Deadlock Detection. Bachelor thesis, Brno University of Technology, 2019.

Requirements for the semestral defence:

- Item 1 and at least some development falling under items 2 and 3 of the assignment.

Detailed formal requirements can be found at <https://www.fit.vut.cz/study/theses/>

Supervisor: **Vojnar Tomáš, prof. Ing., Ph.D.**

Head of Department: Hanáček Petr, doc. Dr. Ing.

Beginning of work: November 1, 2020

Submission deadline: May 19, 2021

Approval date: November 11, 2020

Abstract

Atomer is a *static analyser* based on the idea that if some *sequences of functions* of a *multi-threaded program* are executed *under locks* in some runs, likely, they are *always intended to execute atomically*. *Atomer* thus strives to look for such sequences and then detects for which of them the atomicity may be broken in some other program runs. The author of this master's thesis proposed and implemented the first version of *Atomer* as a plugin of the *Facebook Infer framework* in his bachelor's thesis. In the master's thesis, a new and *significantly improved* version of *Atomer* is proposed. The improvements aim at both increasing *scalability* as well as *precision*. Moreover, support for several initially not supported programming features has been added (including, e.g., the possibility of analysing *C++* and *Java* programs or support for *re-entrant locks* or *lock guards*). Through a number of experiments (including experiments with *real-life code* and *real-life bugs*), it is shown that the new version of *Atomer* is indeed *much more general, scalable, and precise*.

Abstrakt

Nástroj *Atomer* je *statický analyzátor* založený na myšlence, že pokud jsou některé *sekvence funkcí vícevláknového programu* prováděny v některých bězích *pod zámky*, je pravděpodobně zamýšleno, že mají být *vždy provedeny atomicky*. Analyzátor *Atomer* se tudíž snaží takové sekvence hledat a poté zjišťovat, pro které z nich může být v některých jiných bězích programu porušena atomicita. Autor této diplomové práce ve své bakalářské práci navrhl a implementoval první verzi nástroje *Atomer* jako zásuvný modul *aplikačního rámce Facebook Infer*. V této diplomové práci je navržena nová a *výrazně vylepšená* verze analyzátoru *Atomer*. Cílem vylepšení je zvýšení jak *škálovatelnosti*, tak *přesnosti*. Kromě toho byla přidána podpora pro několik původně nepodporovaných programovacích vlastností (včetně např. možnosti analyzovat programy napsané v jazycích *C++* a *Java* nebo podpory pro *reentrantní zámky* nebo *strážce zámek*, tzv. „lock guards“). Prostřednictvím řady experimentů (včetně experimentů s *reálnými programy* a *reálnými chybami*) se ukázalo, že nová verze nástroje *Atomer* je skutečně *mnohem obecnější, přesnější a lépe škáluje*.

Keywords

Facebook Infer, static analysis, abstract interpretation, contracts for concurrency, atomicity violation, concurrent programs, programs analysis, atomicity, incremental analysis, modular analysis, compositional analysis, interprocedural analysis, scalability, *Atomer*, function calls sequence, multi-threaded programs

Klíčová slova

Facebook Infer, statická analýza, abstraktní interpretace, kontrakty pro paralelismus, porušení atomicity, paralelní programy, analýza programů, atomicita, inkrementální analýza, modulární analýza, kompoziční analýza, interprocedurální analýza, škálovatelnost, *Atomer*, sekvence volání funkcí, vícevláknové programy

Reference

HARMIM, Dominik. *Advanced Static Analysis of Atomicity in Concurrent Programs through Facebook Infer*. Brno, 2021. Master's thesis. Brno University of Technology, Faculty of Information Technology. Supervisor prof. Ing. Tomáš Vojnar, Ph.D.

Rozšířený abstrakt

[[Tady nějak zkombinovat a zkrátit úvod + závěr a přeložit do češtiny.]]

Advanced Static Analysis of Atomicity in Concurrent Programs through Facebook Infer

Declaration

Hereby I declare that this master's thesis was prepared as an original author's work under the supervision of professor Tomáš Vojnar. All the relevant information sources used during this thesis's preparation are appropriately cited and included in the reference list.

.....
Dominik Harmim
25th May 2021

Acknowledgements

I thank my colleagues from VeriFIT for their assistance. I would particularly like to thank my supervisor Tomáš Vojnar. I also wish to acknowledge Nikos Gorogiannis from the Infer team at Facebook for valuable discussions about the analyser's development. Lastly, I acknowledge the financial support received from the H2020 ECSEL projects AQUAS, Arrowhead Tools, and VALU3S.

Contents

1	Introduction	1
2	Preliminaries	4
2.1	Concepts in Program Analysis	5
2.2	Abstract Interpretation	11
2.3	Facebook Infer — Static Analysis Framework	16
2.4	Contracts for Concurrency	19
3	Atomer — Atomicity Violations Detector	25
3.1	Related Work	25
3.2	Design of Atomer and Its Principles	26
3.3	Atomer’s Limitations	33
4	Proposal of Enhancements for Atomer	36
4.1	Approximation of Sequences by Sets	36
4.2	Advanced Manipulation with Locks	40
4.3	Analysis’s Parametrisation	45
4.4	Local/Global Atomicity Violations	48
5	Implementation of a New Version of Atomer	50
5.1	Phase 1 — Detection of Atomic Sets	56
5.2	Phase 2 — Detection of Atomicity Violations	60
5.3	Support for New Languages and Locks	64
6	Experimental Evaluation of the New Version of Atomer	66
6.1	Scalability Benchmark	67
6.2	Evaluation on Validation Programs Derived from Gluon	67
6.3	Experiments with Real-Life Programs	69
6.4	Summary of the Evaluation and Future Work	69
7	Conclusion	71
	Bibliography	72
	Appendices	77
A	Contents of the Attached Memory Media	79
B	Installation and User Manual	80
C	Atomer’s Implementation Algorithms	84

Chapter 1

Introduction

Bugs have been present in computer programs ever since the inception of the programming discipline. Unfortunately, they are often hidden in unexpected places, and they can lead to unexpected behaviour, which may cause significant damage. Nowadays, developers have many possibilities of catching bugs in the early development process. *Dynamic analysers* or tools for *automated testing* are often used, and they are satisfactory in many cases. Nevertheless, they can still leave too many bugs undetected because they can analyse only *particular program flows* dependent on the input data. An alternative solution is *static analysis* (which, of course, suffers from some problems too—such as the possibility of reporting many *false alarms*, i.e., *spurious errors*). Quite some tools for static analysis were implemented, e.g., Coverity or CodeSonar. However, they are often proprietary and difficult to openly evaluate and extend.

Recently, Facebook introduced *Facebook Infer*: an *open-source* tool for creating *highly scalable, compositional, incremental, and interprocedural* static analysers. Facebook Infer has grown considerably, but it is still under active development by many teams across the globe. It is employed every day not only in Facebook itself but also in other companies, such as Spotify, Uber, Mozilla, or Amazon. Currently, Facebook Infer provides several analysers that check for various types of bugs, such as buffer overflows, data races and some forms of deadlocks and starvation, null-dereferencing, or memory leaks. However, most importantly, Facebook Infer is a *framework* for building new analysers quickly and easily. Unfortunately, the current version of Facebook Infer still lacks better support for *concurrency bugs*. While it provides a reasonably advanced data race analyser, it is limited to Java and C++ programs only and fails for C programs, which use a *lower-level lock manipulation*. Moreover, the only available checker of *atomicity of call sequences* is the first version of *Atomer* [18] proposed in the bachelor’s thesis of the author.

At the same time, in *concurrent programs*, there are often *atomicity requirements* for the execution of specific sequences of instructions. Violating these requirements may cause many kinds of problems, such as unexpected behaviour, exceptions, segmentation faults, or other failures. *Atomicity violations* are usually not verified by compilers, unlike syntactic or some sorts of semantic rules. Moreover, atomicity requirements, in most cases, are not even documented at all. Therefore, in the end, programmers themselves must abide by these requirements and usually lack any tool support. Furthermore, in general, it is difficult to avoid errors in *atomicity-dependent programs*, especially in large projects, and even more

laborious and time-consuming is finding and fixing them. The papers [13, 15, 37, 41] discuss the importance of *atomicity-related bugs*, and they also show some bugs in *real-world programs*. Unfortunately, tool support for automatically discovering such kinds of errors is currently minimal.

Atomer As already mentioned, within the author’s bachelor’s thesis [18], *Atomer*¹ was proposed — a *static analyser* for finding some forms of *atomicity violations* implemented as a Facebook Infer’s module. In particular, Atomer puts stress on the *atomic execution of sequences of function calls*, which is often required, e.g., when using specific library calls. For example, assume the function `replace` from Listing 1.1 that replaces item `a` in an array by item `b`. It contains an atomicity violation — the index obtained may be outdated when `set` is executed (because, e.g., a *concurrent thread* can modify the array), i.e., `index_of` and `set` should be executed atomically. The analysis is based on the assumption that sequences of function calls executed *atomically once* should probably be executed *always atomically*. Hence, the checker naturally works with sequences of calls.

```

1 void replace(int *array, int a, int b)
2 {
3     int index = index_of(array, a);
4     if (index >= 0) set(array, index, b);
5 }

```

Listing 1.1: An example of an *atomicity violation*

In fact, the idea of checking the atomicity of certain sequences of function calls is inspired by the works of *contracts for concurrency* [13, 37]. In the terminology of [13, 37], the atomicity of specific sequences of calls is the most straightforward (yet very useful in practice) kind of contracts for concurrency. However, while the idea of using sequences in the given context is indeed natural and rather exact, it quite severely limits the *scalability* of the analysis (indeed, even with a few functions, there can appear numerous different orders in which they can be called). Moreover, the implementation of the first version of Atomer targets mainly *C programs* using *PThread locks*. Consequently, there was no support for other languages and their locking mechanisms in the first version of Atomer.

New Developments of Atomer Within this thesis, Atomer has been *significantly improved and extended*. In particular, to improve scalability, working with *sequences* of function calls was *approximated* by working with *sets* of function calls. Furthermore, several new features were implemented: support for *C++* and *Java*, including various *advanced kinds of locks* these languages offer (such as *re-entrant locks* or *lock guards*); or a more precise way of *distinguishing between different lock instances*. Moreover, the analysis has been *parameterised* by function names to concentrate on during the analysis and limits of the number of functions in *critical sections*. These parameters aim to reduce the number of false alarms. Their proposal is based on the author’s analysis of false alarms produced by the first Atomer’s version. Lastly, new experiments were performed to test capabilities of a new version of Atomer.

¹The implementation of a new version of **Atomer** is available at GitHub as an *open-source* repository: <https://github.com/harmim/infer/tree/atomer-v2.0.0>.

Plan of the Thesis The rest of the thesis is organised as follows. Chapter 2 introduces the basics of all the topics related to and essential to this thesis (including *static analysis*, *abstract interpretation*, *Facebook Infer*, and *contracts for concurrency*). The original version of *Atomer*, its limitations, and related work are described in Chapter 3. Subsequently, Chapter 4 presents all the proposed extensions and improvements. The implementation of these extensions is then covered in Chapter 5. The experimental evaluation of the new *Atomer*’s features and other experiments performed within this thesis are discussed in Chapter 6 together with the future work. Finally, the thesis is concluded in Chapter 7. Besides, there are included the following appendices. The content of the attached memory media is listed in Appendix A. Furthermore, Appendix B serves as an installation and user manual. In the end, Appendix C provides several algorithms of the new version of *Atomer*’s implementation that are too long to be placed in the main content of the thesis.

Parts of the thesis concerning the preliminaries and the basic version of *Atomer* are partially taken from the thesis [18]. Moreover, some preliminary results were also published in the Excel@FIT’21 paper [19] written by the author.

Acknowledgement The development of the original *Atomer* started under the H2020 ECSEL projects AQUAS and Arrowhead Tools. The development of its new version is supported by the H2020 ECSEL project VALU3S. It has been discussed with the developers of Facebook Infer too.

Chapter 2

Preliminaries

This chapter provides the background that the thesis builds on. It also explains and describes the existing tools used in the thesis. Note that the author already partially published the contents of the sections in this chapter in his bachelor’s thesis [18] and paper [19]. Indeed, the background of the original and new versions of Atomer are the same. Moreover, the author still tried to improve the description when he found this appropriate.

Multi-Threaded Programs [15, 29] *Multiple threads* of control are commonly used in the software development process because they help reduce *latency*, increase *throughput*, and better utilise *multiple processor computers*. *Threads* or *processes* contain independent *sequences of instructions* that may be performed simultaneously. A process represents a single running program. It has its own address space and a unique identifier. One process can consist of multiple threads, which threads sometimes called *lightweight processes*. All threads of a single process share the same address space for code and data. However, reasoning about the behaviour and correctness of a multi-threaded system is complex due to the need to consider all possible interleavings of the executions of different threads. An integral part of parallel programming is the *synchronisation* of individual threads. Usually, synchronisation mechanisms ensure the *mutual exclusion* of shared resources or synchronise actions that threads perform. Operating systems provide basic synchronisation primitives that can be often used in programming languages in a higher-level manner. Such fundamental mechanisms are *semaphores* (binary semaphores are called *mutexes*), *barriers*, *read-copy-update* techniques, or *monitors*. In this work (and also in practice), for simplicity, the author will use a notion of *locks* for semaphores, mutexes, monitors, etc. as a mechanism that mutually excludes access to a *critical section*.

The use of locks and access to shared data in parallel programs involves the risk of errors not known in sequential programming. Much previous work on checking *thread interference* has focused on *data races*. A data race occurs when two threads simultaneously access the same data variable, and at least one of the accesses is a write. In practice, data races are commonly avoided by guarding each data structure with a lock. Unfortunately, the absence of data races is insufficient to guarantee the absence of errors due to unexpected interference between threads [15, 41].

Atomicity Violations [15] Another possible source of errors due to unexpected interference between threads are *atomicity violations*. A method (or, in general, a code block) is *atomic* iff, for every (arbitrarily interleaved) program execution, there is an equivalent execution with the same overall behaviour where the atomic method is executed serially. In other words, the method’s execution is not interleaved with instructions of other threads. Also, atomicity provides a powerful, indeed maximal, guarantee of non-interference among threads. In short, atomicity is a generally applicable and fundamental *correctness property* of multi-threaded code. Nevertheless, commonly used *testing approaches* are deficient in verifying atomicity. While testing may reveal a concrete interleaving in which an atomicity violation causes erroneous behaviour, the exponentially large number of possible interleavings makes it fundamentally impossible to get suitable *test coverage*. Research results [13, 15, 37, 41] have shown that defects related to atomicity are common, even in well-tested libraries.

In order to be able to check for atomicity violations, it must, of course, be known which blocks of code should be executed atomically. Algorithms for the detection of atomicity violations often deal only with the atomicity of operations on variables. An approach that allows for describing complex *atomicity requirements* is using *contracts* that require atomicity of the execution of certain methods/functions (defined in Section 2.4) [29].

Further, Section 2.1 outlines the fundamental notions and approaches in *program analysis*. In Section 2.2, there is an explanation of *static analysis* by *abstract interpretation*, which is used in *Facebook Infer*, i.e., the key framework used in this thesis. Facebook Infer, its principles, and features are covered in Section 2.3. The concept of *contracts for concurrency* is discussed and defined in Section 2.4.

2.1 Concepts in Program Analysis

This section provides a fundamental intuition about the main principles of *program analysis*. It discusses several standard techniques for reasoning about programs. The section is based on a few first chapters from the book [34] and the overview paper [22].

2.1.1 What to Analyse

The first question is *what programs to analyse*. An obvious characterisation of programs to analyse is the programming language in which the programs are written. Moreover, certain specific families of input programs may be distinguished. However, besides the language and the family of programs to consider, the way input programs are processed can also differ and affect how the analysis works. The most straightforward way to handle input programs is working directly with the *source code* just like a compiler would. Nevertheless, different representations of programs are used in program analysis likewise. The two below types of techniques are often distinguished.

Program-Level Analysis The first possibility is to run the analysis on the *source code* of input programs (e.g., programs written in conventional programming languages like C or Java; or hardware described in VHDL, Verilog, etc.) or on *executable binaries*. This typically involves some *front-end* comparable to a compiler’s one that constructs *syntax trees* of input programs (and possibly other representations such as *control-flow graphs*).

Model-Level Analysis An alternative option is to consider a different input language that aims at modelling the semantics of programs. Then, the analysis takes as an input a description that *models* the program to analyse. Such models either need to be created manually, or some specialised tools are used to extract them from the source code. These models may hide implementation details. Examples of the models are automata of any kind, UML diagrams, Petri nets, Markov chains, or specialised modelling languages like, e.g., Promela.

2.1.2 Static vs Dynamic Approaches

Another critical question in program analysis is *when* the analysis is made. In particular, whether it operates *during* or *before* the program's execution.

One solution is to analyse programs at *run-time*, i.e., *during* the program's execution. This approach is called *dynamic analysis*. It takes place while the program computes, often, over a number of executions.

Another approach is to analyse the program *before* execution, which is called *static analysis*. It is done independently from any execution (at least under the original semantics).

Static and dynamic approaches are significantly distinct in many ways. They come with different benefits and weaknesses. While dynamic approaches are often simpler to design and implement (they may have information about concrete states and thus avoid some false alarms), they often have problems with performance at run-time. They do not force developers to fix bugs before program execution. Moreover, dynamic approaches do have limited coverage of program behaviour since they build on what they see in program executions that were performed. Particular static and dynamic techniques are further discussed in Section 2.1.5.

2.1.3 Automation and Scalability

Automation is a further relevant aspect in program analysis. It would be ideal if program analysis methods were *fully automated* (i.e., no human help is needed). Unfortunately, this is not always possible due to the consequences of the *Rice's theorem* [33]. Thus, sometimes, it is needed to give up on automation and let program analyses ask some *user input* (i.e., some human help is required). In that case, the user is asked to give some information to the analysis, e.g., *invariants*¹. That is, the analysis is partially manual since users need to compute parts of the results themselves. However, having to provide this information may be unwieldy because input programs can be huge or complex.

Scalability is another essential characteristic of program analysis algorithms. Even if a program analysis is fully automatic, it is not guaranteed that it will generate a result within a reasonable time, depending on the complexity of the algorithms. A program analysis tool may not *scale* to extensive programs due to time costs or, e.g., memory constraints.

¹An **invariant** is a *logical property* that can be proved to be inductive for a given program.

2.1.4 Soundness and Completeness

In order to preserve automation and/or scalability, the conditions about program analysis may be relaxed. Namely, the analysis can be proposed to return *inaccurate* results (alternatively, it can return a non-conclusive “do not know” answer). For this purpose, two dual properties (forms of *approximations* or *inaccuracies*) are used. To express these notions, let \mathcal{L} be a *Turing-complete language*, φ be a *non-trivial semantic property* of interest of programs of \mathcal{L} , and T be an analysis tool that decides whether φ holds in a given program. Ideally, if T were absolutely precise, it would be such that:

$$\text{For every program } P \in \mathcal{L} : T(P) = \text{true} \iff P \models \varphi.$$

The above, of course, can be decomposed into a pair of implications:

$$\begin{cases} \forall P \in \mathcal{L} : T(P) = \text{true} \implies P \models \varphi, \\ \forall P \in \mathcal{L} : T(P) = \text{true} \longleftarrow P \models \varphi. \end{cases}$$

Soundness A *sound* program analysis satisfies the first implication.

Definition 2.1.1 (Soundness [34]). The program analyser T is **sound** w.r.t. property φ whenever, for any program $P \in \mathcal{L}$, $T(P) = \text{true} \implies P \models \varphi$.

When a sound analysis terminates and claims that the analysed program has property φ , it guarantees that the program indeed satisfies φ , i.e., no errors are missed. In other words, a sound analysis will refuse all programs that do not satisfy φ . In terms of errors in *binary classification*, there are no *false-negative* errors² when using a sound analysis. On the other hand, there is a chance of *false-positive* errors³.

Completeness A *complete* program analysis satisfies the second implication.

Definition 2.1.2 (Completeness [34]). The program analyser T is **complete** w.r.t. property φ whenever, for any program $P \in \mathcal{L}$, $P \models \varphi \implies T(P) = \text{true}$.

A complete program analysis will accept every program that satisfies property φ . In other words, when a complete analysis refuses an analysed program, it is guaranteed that the program indeed fails to satisfy φ , i.e., there are no false-positive errors. However, there can be false-negative errors.

Due to the *computability* barrier [33], it is not possible to design a general analysis to determine which programs satisfy any non-trivial property for a Turing-complete language that is sound, complete, and fully automated (and ideally also scalable) at the same time. Some of these have to be sacrificed, or the analysis has to be proposed to operate only on a specific set of input programs.

²A **false-negative** error is a real error that is undetected by an analysis tool.

³A **false-positive** error (also called a **false alarm**) is a *spurious error*, i.e., it is detected by an analysis tool, but the error does not exist in the real program.

2.1.5 Program Analysis Techniques

This section describes several standard *program analysis techniques*. Although this thesis focuses on *static analysis* (or maybe rather *bug finding* based on static analysis), it is essential to see the differences between other analysis techniques. Hence, a brief overview is presented in the below sections.

Testing and Dynamic Analysis The *testing* approach checks a finite set of finite program executions. In the development process, several levels of testing are used at various stages of the software/hardware life-cycle, e.g., *unit testing* or *integration testing*. In general, it is challenging to achieve suitable *test coverage* because of infinite program paths when using, e.g., *random testing*. However, several more advanced techniques that improve the coverage have been introduced. These techniques are often combined with other verification approaches. For instance, *concolic testing* combines testing with *symbolic execution*, *dynamic analysis* that observes behaviour in a testing run (such behaviour can be *extrapolated* to behaviour not seen in the given testing run), or *search-based techniques* that can generate test data or parameters. Moreover, to test programs with non-deterministic semantics (i.e., concurrent programs), techniques like *noise injection* are applied. In particular, many advanced dynamic analysers have been proposed to detect data races or deadlocks, including, e.g., Eraser [35] or FastTrack [16].

Testing has the following features. In general, it is simple to automate. It is unsound in almost all cases (besides programs that have a finite number of finite paths). Since failed testing runs provide incorrect concrete executions, the testing is complete. Note, however, that this needs not to hold for extrapolation based dynamic analyses such as the already mentioned Eraser & FastTrack.

Other analysis techniques mentioned below use a static approach. The significant difference between static and dynamic approaches is the following. It is well-known that testing may expose errors, but it cannot prove their absence. It was also famously stated by Edsger W. Dijkstra: “*Program testing can be used to show the presence of bugs, but never to show their absence!*”. However, static approaches may be able to prove their absence—with some *approximation*—they can check *all possible executions* of a program and provide guarantees about its properties. Another static approaches’ benefit is that the analysis can be performed during the development process, so the program does not have to be executable yet, and it already can be analysed. The biggest drawback of static approaches, in general, is that they can produce a lot of false alarms (though this is not the case, e.g., for *theorem proving* or *model checking*), which is often resolved by accepting unsoundness. Another crucial issue of static approaches (this is, however, also an issue of dynamic analyses) is ensuring sufficient scalability—in fact, typically, the more precise the analysis, the less scalable it becomes.

Deductive Verification The (usually) *semi-automated* approach of *deductive verification* (or *theorem proving*) uses *inference systems* for inferring theorems about the analysed system from the facts known about the system and from general theorems of different logical theories. The approach falls under *machine-assisted* techniques, which means that users may be required to provide extra information to the analysis (usually *loop invariants*, procedure *pre-conditions/post-conditions*, assertions, some other invariants, or they may be

required to guide the proof by selecting which inference should be done from what facts and when). This can be demanding and require a very high level of expertise. However, a substantial part of the verification can usually still be carried out in a fully automated fashion. This approach is very general, but there is a problem with generating diagnostic information for incorrect systems. There exist a number of *interactive theorem proving* tools like Coq, Isabelle/HOL, PVS, ACL2, etc. The user usually guides the inference process in these tools.

These techniques also involve *automated decision procedures* (or *satisfiability solvers*) for various logical theories. Such solvers are often used as back-end components for higher-level verification methods, such as *symbolic execution* or *predicate abstraction* in model checking. Commonly used solvers are *SAT-solvers* (e.g., CaDiCaL and Glucose) and *SMT-solvers* (e.g., Z3 and CVC4). Various tools allow the user to provide some *logical annotations* in a code and then automatically attempt to prove specific properties using decision procedures. Examples of such tools are VCC and Dafny.

To sum up, theorem proving techniques have the following properties. They are (usually) not fully automatic, i.e., *high user expertise* is often needed. They are sound w.r.t. the model of the program used in the proof, and they are usually complete.

Model Checking Another technique called (*finite-state*) *model checking* aims primarily at finite systems (extensions to the finite-state systems have also been considered — e.g., model checking with *predicate abstraction*, *regular* model checking, etc.). It automatically verifies whether a system or its model satisfies a particular property based on an *algorithmic exploration of the system’s state space*. Unfortunately, the biggest issue here is the *state space explosion* problem [39]. However, in practice, model checking tools use effective data structures (such as *binary decision diagrams* or *hierarchical storage of states*) to describe program behaviours and avoid enumerating all executions thanks to approaches that reduce the search space. In addition, other techniques are used to cope with this problem. For instance, various abstractions are used (e.g., *predicate abstraction*) or *bounded model checking*⁴. Properties are usually defined using *temporal logics*, such as LTL, CTL, or PCTL. Finite model checking has the following characteristics. It is automatic (up to the need of modelling the system or its environment). It is sound and usually complete. Other advantages are that it is pretty general and provides diagnostic data for incorrect systems.

Model checking is often done at the *model level*, i.e., a model of the program needs to be built, either manually or automatically. Even if the program itself needs not to be modelled, one usually needs a model of its environment. The problem is that the model needs not precisely capture the input program’s behaviours. Thus, the checking of the model may be either incomplete or unsound, w.r.t. the input program. Some model checking techniques can automatically *refine* the model when they fail to prove a property due to a *spurious counterexample*. In practice, model checking tools are often sound and incomplete w.r.t. the input program.

Model checking has found many successful applications, including hardware verification, verification of *concurrent* and *distributed* systems, *probabilistic* systems, *biological* systems, etc. Examples of hardware model checkers are RuleBase, Incisive Verifier, or NuSMV. Model

⁴**Bounded model checking** explores models up to a *fixed depth* only. This technique is often referred to as *bug finding* because it sacrifices completeness and often also soundness. Examples of tools that implement such technique are CBMC, LLBMC, or JBMC.

checkers for concurrent and distributed systems include Spin or DIVINE. Tools, such as CPAchecker or BLAST, are model checkers that use predicate abstraction. PRISM and Storm are state-of-the-art tools for model checking probabilistic systems. Finally, Uppaal is a model checker that verifies temporal logic formulas on *timed automata*.

Static Analysis The last technique (and the most important one for the thesis) is *static analysis*. It relies on other techniques to compute *conservative* (sound but incomplete) descriptions of program behaviours using finite resources rather than building a finite model. The fundamental idea is to finitely *over-approximate* the set of all program behaviours using a particular set of properties.

According to [30], static analysis of programs is reasoning about the behaviour of computer programs without really executing them. It has been used since the 1960s in *optimising compilers* for generating efficient code. More recently, it has proven valuable also for automatic error detection, verification of the correctness of programs, and it is used in other tools that can help programmers. Intuitively, a static program analyser is a program that reasons about the behaviour of other programs by looking for some *syntactic patterns* in the code and/or by assigning the program statements some *abstract semantics* and then deriving a characterisation of the behaviour in terms of the abstract semantics. Nowadays, static analysis is one of the leading concepts in *formal verification*. It can even automatically answer considerably complicated questions about a given program, such as [30]:

- Does the program *terminate* on every input?
- Do two pointers refer to *disjoint data structures* in memory?
- Are *data races* possible? Can the program *deadlock*?
- Does there exist an input that leads to a *null-pointer dereference*, *division-by-zero*, or *arithmetic overflow*?
- Are arrays always accessed *within their bounds*?
- **Are certain operations executed *atomically*?**

In general, conservative static analysis has the following characteristics. It is automatic, it can often handle large systems, and it does not ordinarily need a model of the environment. It produces sound results, and it is generally incomplete, i.e., it can produce many false alarms. However, it is possible to drop soundness to minimise the number of false alarms and preserve automation, which, in fact, is done by many tools. Such techniques are often instead called *bug finding* or *bug hunting*, which is based not only on static analysis but also, e.g., on model checking. Since the primary motivation of these approaches is to discover bugs (and not to prove their absence), they are neither sound nor complete, and they aim at swiftly catching bugs.

The above explanation of static analysis is quite general. Thus, even model checking or deductive verification may be considered as static analysis. However, when setting these approaches aside, the most traditional techniques in the static analysis include:

- *Syntactic checks* looking for various *bug patterns* (*anti-patterns*) — implemented, e.g., in Lint, Cppcheck, or in analyses in VisualStudio, Clang, GCC.

- *Data-flow analysis* — can be sound, but it is often combined with bug pattern searching and used in an unsound way in the end. Many successful tools of this kind have been implemented. They are often proprietary and hard to openly evaluate or extend. Examples of some well-known tools (however, some of them also involve *abstract interpretation*) are Coverity, CodeSonar, Klocwork, FindBugs/SpotBugs, SonarQube.
- *Constraint-based analysis, type-based analysis* (type and effect systems).
- *Symbolic execution* — implemented, e.g., in KLEE, Pex, Symbiotic.
- *Abstract interpretation* (see Section 2.2) — examples of sound tools are Polyspace, AbsInt/Astrée, or Sparrow. State-of-the-art frameworks that allow creating sound analysers include Facebook Infer (see Section 2.3), Facebook SPARTA, or Frama-C. These tools also offer already implemented sound as well as unsound checkers.

Properties of the presented program analysis techniques are summarised in Table 2.1.

Table 2.1: A summary of *program analysis techniques* [34] extended by *dynamic analysis*

Technique	Automatic	Sound	Complete	Object	Approach
Testing	Yes	No	Yes	Program + Test suite	Dynamic
Extrapolated dynamic analysis	Yes	No	No	Program + Test suite	Dynamic
Deductive verification	No	Yes	Yes/No	Model/Program	Static
Model checking	Yes	Yes	Yes/No	Model/Program	Static
Conservative static analysis	Yes	Yes	No	Program	Static
Bug finding	Yes	No	No	Program	Static

2.2 Abstract Interpretation

This section explains and defines the basics of the static analysis technique of *abstract interpretation*. The description is based on [8, 9, 10, 11, 12, 22, 24, 30, 31, 34]. In these works, there can be found a more detailed and formal explanation.

Abstract interpretation was introduced and formalised by the French computer scientist Patrick Cousot and his wife Radhia Cousot at POPL’77 [11]. It is a generic *framework* for static analyses. It allows one to create particular analyses by providing specific *components* (discussed in Section 2.2.1) to the framework. The obtained analysis is guaranteed to be *sound* if specific properties of the components are met.

In general, in *set theory*, which is independent of the application setting, abstract interpretation is considered a theory for *approximating* sets and set operations. A more restricted formulation of abstract interpretation is to interpret it as a theory of approximation of the behaviour of the *formal semantics* of programs. *Fixpoints* may characterise those behaviours (see Section 2.2.3), which is why the primary part of the theory provides efficient techniques for *fixpoint approximation* [31]. Hence, abstract interpretation is used to derive

the approximate abstract semantics over an *abstract domain* (explained in Section 2.2.1). The abstract semantics obtained from program analysis can then be used for verification, optimisation, code generation, etc. [10].

To be sound, it is essential that an approximation performed using abstract interpretation should be an *over-approximation*. This means that the considered abstract semantics should be a *superset* of the concrete semantics, i.e., the abstract semantics should cover all possible cases. Whence, if the abstract semantics satisfies a given property, then the concrete semantics satisfies it too. Moreover, the consequence of the over-approximation of the possible executions is that inexistent executions are considered, leading to *false alarms*. Thus, due to a lack of accuracy, abstract interpretation is usually *incomplete* [9].

Patrick Cousot intuitively and informally illustrates abstract interpretation in [9] as follows. Figure 2.1a shows the concrete semantics of a program by a set of curves, which represents the set of all possible executions of the program in all possible execution environments. Each curve shows the evolution of the vector $x(t)$ of input, state, and output values of the program as a function of the time t . *Forbidden zones* on this figure represent a set of erroneous states of the program execution. Proving that the intersection of the concrete semantics of the program with the forbidden zones is empty may be *undecidable* because the program concrete semantics is, in general, *not computable*. As demonstrated in Figure 2.1b, abstract interpretation deals with the abstract semantics, i.e., a superset of the concrete program semantics. The abstract semantics includes all possible executions. That implies that if the abstract semantics is safe—i.e., it does not intersect the forbidden zones—the concrete semantics is safe as well. However, the over-approximation of the possible program executions causes that inexistent program executions are considered, which leads to false alarms, as demonstrated by Figure 2.1b. It is the case when the abstract semantics intersects the forbidden zone, whereas the concrete semantics does not intersect it.

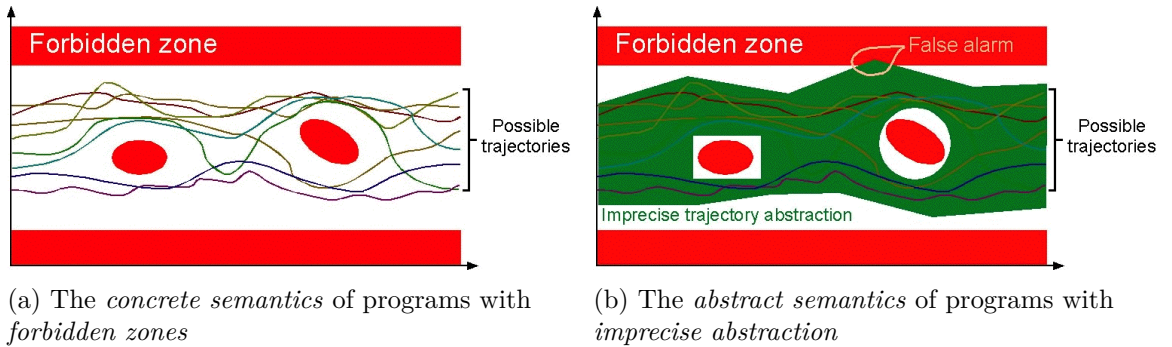


Figure 2.1: An illustration of *abstract interpretation* [9] (horizontal axes: the time t ; vertical axes: the vector $x(t)$ of input, state, and output values of the considered program)

2.2.1 Components of Abstract Interpretation

Before the formal definition of abstract interpretation is given, it is first intuitively described the basic components of the framework in accordance with [8, 10, 24, 31]:

- **Abstract Domain Q :**
 - An abstraction of the possible concrete program states (or their parts) in the form of *abstract properties*⁵.
 - In other words, it is a set of *abstract states* (or *abstract contexts*), where an abstract state represents a set of program states reachable at a particular program location.
 - Various practical domains have been defined — for example: numerical intervals, polyhedra, octagons, congruences, or various heap domains.
- **Abstract Transformers τ :**
 - There is a *transform function* τ for each program operation (instruction) representing the impact of the operation executed on an abstract state.
- **Join Operator \sqcup :**
 - Joins abstract states from individual program branches into a single one.
- **Widening Operator ∇ :**
 - Enforces *termination* of the abstract interpretation.
 - It is used to over-approximate the *least fixed points* of program semantics (it is performed on a sequence of abstract states at a certain location).
 - Usually, the later in the analysis this operator is applied, the more accurate the result is (but the analysis takes more time).
- **Narrowing Operator Δ :**
 - The approximation obtained by the widening operator can be *refined* using the narrowing operator, i.e., it may be used to refine the result of widening.
 - It is used when a fixpoint is approximated by the widening operator.

Note that neither the widening operator nor narrowing operator are required. However, at least the widening operator is frequently used. The narrowing operator can be sometimes dropped.

2.2.2 Formal Definition of Abstract Interpretation

The definitions below consider notions from *lattice theory*. More information about lattices, functions on lattices, *partial orders*, and set theory can be found, e.g., in [30, 31, 34].

⁵Abstract properties approximate *concrete properties* of behaviours [10].

Definition 2.2.1 (Abstract Interpretation). According to [11, 24], **abstract interpretation** I of a program P with the instruction set $Instr$ is a tuple

$$I = (Q, \sqcup, \sqsubseteq, \top, \perp, \tau)$$

where

- Q is the *abstract domain* (the set of *abstract states*),
- $\sqcup : Q \times Q \rightarrow Q$ is the *join operator* for accumulation of abstract states,
- $\sqsubseteq \subseteq Q \times Q$ is an *ordering* defined as $x \sqsubseteq y \iff x \sqcup y = y$ where
 - $\langle Q, \sqsubseteq \rangle$ is a *complete \sqcup -semilattice*,
- $\top \in Q$ is the *supremum* of $\langle Q, \sqsubseteq \rangle$,
- $\perp \in Q$ is the *infimum* of $\langle Q, \sqsubseteq \rangle$ — thus, $\langle Q, \sqsubseteq \rangle$ is, in fact, a *complete lattice*,
- $\tau : Instr \times Q \rightarrow Q$ defines the *abstract transformers* for specific instructions,
 - τ must be *monotone*⁶ on Q for each instruction from $Instr$.

Using a so-called *Galois connection* [11, 12, 24, 30, 31, 34], one can ensure *soundness* of abstract interpretation, i.e., the correspondence between the *concrete semantics* and its *abstract semantics* can be formalised by a Galois connection (also referred to as a pair of *adjointed functions*).

Definition 2.2.2 (Galois Connection). If $\mathcal{P} = \langle P, \preceq \rangle$ and $\mathcal{Q} = \langle Q, \sqsubseteq \rangle$ are *partially-ordered sets*, then a quadruple $\pi = (\mathcal{P}, \alpha, \gamma, \mathcal{Q})$ is a **Galois connection**, written

$$\langle P, \preceq \rangle \xrightleftharpoons[\alpha]{\gamma} \langle Q, \sqsubseteq \rangle$$

iff $\alpha : P \rightarrow Q$ and $\gamma : Q \rightarrow P$ are *monotone* functions such that

$$\forall p \in P : \forall q \in Q : \alpha(p) \sqsubseteq q \iff p \preceq \gamma(q)$$

where P is the concrete domain and Q is the abstract domain. Furthermore, α is an *abstraction function*, and γ is a *concretisation function*. Consequently, $\alpha(p)$ is the abstraction of p , i.e., the most precise approximation of $p \in P$ in Q , and $\gamma(q)$ is the concretisation of q , i.e., the most imprecise element of P which $q \in Q$ can soundly approximate. A Galois connection is illustrated in Figure 2.2.

If the abstraction and concretisation functions of the abstract interpretation form a Galois connection, then applying the abstraction function and concretising the result back yield a less (or equally) precise result. However, it is a safe (*conservative*) approximation. That is, formally written, $\forall p \in P : p \preceq \gamma(\alpha(p))$. Finally, the abstract interpretation may only over-approximate the concrete semantics (i.e., it is sound) if for each *concrete transformer* (instruction) $i : P \rightarrow P$ from the instruction set $Instr$, the appropriate *abstract transformer* $\tau : Instr \times Q \rightarrow Q$ respects the considered Galois connection, this is, $\forall p \in P : \alpha(i(p)) \sqsubseteq \tau(i, \alpha(p))$.

⁶Assume two *partially-ordered sets* $\langle E, \preceq_E \rangle$ and $\langle F, \preceq_F \rangle$, and a function $f : E \rightarrow F$. f is **monotone** iff $\forall x, y \in E : x \preceq_E y \implies f(x) \preceq_F f(y)$ [34].

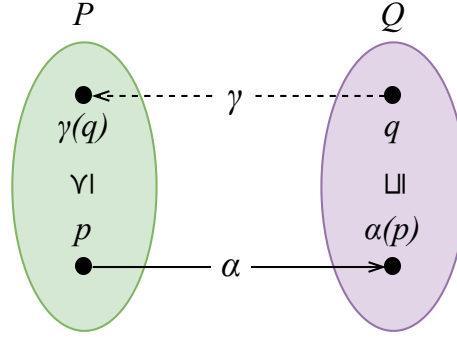


Figure 2.2: An illustration of a *Galois connection*

2.2.3 Fixpoint Approximation [11, 12, 24, 30, 31, 34]

Most program properties can be represented as *fixpoints*⁷. This reduces program analysis to *fixpoint approximation*. The complete analysis of a program using abstract interpretation can be then viewed as finding the *least/greatest* fixpoint of the equation $\bar{Q} = \bar{\tau}(\bar{Q})$, where \bar{Q} is a vector of abstract states and $\bar{\tau}$ is an extension of τ to the whole program. This fixpoint equation is then solved iteratively. *Knaster-Tarski theorem* [38] guarantees these fixpoints' existence (provided that $\bar{\tau}$ is monotone & complete lattices are considered).

The computation of the most precise abstract fixpoint is not generally guaranteed to *terminate*, in particular, when a given program contains a loop (or recursion) and uses an infinite domain (or even finite but very large). In order to enforce or accelerate the convergence, the fixpoint is often over-approximated using the *widening* operator ∇ . The approximation may be later refined using the *narrowing* operator Δ . These two operators are defined below. In practice, the analysis is usually done by iterating the abstract transformers over the *control-flow graph*⁸ (CFG). The join operator \sqcup is applied at program locations where different branches come across. Moreover, the widening operator ∇ is applied at loop junctions (and afterwards, the narrowing operator Δ may be used).

Definition 2.2.3 (Widening). Assume the abstract interpretation from Definition 2.2.1. The **widening operator** over the abstract domain Q is a binary operator $\nabla : Q \times Q \rightarrow Q$ such that $\forall x, y \in Q : x \sqcup y \sqsubseteq x \nabla y$, and for all *increasing chains* $x_0 \sqsubseteq x_1 \sqsubseteq \dots$, the increasing chain defined by $y_0 = x_0, \dots, y_{i+1} = y_i \nabla x_{i+1}, \dots$ is not *strictly increasing*. Note that the chain eventually stabilises since the result of ∇ is an *upper bound*.

Definition 2.2.4 (Narrowing). Assume the abstract interpretation from Definition 2.2.1. The **narrowing operator** over the abstract domain Q is a binary operator $\Delta : Q \times Q \rightarrow Q$ such that $\forall x, y \in Q : y \sqsubseteq x \implies y \sqsubseteq x \Delta y \sqsubseteq x$, and for all *decreasing chains* $x_0 \sqsupseteq x_1 \sqsupseteq \dots$, the decreasing chain defined by $y_0 = x_0, \dots, y_{i+1} = y_i \Delta x_{i+1}, \dots$ is not *strictly decreasing*. Note that the chain eventually stabilises since the result of Δ is a *lower bound*.

⁷Consider a function $f : L \rightarrow L$ on a set L . A **fixed point** (or **fixpoint**) of f is an element $l \in L$ such that $f(l) = l$. The set of all fixpoints of f is denoted as $Fix(f) = \{l \in L \mid f(l) = l\}$ [31].

⁸A **control-flow graph** (CFG) is a *directed graph* in which the nodes represent *basic blocks* of the program, and the edges represent control-flow paths among them [1]. A basic block is a *maximal sequence of statements* such that all transfers to it are to the first statement in the sequence, and all statements in the sequence are executed sequentially [31].

2.3 Facebook Infer — Static Analysis Framework

This section describes the principles and features of *Facebook Infer*. The description is based on information provided at the Facebook Infer’s website⁹ and in [4].

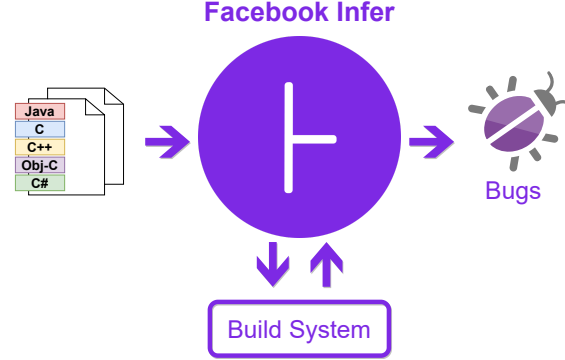


Figure 2.3: Static analysis in *Facebook Infer*

Facebook Infer is an *open-source*¹⁰ *static analysis framework*, which can discover various kinds of software bugs and which stresses the *scalability* of the analysis. The Facebook Infer’s basic usage is illustrated in Figure 2.3. A more detailed explanation of its architecture is given in Section 2.3.2. Facebook Infer is implemented in *OCaml*¹¹ — a *functional* programming language, also supporting *imperative* and *object-oriented* paradigms. Further details about OCaml can be found in the book [28]. Infer has initially been a rather specialised tool focused on *sound verification* of the absence of *memory safety violations*, which was first published in the well-known paper [7]. Once Facebook has purchased it, its scope significantly widened and abandoned the focus on sound analysis only.

Facebook Infer can analyse programs written in the following languages: C, C++, Java, Objective-C (and support for C# has been recently implemented [40]). Moreover, it is possible to extend Facebook Infer’s *frontend* for supporting other languages. Currently, Facebook Infer contains many analyses focusing on various kinds of bugs, e.g., *Inferbo* [42] (buffer overruns); *RacerD* [5, 6, 17] (data races); and other analyses that check for some forms of deadlocks and starvation, null-dereferencing, memory leaks, resource leaks, etc. Since Facebook Infer is a relatively popular and open-source framework, many experimental analysers arise pretty often. For instance, the promising experimental deadlock checker L2D2 [26] has been implemented at FIT BUT not long ago.

2.3.1 Abstract Interpretation in Facebook Infer

Facebook Infer is a general framework for static analysis of programs, and it is based on *abstract interpretation*. Despite the original approach taken from [7], Facebook Infer aims to find bugs rather than perform *formal verification*. It is still possible to propose *sound* and *complete* analyses in the framework. However, the majority of the checkers already implemented in Facebook Infer are both unsound and incomplete.

⁹Facebook Infer’s website: <https://fbinfer.com>.

¹⁰Facebook Infer’s open-source repository at GitHub: <https://github.com/facebook/infer>.

¹¹OCaml’s website: <https://ocaml.org>.

It can be used to develop new sorts of *compositional* and *incremental* analysers quickly (both *intraprocedural* and *interprocedural* [31]) based on the concept of function *summaries*. In general, a summary χ represents a set of *pre-condition/post-condition* pairs for a function. In particular, it records under which pre-condition a function can be performed leading to a given post-condition. Formally, it can be described using *Hoare triples* [20], because it can be viewed as a triple $\{P\} S \{R\}$, where P is a pre-condition, S is a program (or a single statement/command)¹², and R is a description of the result of the execution of S (i.e., a post-condition). In theory, P and R are considered formulae in a suitable logic. The triple may be interpreted as follows. Suppose S is executed from a state satisfying P , and the execution of S terminates. In that case, the program state after S terminates satisfies R . However, in practice, a summary is a custom data structure that may be used for storing any information resulting from the analysis of particular functions. Usually, a summary consists of the relevant parts of *abstract states* for a particular analysis.

Facebook Infer generally does not compute the summaries during the analysis along the CFG as it is done in classical analyses based on the concepts from [32, 36]. Instead, Facebook Infer performs the analysis of a program *function-by-function along the call tree*, starting from its leaves (demonstrated later in Example 2.3.1). Therefore, a function is analysed, and a summary is computed without knowledge of the call context. Then, the summary of the function is used at all its call sites. It is needed to deduce under which pre-conditions a function can produce post-conditions appropriate for the given analysis. Since the summaries do not differ for different contexts, each function is analysed precisely once, and the analysis becomes more scalable, but it can lead to a loss of accuracy. However, of course, it is more troublesome to propose such a *bottom-up* analysis.

In order to create a new intraprocedural analyser in Facebook Infer, it is required to define the following (the listed items are described in more detail in Section 2.2):

1. The *abstract domain* Q , i.e., the type of *abstract states*.
2. The *ordering operator* \sqsubseteq , i.e., an ordering of abstract states.
3. The *join operator* \sqcup , i.e., the way of joining two abstract states.
4. The *widening operator* ∇ , i.e., the way how to enforce termination of the computation.
5. The *transfer functions* τ , i.e., transformers that take an abstract state as an input and produce an abstract state as an output.

Furthermore, to create an interprocedural analyser, it is required to define additionally:

1. The type of function summaries χ .
2. The logic for using summaries in transfer functions and the logic for transforming an intraprocedural abstract state to a summary.

An important Facebook Infer’s feature, which improves its scalability, is the *incrementality* of the analysis. It allows one to analyse separate code changes only, instead of analysing the whole codebase. It is more suitable for extensive and variable projects where ordinary

¹²In Facebook Infer’s summaries considered as Hoare triples, S is usually one function.

analysis is not feasible. The incrementality is based on *re-using summaries* of functions for which there is no change in them neither in the functions transitively invoked from them, as shown in Example 2.3.1 later on.

2.3.2 Architecture of the Infer.AI Framework

The architecture of the abstract interpretation framework of Facebook Infer (**Infer.AI**) may be split into three major parts, as demonstrated in Figure 2.4: the *frontend*, an *analysis scheduler* (and the *results database*), and a set of *analyser plugins*.

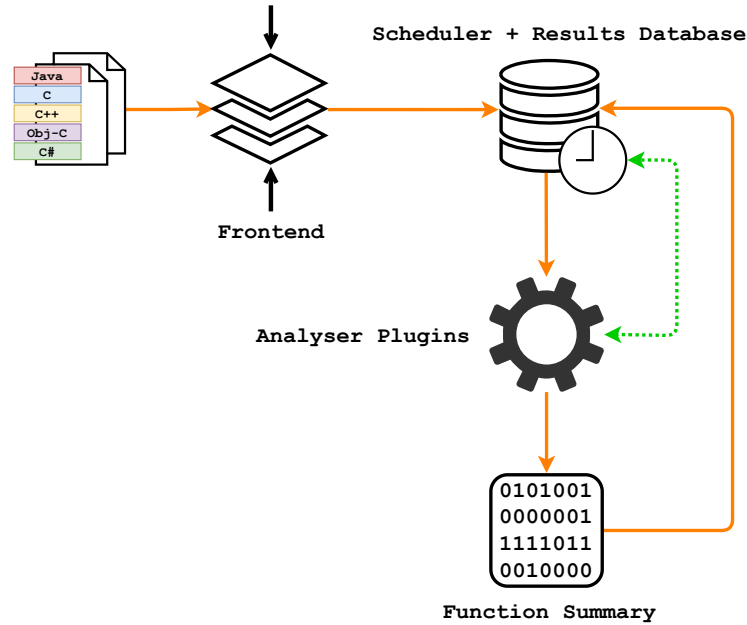


Figure 2.4: The *architecture* of the Facebook Infer's *abstract interpretation framework* [18]

The frontend compiles input programs into the *Smallfoot Intermediate Language* (SIL) and represents them as a CFG. There is a separate CFG representation for each analysed function. Nodes of this CFG are formed as SIL instructions. The SIL language consists of the following underlying instructions:

- **LOAD** — reading into a temporary variable;
- **STORE** — writing to a program variable, field of a structure, or array;
- **CALL** — a function call;
- **PRUNE e** — the evaluation of an expression **e** in a condition, cycle, etc.

The frontend allows one to propose *language-independent* analyses (to a certain extent) because it supports input programs to be written in multiple languages.

The next part of the architecture is the scheduler, which defines the order of the analysis of single functions according to the appropriate *call graph*¹³. The scheduler also checks if it is possible to simultaneously analyse some functions, allowing Facebook Infer to run the analysis in parallel.

Example 2.3.1. For demonstrating the order of the analysis in Facebook Infer and its incrementality, assume the call graph given in Figure 2.5. At first, leaf functions F5 and F6 are analysed. Further, the analysis goes on towards the root of the call graph—F_{MAIN}, while considering the dependencies denoted by the edges. This order ensures that a summary is available once a nested function call is abstractly interpreted within the analysis. When there is a subsequent code change, only directly changed functions and all the functions up the call path are re-analysed. For instance, if there is a change of source code of function F4, Facebook Infer triggers re-analysis of functions F4, F2, and F_{MAIN} only.

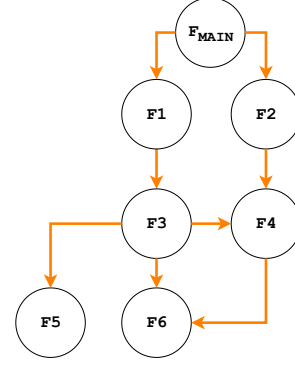


Figure 2.5: A *call graph* for an illustration of Facebook Infer’s *analysis process* [18]

The last part of the architecture consists of analyser plugins. Each plugin performs some analysis by interpreting SIL instructions. The result of the analysis of each function (function summary) is stored in the results database. The interpretation of SIL instructions (*commands*) is made using the *abstract interpreter* (also called the *control interpreter*) and *transfer functions* (also called the *command interpreter*). The transfer functions take a previously generated abstract state of an analysed function as an input, and by applying the interpreting command, produce a new abstract state. The abstract interpreter interprets the command in an abstract domain according to the CFG. This workflow is shown in a simplified form in Figure 2.6.

2.4 Contracts for Concurrency

This section introduces the concept of *contracts for concurrency* [13, 37]. Examples and listings in the section are pieces of programs written in the C language.

Respecting the *protocol* of a software module—that defines which *sequences of functions* are legal to invoke—is one of the requirements for the correct behaviour of the module. For example, a module that deals with a file system typically requires that a programmer using this module should call a function `open` at first, followed by an optional number of functions `read` and `write`, and at last, call a function `close`. A program utilising such a module that does not follow this protocol is erroneous. The methodology of *design by contract* [27] requires programs to meet such well-defined behaviours [37].

In *concurrent programs*, contracts for concurrency allow one—in the simplest case—to specify sequences of functions that are needed to be *executed atomically* in order to avoid

¹³A **call graph** is a *directed graph* describing call dependencies among functions.

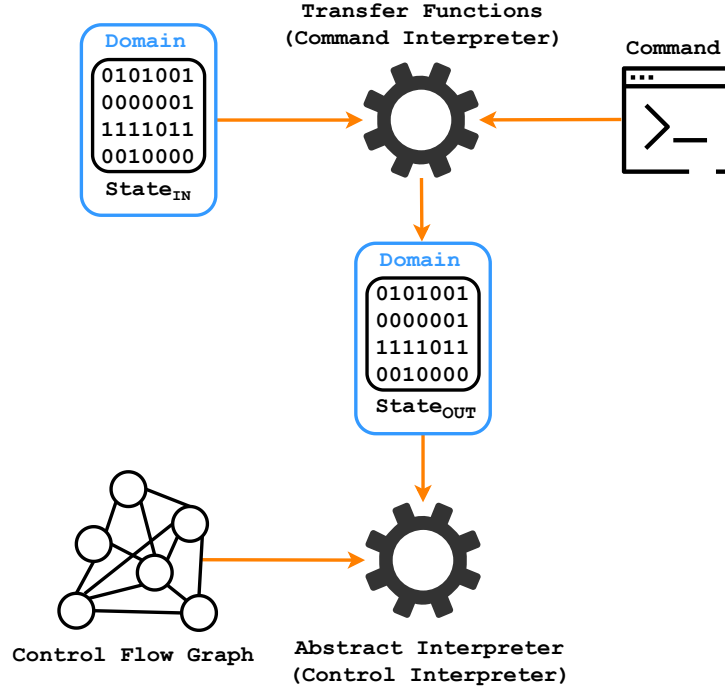


Figure 2.6: The Facebook Infer’s *abstract interpretation process* [18]

atomicity violations. In general, contracts for concurrency specify sets of sequences of calls called *spoilers* and sets of sequences of calls called *targets*. It is then required that no target overlaps fully with any spoiler. A developer may manually specify such contracts or automatically generate them by a program (analyser). These contracts can be used to verify the correctness of programs as well as they can serve as helpful documentation.

Section 2.4.1 defines the notion of *basic contracts for concurrency*. Further, Section 2.4.2 defines contracts extended to consider the *data flow* between functions, where a sequence of function calls must be atomic only if they handle the same data. Finally, Section 2.4.3 presents those mentioned above more general contracts for concurrency with spoilers and targets, which essentially extend the basic contracts with some *contextual information*.

2.4.1 Basic Contracts for Concurrency

In [13, 37], a *basic contract for concurrency* is formally defined as follows. Let $\Sigma_{\mathbb{M}}$ be a set of all function names of a software module \mathbb{M} . A contract is a set \mathbb{R} of *clauses*, where each clause $\varrho \in \mathbb{R}$ is a *star-free regular expression*¹⁴ over $\Sigma_{\mathbb{M}}$. A *contract violation* occurs if any of the sequences expressed by the contract clauses are interleaved with the execution of functions from $\Sigma_{\mathbb{M}}$. In other words, each sequence specified by any clause ϱ must be executed atomically. Otherwise, there is a violation of the contract. The number of sequences defined by a contract is finite since the contract is a union of *star-free languages*.

¹⁴**Star-free regular expressions** are regular expressions that use only the *concatenation operator* and the *alternative operator* ($|$), without the *Kleene star operator* ($*$).

Example 2.4.1. Consider the following example from [13, 37]. Assume a module \mathbb{M} implementing a resizable array of integers with the following interface functions $\Sigma_{\mathbb{M}} = \{\text{add}, \text{contains}, \text{index_of}, \text{get}, \text{set}, \text{remove}, \text{size}\}$ defined as:

```
void add(int *array, int element)

bool contains(int *array, int element)

int index_of(int *array, int element)

int get(int *array, int index)

void set(int *array, int index, int element)

void remove(int *array, int index)

int size(int *array)
```

The module’s contract contains clauses $\mathbb{R} = \{\varrho_1, \varrho_2, \varrho_3, \varrho_4\}$ such that:

(ϱ_1) `contains index_of`

The execution of `contains` followed by the execution of `index_of` should be atomic. Otherwise, the program may fail to get the index because after checking the presence of an element in an array, it can be removed by some concurrent thread.

(ϱ_2) `index_of (get | set | remove)`

The execution of `index_of` followed by the execution of `get`, `set`, or `remove` should be atomic. Otherwise, the received index may be outdated when applied to address an array element because a concurrent modification of the array may shift the element’s position.

(ϱ_3) `size (get | set | remove)`

The execution of `size` followed by the execution of `get`, `set`, or `remove` should be atomic. Otherwise, an array may be empty when accessing it because of a concurrent change of the array. This can be an issue since a given index is not in a valid range anymore (e.g., testing `index < size`).

(ϱ_4) `add index_of`

The execution of `add` followed by the execution of `index_of` should be atomic. Otherwise, the added element needs no longer exist in an array.

2.4.2 Contracts for Concurrency with Parameters

The above definition of basic contracts for concurrency is quite limited in some circumstances and can consider valid programs as erroneous (i.e., *false alarms* may be reported). Hence, in this section, there is introduced an extension of basic contracts—*contracts for concurrency with parameters* (defined in [13, 37])—which takes into consideration the *data flow* within function calls.

Example 2.4.2. Consider the following example from [13, 37], given in Listing 2.1. There is the function `replace` that replaces item `a` in an array by item `b`. The implementation of this function comprises two atomicity violations:

1. when `index_of` is invoked, item `a` does not need to be in the array anymore; and
2. the acquired index can be obsolete when `set` is executed.

A basic contract could cover this scenario by the following clause:

(ϱ_5) contains index_of set

It can be obtained from the composition of clauses ϱ_1 and ϱ_2 . Nevertheless, this definition is too restrictive because the functions are required to be executed atomically only if `contains` and `index_of` have the same arguments `array` and `element`; `index_of` and `set` have the same argument `array`; and the returned value of `index_of` is used as the argument `index` of the function `set`.

```

1 void replace(int *array, int a, int b)
2 {
3     if (contains(array, a))
4     {
5         int index = index_of(array, a);
6         set(array, index, b);
7     }
8 }

```

Listing 2.1: An example of an atomicity violation with *data dependencies* [13, 37]

To respect function call *parameters* and *return values* of functions in contracts, the basic contracts are extended by dependencies among functions in [13, 37] as follows. Function call parameters and return values are expressed as *meta-variables*. Further, if a contract is required to be respected exclusively if the same object emerges as an argument or as the return value of multiple calls in a given call sequence, it may be denoted by using the same meta-variable at the position of all these occurrences of parameters and return values.

Clause ϱ_5 can be extended as follows (repeated application of meta-variables `X/Y/Z` requires the same objects $o_1/o_2/o_3$ to be used at the positions of `X/Y/Z`, respectively; and the underscore indicates a *free meta-variable* that does not restrict the contract clause):

(ϱ'_5) contains(X, Y) Z=index_of(X, Y) set(X, Z, _)

Example 2.4.3. With the extension described above, it is possible to extend the contract from Example 2.4.1 to capture the dependencies between function calls as follows:

(ϱ'_1) contains(X, Y) index_of(X, Y)

(ϱ'_2) Y=index_of(X, _) (get(X, Y) | set(X, Y, _) | remove(X, Y))

(ϱ'_3) size(X) (get(X, _) | set(X, _, _) | remove(X, _))

(ϱ'_4) add(X, Y) index_of(X, Y)

2.4.3 Contracts for Concurrency with Spoilers

Interleaving a sequence of function calls from a contract clause (considering the basic contracts for concurrency or contracts with parameters) with some function calls of a given module can cause atomicity violations. At the same time, this is not the case for some other function calls from the module. However, this is not possible to describe in the contracts introduced so far. For instance, clause (ϱ_1) **contains index_of** requires that this sequence of calls must always be performed atomically, i.e., it does not matter which functions are executed by other threads. Thus, interleaving the execution of this sequence with, e.g., **remove** or **get** is a contract violation. However, in effect, only the execution of **remove** may be problematic, while the execution of **get** may not.

The paper [13] proposes a solution to the above issue — an extension of the basic contracts for concurrency with *contextual information* — allowing one to describe in which context the contract clauses should be enforced. Each clause of the basic contract is now called *target*. For each target, there is a set of so-called *spoilers* that restrict its application. That is, a spoiler is a sequence of function calls that can violate its target. In the end, it has to be ensured that each target is executed atomically w.r.t. its spoilers. Assuming the earlier example, let **contains index_of** be the target clause. Then, a possible spoiler is **remove**. The syntax for this description is as follows: **contains index_of** \Leftarrow **remove**.

Formally, let $\Sigma_{\mathbb{M}}$ be a set of all function names of a software module \mathbb{M} . Further, let \mathbb{R} be a set of *target clauses*, where each clause $\varrho \in \mathbb{R}$ is a star-free regular expression over $\Sigma_{\mathbb{M}}$. Let \mathbb{S} be a set of *spoiler clauses*, where each clause $\sigma \in \mathbb{S}$ is a star-free regular expression over $\Sigma_{\mathbb{M}}$. Moreover, let $\Sigma_{\mathbb{R}} \subseteq \Sigma_{\mathbb{M}}$ and $\Sigma_{\mathbb{S}} \subseteq \Sigma_{\mathbb{M}}$ be the alphabets of function names used in the target or spoiler clauses, respectively. Then, a contract is a relation $\mathbb{C} \subseteq \mathbb{R} \times \mathbb{S}$, which states for each target the spoilers that can cause atomicity violations. One spoiler may violate more than one target, and, on the contrary, one target may be violated by more than one spoiler. A contract is violated iff any executed sequence expressed by a target $\varrho \in \mathbb{R}$ is *completely interleaved* with the execution of the sequence representing its spoiler, i.e., $\sigma \in \mathbb{C}(\varrho)$. A target sequence r is completely interleaved by a spoiler sequence s iff the execution of r *starts before* the execution of s , and the execution of s *ends before* the execution of r . *Partial interleavings* of targets and spoilers are here not taken into account as causing an error. However, if needed, this can be resolved by adding a new target clause with an appropriate spoiler. Whether a sequence is executed before another one, is defined using the so-called “*happened before*” relation ($\xrightarrow{\text{hb}}$) [23]. A complete interleaving is illustrated in Figure 2.7a, and partial interleavings are shown in Figures 2.7b and 2.7c.

Example 2.4.4. With the extension of spoilers, the contract from Example 2.4.1 can be refined with contextual information to refuse *unsafe interleavings* as follows (all other interleavings, not described by spoilers, are deemed safe):

```
( $\varrho_1''$ ) contains index_of  $\Leftarrow$  remove | set
( $\varrho_2''$ ) index_of (get | set | remove)  $\Leftarrow$  remove | set
( $\varrho_3''$ ) size (get | set | remove)  $\Leftarrow$  remove
( $\varrho_4''$ ) add index_of  $\Leftarrow$  remove | set
```

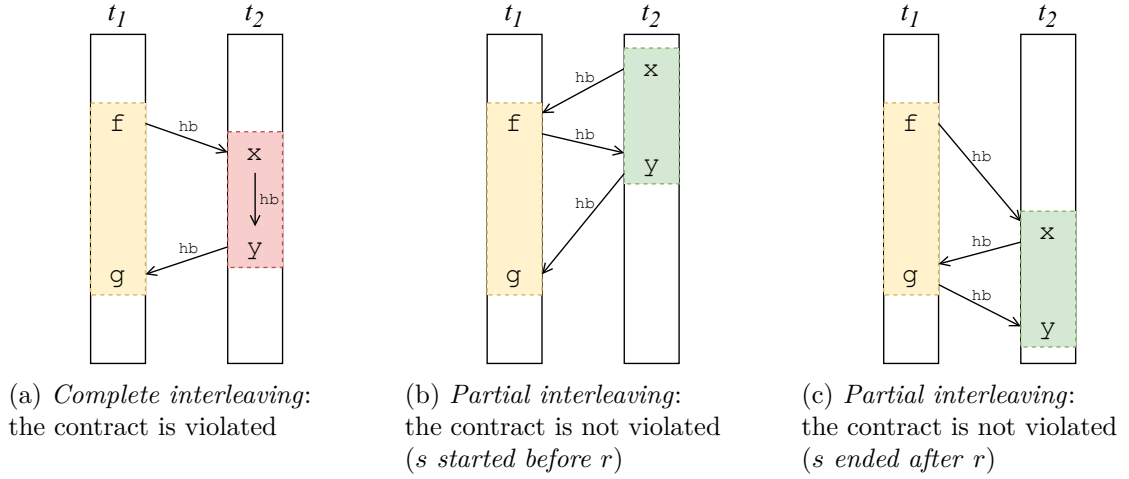


Figure 2.7: A *contract violation* demonstration with *target* and *spoiler interleavings*. In each sub-figure, a thread t_1 executes a target sequence $r = \mathbf{f} \ \mathbf{g}$ and a thread t_2 executes a spoiler sequence $s = \mathbf{x} \ \mathbf{y}$. Consequently, the following contract is considered: $\mathbf{f} \ \mathbf{g} \leftarrow \mathbf{x} \ \mathbf{y}$ [29]

Lastly, it is possible to combine the extension of contracts for concurrency with spoilers and with parameters. The following clause can demonstrate it:

$$\mathbf{Y} = \mathbf{index_of}(\mathbf{X}, _) \ \mathbf{get}(\mathbf{X}, \mathbf{Y}) \leftarrow \mathbf{remove}(\mathbf{X}, _)$$

The clause requires sequences of `index_of` and `get` to be performed atomically, but only when working with the same element \mathbf{X} and index \mathbf{Y} , and only w.r.t. the concurrent execution of `remove`. Note that the argument `index` of `remove` is not restricted because any concurrent removal may produce an atomicity violation — by either removing an element on index \mathbf{Y} or by shifting its position.

Chapter 3

Atomer — Atomicity Violations Detector

This chapter describes the principles and limitations of the basic version of the *Atomer static analyser* proposed as a module of *Facebook Infer* (introduced in Section 2.3) for finding some forms of *atomicity violations*. *Atomer* was proposed and in detail described in the bachelor’s thesis [18] of the author of this thesis. Therefore, naturally, the description in Section 3.2 is based on the mentioned thesis. Already existing solutions in this area (besides *Atomer*) are discussed in Section 3.1. In particular, it deals with other existing approaches and tools for finding atomicity violations, their advantages, disadvantages, features, availability, and so on. In Section 3.3, there are discussed limitations and shortcomings of *Atomer*. Some of the thoughts mentioned in this section are taken into consideration already in [18].

3.1 Related Work

Atomer is slightly inspired by ideas from [13, 37]. In these papers, there is described a proposal and implementation of a *static* approach for finding *atomicity violations of sequences of function calls*, which is based on *grammars* and *parsing trees*. Note that in the paper [13], there is also described and implemented a *dynamic* approach for the validation. The authors of these papers implemented a stand-alone prototype static analyser called *Gluon*¹ for analysing programs written in Java. To the best author’s knowledge, *Gluon* is the only static analyser that tries to go in a similar direction as *Atomer* does. *Gluon* led to some promising experimental results, but the *scalability* of *Gluon* was still limited.² Moreover, *Gluon* is no more actively developed, and it is not easy to use. Despite all author’s efforts, it was not put into operation. Above that, the authors themselves note that the code of *Gluon* is very ad hoc, and many things are hard-coded in it. These facts, in fact, inspired the decision that eventually led to the implementation of the first version of *Atomer*, namely, to get inspired by the ideas of [13, 37] but reimplement them in *Facebook Infer*, redesigning it in accordance with the principles of *Facebook Infer*, which should make

¹**Gluon** is a tool that implements a *static* approach for finding *atomicity violations of sequences of function calls* in Java programs. It is available at <https://github.com/trxsys/gluon>.

²Some of the experiments performed by *Gluon* are also similarly performed by *Atomer*, which is discussed in Sections 6.2, 6.3.

the resulting tool more scalable. In the end, however, due to adapting the analysis to the context of Facebook Infer, the implementation of Atomer’s analysis is significantly different from [13, 37], as is presented in Chapter 4 of [18]. Furthermore, unlike Gluon, a new version of Atomer is capable of analysing a much wider range of programs because it also supports other languages than Java, and it supports more advanced *locking mechanisms*. On the other hand, Gluon implements *extended contracts for concurrency* (see Sections 2.4.2, 2.4.3) that consider *data flow* within functions and *contextual information*, while Atomer implements only the idea of *basic contracts for concurrency* (see Section 2.4.1). These extended contracts should improve the precision of the analysis. Nonetheless, it is the author’s future work to implement the extended contracts in Atomer as well.

In Facebook Infer, there is already implemented analysis called *Lock Consistency Violation*³. It is a part of *RacerD* [5, 6, 17]. This analysis finds atomicity violations in C++ and Objective-C programs for *reads/writes* on single variables required to be executed atomically. Atomer is different; it finds atomicity violations for *sequences of functions* required to be executed atomically. Moreover, Atomer tries to automatically determine which sequences should indeed be executed atomically.

3.2 Design of Atomer and Its Principles

Atomer concentrates on checking the *atomicity of the execution of certain sequences of function calls*, which is often required for *concurrent programs*’ correct behaviour. In principle, Atomer is based on the assumption that sequences of function calls executed *atomically once* should probably be executed *always atomically*.

The proposal of Atomer is based on the concept of *contracts for concurrency* described in Section 2.4. In particular, the proposal considers the *basic contracts* described in Section 2.4.1. Neither the contracts extended by *parameters* explained in Section 2.4.2 nor the contracts extended by *spoilers* and *targets* described in Section 2.4.3 are so far taken into account.

In general, basic contracts for concurrency allow one to define sequences of functions required to be executed atomically, as explained in more detail in Section 2.4. Atomer is able to automatically derive candidates for such contracts and then verify whether the contracts are fulfilled. In other words, Atomer can both automatically derive those sequences that are sometimes executed atomically as well as subsequently check whether they are indeed always executed atomically. Both of these steps are done statically. The proposed analysis is thus divided into two parts (*phases of the analysis* that are in-depth described in the sections below and illustrated in Figure 3.1):

Phase 1: Detection of (likely) *atomic sequences*.

Phase 2: Detection of *atomicity violations* (violations of the atomic sequences).

This section provides a high-level view of the *abstract interpretation* underlying Atomer. The concrete types of the *abstract states* (i.e., elements of the *abstract domain Q*) and

³**Lock Consistency Violation** is an *atomicity violations* analysis implemented in *Facebook Infer*. It is described at https://fbinfer.com/docs/all-issue-types#lock_consistency_violation.

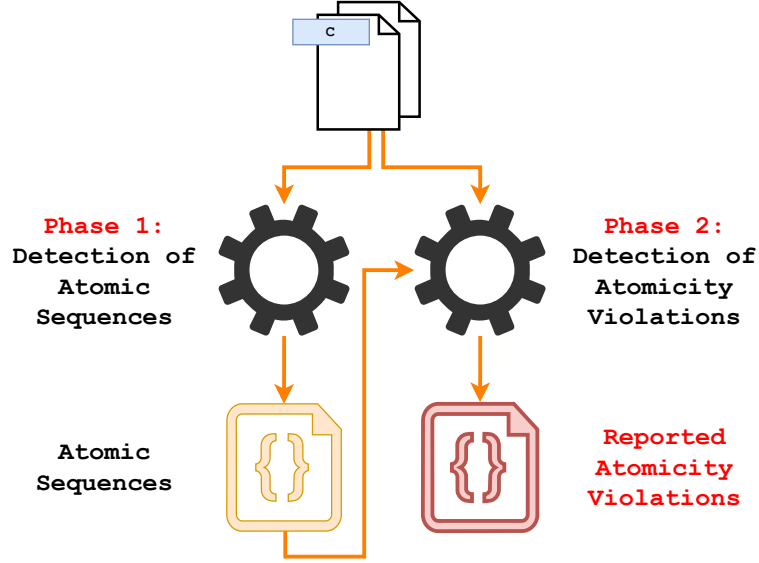


Figure 3.1: *Phases of the analysis of the Atomer’s first version and the analysis high-level process illustration* [18]

the *summaries* χ , along with the implementation of all necessary *abstract interpretation operators* are stated in Chapter 4 of [18]. However, actually, the abstract states $s \in Q$ of both phases of the analysis are proposed as *sets*. So, in fact, the *ordering operator* \sqsubseteq is implemented using testing for a *subset* (i.e., $s \sqsubseteq s' \iff s \subseteq s'$, where $s, s' \in Q$), the *join operator* \sqcup is implemented as the *set union* (i.e., $s \sqcup s' \iff s \cup s'$), and the *widening operator* ∇ is implemented using the join operator (i.e., $s \nabla s' \iff s \sqcup s'$) since the domains are *finite*. Furthermore, it is essential to note that the proposed analysis is neither *sound* nor *complete*. Its goal is to effectively find bugs rather than formally verify the absence of atomicity violations.

Function summaries are in the below sections reduced to the output parts only (*post-conditions* R). The input parts of summaries (*pre-conditions* P) are in case of the proposed analysis always empty because, so far, it is not necessary to have any pre-conditions for analysed functions. Thus, in this case, *Hoare triples* — *true* S $\{R\}$ — are used, where S is an analysed function, i.e., $P = \text{true}$.

Listings in the below sections are pieces of programs written in the C language (assuming *lock/unlock* functions for *mutual exclusion* to *critical sections*).

3.2.1 Phase 1 — Detection of Atomic Sequences

Before detecting *atomicity violations* may begin, it is required to have *contracts for concurrency* introduced earlier. Phase 1 of Atomer is able to produce such contracts, i.e., it detects *sequences of functions* that should be *executed atomically*. Intuitively, the detection is based on looking for sequences of functions executed atomically, in particular, under some *lock*, on some path through a program. The assumption is that if it is *once needed to execute a sequence atomically*, it should probably be *always executed atomically*.

For a description of the analysis, it is first needed to introduce a notion of a *reduced sequence* of function calls. Such a sequence denotes a sequence in which the *first appearance* of each function is recorded only. It is needed to ensure the *finiteness* of the sequences derived by the analysis, and hence the analysis's *termination*. The detection of sequences of calls to be executed atomically is based on analysing all paths through the CFG of a function and generating all pairs $(A, B) \in \Sigma^* \times \Sigma^*$ (where Σ is the set of functions of a given program) of reduced sequences of function calls for each path such that: A is a reduced sequence of function calls that appear between the beginning of the function being analysed and the first lock; between an unlock and a subsequent lock; or between an unlock and the end of the function being analysed. B is a reduced sequence of function calls that follow the calls from A , and that appear between a lock and unlock (or between a lock and the end of the function being analysed). Thus, the abstract states of the analysis are elements of the set $2^{2^{\Sigma^* \times \Sigma^*}}$ because there is a set of the (A, B) pairs for each program path.

It would be more precise to generate longer sequences of type $A_1 \cdot B_1 \cdot A_2 \cdot B_2 \dots$ instead of the sets of the pairs (A, B) . Nevertheless, it would be more challenging to ensure the above longer sequences' finiteness and the sets of these sequences' finiteness. Moreover, there would be a significantly larger *state space explosion problem* [39]. So, the proposed representation of the sets of pairs of sequences has been chosen to compromise accuracy and efficiency. However, the experiments described in Chapter 5 of [18] show that it needs even more pronounced abstraction for appropriate *scalability*.

Formally, the *initial abstract state* of a function is defined as $s_{init} = \{ \{(\varepsilon, \varepsilon)\} \}$, where ε indicates an empty sequence. To formalise the analysis of a function, let \mathbf{f} be a called leaf function. Further, let s_g be the abstract state of a function \mathbf{g} being analysed before the function \mathbf{f} is called. After the call of \mathbf{f} , the abstract state will be changed as follows:

$$s_g = \{ p' \in 2^{\Sigma^* \times \Sigma^*} \mid \exists p \in s_g : p' = \{ (A', B') \in \Sigma^* \times \Sigma^* \mid \exists (A, B) \in p : \\ \textcolor{blue}{[\neg actual(p, (A, B)) \wedge (A', B') = (A, B)]} \vee \textcolor{red}{[actual(p, (A, B))]} \\ \wedge \textcolor{green}{[(lock \wedge (A', B') = (A, B \cdot \mathbf{f})) \vee (\neg lock \wedge (A', B') = (A \cdot \mathbf{f}, B))]} \} \}$$

where *actual* is a Boolean function that determines whether a given (A, B) pair is the most recent pair of sequences of the current program state for a given program path. Furthermore, *lock* is a predicate indicating whether the current program state is inside an *atomic block*. Further, let s_g be the abstract state of a function \mathbf{g} being analysed before an unlock is called. After the unlock is called, a new (A, B) pair is created and labelled as an actual using the function *setActual* as follows:

$$s_g = \{ p' \in 2^{\Sigma^* \times \Sigma^*} \mid \exists p \in s_g : p' = \{ (A, B) \in \Sigma^* \times \Sigma^* \mid \\ \textcolor{red}{[(A, B) = (\varepsilon, \varepsilon) \wedge setActual(p, (A, B))]} \vee (A, B) \in p \} \}$$

Example 3.2.1. For an explanation of the computation of the sets of the pairs (A, B) , assume that a state of the analysis of a program S is the following sequence of function calls: $\mathbf{f} \cdot \mathbf{g}$; and a state of the analysis of a program S' is the following sequence of function calls: $\mathbf{f} \cdot \mathbf{g} \text{ [} \mathbf{m} \cdot \mathbf{n}$. The square brackets are used to indicate an *atomic sequence* (the closing square bracket is missing in the case of the program S' , which means that the program state is currently inside an atomic block). The computed abstract state for the program S is $s_S = \{ \{ (\mathbf{f} \cdot \mathbf{g}, \varepsilon) \} \}$, and for the program S' , it is $s_{S'} = \{ \{ (\mathbf{f} \cdot \mathbf{g}, \mathbf{m} \cdot \mathbf{n}) \} \}$. Now, if the next instruction is a call of a function \mathbf{x} , in the case of the program S , the call will be added to the first A sequence, and in the case of the program S' , the call will be added to the

first B sequence as follows: $s_S = \{\{(f \cdot g \cdot x, \varepsilon)\}\}$, $s_{S'} = \{\{(f \cdot g, m \cdot n \cdot x)\}\}$. Subsequently, if the next step in the program S is a lock call, the next function calls will be added to the first B sequence of the set s_S . And if the next step in the program S' is an unlock call, it will be created a new element of the first set of the set $s_{S'}$, and the next function calls will be added to the A sequence of this element. Finally, if a function y is called, the resulting sets will look like follows: $s_S = \{\{(f \cdot g \cdot x, y)\}\}$, $s_{S'} = \{\{(f \cdot g, m \cdot n \cdot x), (y, \varepsilon)\}\}$. Note that the final sequences of function calls look like follows: $f \cdot g \cdot x [y$ and $f \cdot g [m \cdot n \cdot x] y$ for the programs S and S' , respectively.

A *summary* $\chi_f \in 2^{\Sigma^*} \times \Sigma^*$ of a function f is then a pair $\chi_f = (\mathcal{B}, C)$, where:

- \mathcal{B} is a set constructed such that it contains all the B sequences that appear on program paths through f , i.e., those computed within the (A, B) pairs at the exit of f . Formally, $\mathcal{B} = \{B' \in \Sigma^* \mid \exists p \in s_f : \exists (A, B) \in p : B \neq \varepsilon \wedge B' = B\}$, where s_f is the abstract state at the end of the abstract interpretation of f . In other words, this component of the summary is a set of sequences of atomic function calls appearing in f .
- C is a *concatenation* of all the A and B sequences with removed duplicates of function calls. In particular, assume that the following set of (A, B) pairs is computed at the exit of f : $\{(A_1, B_1), (A_2, B_2), \dots, (A_n, B_n)\}$, then the result is the sequence $A_1 \cdot B_1 \cdot A_2 \cdot B_2 \cdot \dots \cdot A_n \cdot B_n$ with removed duplicates. Formally,

$$C = \text{reduce}(\bigcup_{c \in C'} c)$$

where $C' = \{c \in \Sigma^* \mid \exists p \in s_f : \exists (A, B) \in p : c = A \cdot B\}$, \bigcup concatenates all sequences of a set, and *reduce* is a function that removes duplicates of function calls. Intuitively, in this component of the summary, the analysis gathers occurrences of all called functions within the analysed function obtained by a concatenation of all the A and B sequences. C is recorded to facilitate the derivation of atomic call sequences that show up higher in the *call hierarchy*. Indeed, while locks/unlocks can appear in such a *higher-level* function, parts of the call sequences can appear lower in the call hierarchy.

Example 3.2.2. For instance, the analysis of the function f from Listing 3.1 produces the following sequences:

$$\overbrace{x \cdot \bar{x} \cdot y}^{A_1} \overbrace{[a \cdot b \cdot b]}^{B_1} \overbrace{y \cdot \bar{y}}^{A_2} \overbrace{[a \cdot c]}^{B_2} \overbrace{\bar{y}}^{A_3} \overbrace{[a \cdot a \cdot c]}^{B_3}$$

The functions a , b , c , x , y are not deeper analysed because it is assumed that these functions are leaf nodes of the call graph. The strikethrough of the functions b , x , y denotes removing already recorded function calls in the A and B sequences to get the reduced form. The strikethrough of the entire sequence $y [a \cdot a \cdot c]$ means discarding sequence already seen before. For the above, the abstract state at the end of the abstract interpretation of the function f is $s_f = \{\{(x \cdot y, a \cdot b), (y, a \cdot c), (\varepsilon, \varepsilon)\}\}$. The derived summary χ_f for the function f is $\chi_f = (\mathcal{B}, C)$, where $\mathcal{B} = \{a \cdot b, a \cdot c\}$, i.e., B_1 and B_2 ; and $C = x \cdot y \cdot a \cdot b \cdot c$, i.e., the concatenation of A_1, B_1, A_2, B_2 from which duplicate function calls were removed.

```

1 void f()
2 {
3     x(); x(); y();
4     lock(&L); // a . b
5     a(); b(); b();
6     unlock(&L);
7     y(); y();
8     lock(&L); // a . c
9     a(); c();
10    unlock(&L);
11    y();
12    lock(&L); // a . c
13    a(); a(); c();
14    unlock(&L);
15 }

```

Listing 3.1: A code snippet used for an illustration of the derivation of *sequences of functions called atomically*

Further, it is demonstrated how the results of the analysis of *nested functions* are used during the detection of atomic sequences. The result of the analysis of a nested function is used as follows. When calling an already analysed function, one plugs the sequence from the second component of its summary (i.e., the C sequence) into the most recent A or B sequence of all the program paths (where a program path corresponds to a single element of an abstract state, i.e., a set of the (A, B) pairs). In particular, assume that (A, B) is the most recent pair of sequences of the program state of a path being analysed. Subsequently, it is called a function f with a non-empty summary (i.e., $C \neq \varepsilon$). If the current program state of the analysed function is inside an atomic block, the analysis in this step will transform the pair (A, B) to a new (A', B') pair as follows: $(A', B') = (A, B \cdot f \cdot C)$. Otherwise, $(A', B') = (A \cdot f \cdot C, B)$. In such cases where a summary is empty, i.e., there are no function calls in a called function, or it is a leaf node of the call graph, just the function name is appended to the most recent A or B sequences of all the program paths. To formalise this process, let f be a called function that was already analysed, and the second component of its summary is C . Further, let s_g be the abstract state of a function g being analysed before the function f is called. After the call of f , the abstract state will be changed as follows:

$$\begin{aligned}
s_g = \{p' \in 2^{\Sigma^* \times \Sigma^*} \mid \exists p \in s_g : p' = \{ & (A', B') \in \Sigma^* \times \Sigma^* \mid \exists (A, B) \in p : \\
& [\neg actual(p, (A, B)) \wedge (A', B') = (A, B)] \vee [actual(p, (A, B)) \\
& \wedge [(lock \wedge (A', B') = (A, B \cdot f \cdot C)) \vee (\neg lock \wedge (A', B') = (A \cdot f \cdot C, B))]\} \}
\end{aligned}$$

Example 3.2.3. This example shows how the function g from Listing 3.2 would be analysed using the result of the analysis of the function f from Listing 3.1. The second component of χ_f is $C = x \cdot y \cdot a \cdot b \cdot c$. The analysis of the function g produces the following sequence:

$$x \cdot f \cdot x \cdot y \cdot a \cdot b \cdot c \cdot z \ [f \cdot x \cdot y \cdot a \cdot b \cdot c]$$

For the above, the abstract state at the end of the abstract interpretation of the function g is $s_g = \{ \{(x \cdot f \cdot y \cdot a \cdot b \cdot c \cdot z, f \cdot x \cdot y \cdot a \cdot b \cdot c), (\varepsilon, \varepsilon)\} \}$. The derived summary χ_g for the function g is $\chi_g = (\{f \cdot x \cdot y \cdot a \cdot b \cdot c\}, x \cdot f \cdot y \cdot a \cdot b \cdot c \cdot z)$.

```

1 void g()
2 {
3     x(); f(); z();
4     lock(&L); // f . x . y . a . b . c
5     f();
6     unlock(&L);
7 }

```

Listing 3.2: A code snippet used to illustrate the derivation of sequences of functions called atomically with a *nested function call* (function `f` is defined in Listing 3.1)

Cases Where Lock/Unlock Calls Are Not Paired in a Function

For treating cases where *lock/unlock calls are not paired* in a function—as demonstrated in Listing 3.3—the following solution is implemented in the basic version of Atomer.

Everything is unlocked at the end of a function, i.e., one *virtually appends* an unlock to the end of the function if it is necessary. Then, for the function `m` from Listing 3.3, the atomic section is virtually closed. Hence, there is detected an atomic sequence `a`. In particular, the summary is as follows: $\chi_m = (\{a\}, a)$.

Moreover, *all unlock calls not preceded by a lock are ignored*. Thus, in the function `n` from Listing 3.3, there are not detected any atomic sequences: $\chi_n = (\emptyset, a)$.

```

1 void m()
2 {
3     lock(&L); // a
4     a();
5 }
6 void n()
7 {
8     a();
9     unlock(&L);
10 }

```

Listing 3.3: A code snippet used to illustrate treating cases where *lock/unlock calls are not paired* in a function

Summary of Phase 1 The derived sequences of calls assumed to execute atomically—the \mathcal{B} sequences—from the summaries of all analysed functions are stored into a file used during Phase 2, which is described later on.

3.2.2 Phase 2—Detection of Atomicity Violations

In the second phase of the analysis, i.e., when *detecting violations* of the atomic sequences obtained from Phase 1, the analysis looks for *pairs of functions* that *should be called atomically* (or just for single functions if there is only one function call in an atomic sequence) and

that are not executed atomically (i.e., under a lock) on some path through the CFG. The pairs of function calls to be checked for atomicity are obtained as follows. For each function f with a Phase 1 summary $\chi_f = (\mathcal{B}, C)$ in a given program S , where $\mathcal{B} = \{B_1, B_2, \dots, B_n\}$, the analysis considers *every pair* $(x, y) \in \Sigma \times \Sigma$ of functions that appear as a *substring* in some of the B_i sequences, i.e., $B_i = w \cdot x \cdot y \cdot w'$ for some sequences w, w' . Note that x could be ε (an empty sequence) if some B_i consists of a single function. All these “atomic pairs” are put into the set $\Omega \in 2^{\Sigma \times \Sigma}$. More formally,

$$\Omega = \{(x, y) \in \Sigma \times \Sigma \mid \exists (\mathcal{B}, C) \in X_S : \exists B \in \mathcal{B} : [|B| = 1 \wedge (x, y) = (\varepsilon, B)] \\ \vee [|B| > 1 \wedge \exists w, w' \in \Sigma^* : B = w \cdot x \cdot y \cdot w' \wedge (x, y) \neq (\varepsilon, \varepsilon)]\}$$

where $X_S \in 2^{2^{\Sigma^*} \times \Sigma^*}$ is a set of all Phase 1 summaries of the program S .

Example 3.2.4. For instance, assume that in Phase 1, there was analysed a function f . It produced the summary $\chi_f = (\mathcal{B}, C)$, where $\mathcal{B} = \{a \cdot b \cdot c, a \cdot c \cdot d\}$, i.e., a set of sequences of functions that should be called atomically. The analysis will then look for the following pairs of functions that are not called atomically: $\Omega = \{a \cdot b, b \cdot c, a \cdot c, c \cdot d\}$.

An element of this phase’s abstract state is a triple $(x, y, \delta) \in \Sigma \times \Sigma \times 2^{\Sigma \times \Sigma \times \mathbb{N}}$ where (x, y) is a pair of the most recent calls of functions performed on the program path being explored, and δ is a set of so far detected *pairs that violate atomicity* on particular lines of code. Thus, the abstract states are elements of the set $2^{\Sigma \times \Sigma \times 2^{\Sigma \times \Sigma \times \mathbb{N}}}$. Whenever a function f is called on some path that led to an abstract state (x, y, δ) , a new pair (x', y') of the most recent function calls is created from the previous pair (x, y) such that $(x', y') = (y, f)$. Further, when the current program state is not inside an atomic block, the analysis checks whether the new pair (or just the last call) violates atomicity (i.e., $(x', y') \in \Omega \vee (\varepsilon, y') \in \Omega$). When it does, it is added to the set δ of pairs that violate atomicity.

Formally, the initial abstract state (in this phase) of a function is defined as $s_{init} = \{(\varepsilon, \varepsilon, \emptyset)\}$. To formalise the analysis of a function, let f be a called leaf function on a line $c \in \mathbb{N}$. Further, let s_g be the abstract state of a function g being analysed before the function f is called. After the call of f , the abstract state will be changed as follows:

$$s_g = \{(x', y', \delta') \in \Sigma \times \Sigma \times 2^{\Sigma \times \Sigma \times \mathbb{N}} \mid \exists (x, y, \delta) \in s_g : (x', y') = (y, f) \\ \wedge [(\neg lock \wedge \delta' = \{(x'', y'', c') \in \Sigma \times \Sigma \times \mathbb{N} \mid (x'', y'', c') \in \delta \\ \vee [(x'', y'') = (x', y') \vee (x'', y'') = (\varepsilon, y')] \wedge (x'', y'') \in \Omega \wedge c' = c\})] \vee (lock \wedge \delta' = \delta)]\}$$

A *summary* of the second phase of the analysis — $\chi_f \in 2^{\Sigma \times \Sigma \times \mathbb{N}}$ — of a function f is then $\chi_f = \Delta$, where Δ is a set of pairs that violate atomicity within the function f . Δ is constructed such that it contains a union of all the δ sets that appear on program paths through f . Formally, $\Delta = \bigcup_{\delta' \in \Delta'} \delta'$, where $\Delta' = \{\delta' \in 2^{\Sigma \times \Sigma \times \mathbb{N}} \mid \exists p \in s_f : \exists (x, y, \delta) \in p : \delta \neq \emptyset \wedge \delta' = \delta\}$ and s_f is the abstract state at the end of the abstract interpretation of f .

The analysis of functions with *nested function calls* and cases where *lock/unlock calls are not paired* within functions are handled analogically as in Phase 1. For a detailed explanation, see [18].

Example 3.2.5. To demonstrate the detection of an atomicity violation, assume the functions f and g from Listing 3.4. The set of atomic sequences of the function f with the first

phase's summary $\chi_f = (\mathcal{B}, C)$ is $\mathcal{B} = \{a \cdot b \cdot c\}$, thus, $\Omega = \{(a, b), (b, c)\}$. In the function g , an atomicity violation is detected because the pair of functions b, c is not called atomically on line 12, i.e., $(b, c) \in \Omega$. Consequently, the derived summary χ_g for the function g for the second phase of the analysis is $\chi_g = \{(b, c, 12)\}$.

```

1 void f()
2 {
3     x();
4     lock(&L); // a . b . c
5     a(); b(); c();
6     unlock(&L);
7     y();
8 }
9 void g()
10 {
11     x();
12     b(); c(); // ATOMICITY_VIOLATION: (b, c)
13     y();
14 }

```

Listing 3.4: An example of an *atomicity violation*

Summary of Phase 2 The sets of pairs that violate atomicity — the Δ sets — from the summaries of all analysed functions are finally reported to the user.

3.3 Atomer's Limitations

The basic version of Atomer has been proposed as it is detailed in Section 3.2. The first version of the analyser has also been implemented in [18], and it works as expected. Moreover, it can be used in practice to analyse various kinds of programs, and it may find *real-world atomicity related bugs*. Nevertheless, there are still several *limitations* and cases where the original version of Atomer would not work correctly, i.e., cases not addressed during the original proposal. Some of these cases are briefly discussed already in [18] and further described in [19].

So far, Atomer does not distinguish *different lock instances* used simultaneously in a program. Only calls of locks/unlocks are identified, and the parameters of these calls — *lock objects* — are not considered at all. Therefore, if there are several lock objects used, the analysis does not work correctly. Although this may happen in *real-life programs*, inasmuch as one could have, e.g., another (smaller) atomic section inside a current atomic section (this does not have to be evident at first because the *inner atomic section* could be, e.g., inside a nested function). For example,

```
... lock(L1); ... lock(L2); ... unlock(L2); ... unlock(L1); ...
```


Another possibility is an *alternating sequence of locks*, e.g., two locks are locked at first, and then, they are unlocked in the same order, i.e.,

```
... lock(L1); ... lock(L2); ... unlock(L1); ... unlock(L2); ...
```

Another limitation of Atomer’s basic version is that it supports only the analysis of programs written in the *C language* that use *PThread* locks to *synchronise concurrent threads*. Of course, in practice, many other *types of locks* for synchronisation of concurrent threads or even synchronisation of *concurrent processes* are used. Although the first version of Atomer can analyse C programs with other types of locks, these locks are not recognised as locks. Thus, the analysis would not work as expected. It would definitely be helpful also to analyse other languages than just C. As described in Section 2.3, Facebook Infer is capable of analysing programs written in C, C++, Objective-C, Java, and C#. The analysis algorithm could then be the same for all these languages because Facebook Infer’s *intermediate language* is analysed instead of directly analysing the input languages. Again, the first version of Atomer should be able to analyse the above languages, but it has not been tested within [18]. However, most importantly, other languages may use *very different lock types*, which would not be recognised. Examples of some advanced locking mechanisms (not supported by the basic version of Atomer) are *lock guards*, *re-entrant locks*, or *try-locks*.

Regarding *scalability*, the basic version of Atomer can have problems with more *extensive* and *complex* programs, as mentioned in [18] (problems with *memory* as well as problems with the *analysis time*). The problem is working with the sets of (A, B) pairs of *sequences* in the abstract states of Phase 1 and working with *sequences* of calls in the summaries of this phase. It may be necessary to store many of these sequences, and they can be very long (due to all different paths through the CFG of an analysed program). It may lead to the *state space explosion problem* [39].

One of the main reasons that Atomer’s first version reports *false alarms* is that in *critical sections*, in real-world programs, there are sometimes called *generic functions* that do not influence atomicity violations (such as functions for printing to the standard output, functions for recasting variables to different types, functions related to iterators, and whatever other “safe” functions for particular program types). Often, to find some atomicity violations, it is sufficient to focus only on certain “critical” functions. In practice, another issue is that in an analysed program, there can be “large” critical sections or critical sections in which appear function calls with a *deep hierarchy of nested function calls*. All the above cases may cause massive and “imprecise” atomic sequences that are the source of false alarms. However, regardless of the above issues, Atomer can still report quite some false alarms. It is due to the assumption that *sequences called atomically once* should *always be called atomically*, but this does not always have to hold. None of the above reasons that can generate false alarms is resolved in the first version of Atomer.

The next source of false alarms is something that the author of this work calls *local atomicity violations*. Imagine a function **f** that contains non-atomic calls of functions **a**, **b**, and these functions should always be invoked atomically. Obviously, this is an atomicity violation. However, suppose that **f** is called exclusively from atomic sections of other functions higher in the call hierarchy. In this case, in effect, that is not a real atomicity violation (it can be considered as a local atomicity violation within a single function, but globally, it is not). As a consequence, a false alarm would be reported by the basic version of Atomer. In real-life programs, this situation may be fairly common due to *complicated call graphs*.

Atomer considers only the *basic contracts for concurrency*, defined in Section 2.4.1. It is pretty limited in some circumstances, and therefore, Atomer can report *false alarms*. The basic contracts do not take into consideration the *data flow* within function calls. However, a better idea is to work with the assumption that a sequence of function calls must be atomic only if it *handles the same data*. Assume that functions f , g are manipulating with the same container C as follows: $f(C); g(C);$. These are called atomically. Somewhere else—where f , g are not called atomically—it does not necessarily cause an atomicity violation because they can be invoked with different arguments, which could be valid. This behaviour corresponds to the *extended contracts with parameters* (see Section 2.4.2). Another (more complex) limitation is that basic contracts do not consider any *contextual information*. It would be more precise to consider as atomicity violations such sequences that could be violated only by particular (“dangerous”) function calls, not by any calls. For example, suppose that there is the following sequence of functions called atomically: $f(); g();$. While somewhere else, these functions are not called atomically, it does not necessarily cause that it is an atomicity violation because, in this particular context, none of the “dangerous” functions can be executed by any concurrent thread. The *extended contracts with spoilers* formally describe these scenarios in Section 2.4.3.

A remarkable problem (though it is not directly a problem of Atomer) is identifying whether a reported atomicity violation is a *real bug* or whether it is just a false alarm. It could be really challenging, especially in *extensive real-life* programs.

Solutions for some of the above issues and limitations are proposed in Chapter 4 and further implemented in a new version of Atomer, detailed in Chapter 5.

Chapter 4

Proposal of Enhancements for Atomer

In this chapter, the author proposes solutions for some of the Atomer’s limitations stated in Section 3.3. The solutions enhance the analysis’s *precision* and *scalability*. In order to formally define these enhancements, the notions and symbols introduced in Section 3.2 are used. Some of the enhancements were briefly discussed already in [19].

Section 4.1 proposes an optimisation of Atomer’s scalability. The following Sections 4.2, 4.3, 4.4 cover precision improvements, i.e., extensions of Atomer with additional features that improve its ability to cope with cases that were not supported in the first version of Atomer, and that can be seen in *real-life code*. Furthermore, Chapter 5 provides an overview of the implementation of all the below improvements in a new version of Atomer.

In the following sections, to give an intuition, there are used listings with pieces of C programs (assuming `lock/unlock` functions for *mutual exclusion* to *critical sections*). In addition, there are used C++ and Java programs to illustrate the *locking mechanisms* in these languages.

4.1 Approximation of Sequences by Sets

Regarding *scalability*, the basic version of Atomer can have problems with more *extensive* and *complex* programs, which can manifest both in its *time* and *memory* consumption. The problems arise primarily due to working with the sets of (A, B) pairs of **sequences** of function calls in *abstract states* (during Phase 1). It may be necessary to store many of these sequences, and they could be very long (due to all different paths through the CFG of an analysed program). The author’s idea is to *approximate* these sets by working with sets of (A, B) pairs of **sets** of function calls. Apart from representing the abstract states of the first phase of the analysis, elements of these pairs do also appear in the first phase’s *summaries*, and they are then used during Phase 2 as well. Thus, it is needed to make a certain approximation in the summaries and their subsequent usage too. The approximated phases of the analysis and their collaboration are illustrated in Figure 4.1 (one can compare that with the illustration of the first version of Atomer in Figure 3.1).

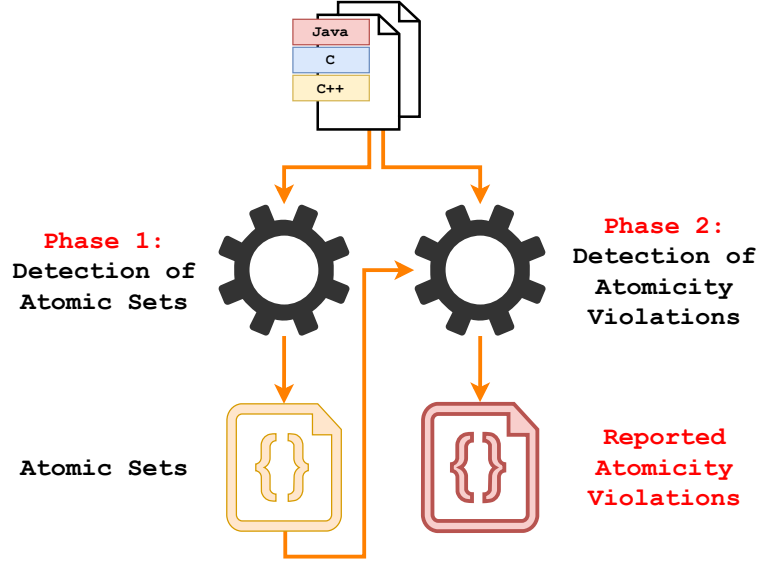


Figure 4.1: An illustration of the *phases* of the Atomer’s analysis and the *high-level analysis process* with an *approximation* of working with *sequences* by working with *sets* (moreover, note that a new version of Atomer accepts programs also written in C++ and Java languages, which is described in Section 5.3)

In particular, the proposed solution is more scalable because the ordering of function calls that appear in the pairs is not relevant anymore. Therefore, less memory is required because different sequences of function calls can map the same set. The analysis is also faster since there are stored fewer sets of function calls to work with. On the other hand, the analysis is less accurate because the new approach causes some loss of information. In practice, this loss of information could eventually lead to *false alarms*. However, the number of such false alarms is typically not that high as this thesis’s experimental evidence shows. Moreover, later, there are discussed some techniques that may rid of these false alarms.

4.1.1 Approximation with Sets in Phase 1

The *detection of sequences of calls to be executed atomically* now generates all (A, B) pairs of **sets** of function calls for each path instead of pairs of **sequences**, i.e., $(A, B) \in 2^\Sigma \times 2^\Sigma$. Here, A, B are not *reduced sequences* (the notion of a reduced sequence is not needed anymore) but sets. The purpose of the pairs is preserved. Hence, the abstract states are elements of the set $2^{2^\Sigma \times 2^\Sigma}$. In all the implemented algorithms and definitions, it is sufficient to work with:

- *sets* 2^Σ of functions, instead of *sequences* Σ^* of functions;
- the *empty sets* \emptyset , instead of the *empty sequences* ε ; and
- *unions* \cup of sets, instead of the *concatenation* \cdot of sequences.

The above implies that the *initial abstract state* of a function is changed to $s_{init} = \{(\emptyset, \emptyset)\}$. During the analysis of a function g with an abstract state s_g , when a leaf function f is called,

the abstract state's transformation is changed as follows:

$$s_g = \{p' \in 2^{2^\Sigma \times 2^\Sigma} \mid \exists p \in s_g : p' = \{(A', B') \in 2^\Sigma \times 2^\Sigma \mid \exists (A, B) \in p : \\ \textcolor{brown}{[\neg actual(p, (A, B)) \wedge (A', B') = (A, B)]} \vee \textcolor{red}{[actual(p, (A, B))]} \\ \wedge \textcolor{green}{[(lock \wedge (A', B') = (A, B \cup \{f\})) \vee (\neg lock \wedge (A', B') = (A \cup \{f\}, B))]\textcolor{blue}{]}}\}$$

Further, when an unlock is called, a new (A, B) pair is created as follows:

$$s_g = \{p' \in 2^{2^\Sigma \times 2^\Sigma} \mid \exists p \in s_g : p' = \{(A, B) \in 2^\Sigma \times 2^\Sigma \mid \\ \textcolor{brown}{[(A, B) = (\emptyset, \emptyset) \wedge setActual(p, (A, B))]} \vee (A, B) \in p\}\}$$

Other algorithms (e.g., calling an already analysed *nested* function) are modified analogically.

Another approximation was made in the summaries. The first component of the summary has to be changed to a set of sets of function calls because it is constructed from the B items from the abstract states, which are now sets. The second component of the summary can be changed to a set of function calls because even before, it was a reduced sequence of all the (A, B) pairs. Therefore, the ordering of function calls was significantly approximated even so. Moreover, it is used to analyse functions higher in the *call hierarchy* where it is appended to A or B , which are now sets. Thus, it would make no sense to store it in summaries as a sequence. Formally, the form of summaries χ changes from $2^{\Sigma^*} \times \Sigma^*$ to $2^{2^\Sigma} \times 2^\Sigma$. In particular, a summary $\chi_f \in 2^{2^\Sigma} \times 2^\Sigma$ of a function f is redefined as $\chi_f = (\mathcal{B}, C)$, where:

- $\mathcal{B} = \{B' \in 2^\Sigma \mid \exists p \in s_f : \exists (A, B) \in p : B \neq \emptyset \wedge B' = B\}$, where s_f is the abstract state at the end of the abstract interpretation of f .
- $C = \bigcup_{c \in C'} c$, where $C' = \{c \in 2^\Sigma \mid \exists p \in s_f : \exists (A, B) \in p : c = A \cup B\}$.

Example 4.1.1. For demonstrating the approximation of the analysis to sets, assume functions f and g from Listing 4.1. Further, assume that a, b, x, y are leaf nodes of the call graph. Before the approximation, when the analysis was working with sequences of function calls, Phase 1 of the analysis produced the following abstract states and summaries while analysing the functions:

- $s_f = \{((x \cdot y, a \cdot b), (y \cdot x, b \cdot a), (\varepsilon, \varepsilon))\}$, $\chi_f = (\{a \cdot b, b \cdot a\}, x \cdot y \cdot a \cdot b)$;
- $s_g = \{((y \cdot x, b \cdot a), (\varepsilon, \varepsilon))\}$, $\chi_g = (\{b \cdot a\}, y \cdot x \cdot b \cdot a)$.

Whereas, after the approximation, the produced abstract states and summaries are as follows: $s_f = s_g = \{((\{x, y\}, \{a, b\}), (\emptyset, \emptyset))\}$, $\chi_f = \chi_g = (\{\{a, b\}\}, \{a, b, x, y\})$. They are the same for both functions because there are the same locked/unlocked function calls; only the order of calls differs.

4.1.2 Approximation with Sets in Phase 2

The *detection of atomicity violations* in Phase 2 then works almost the same way as before the approximation. However, there is one difference. Before, the analysis implemented in

```

1 void f()
2 {
3     x(); y();
4     lock(&L); // a . b -> {a, b}
5     a(); b();
6     unlock(&L);
7     y(); x();
8     lock(&L); // b . a -> {a, b}
9     b(); a();
10    unlock(&L);
11 }
12 void g()
13 {
14     y(); x();
15     lock(&L); // b . a -> {a, b}
16     b(); a();
17     unlock(&L);
18 }

```

Listing 4.1: A code snippet used to illustrate the proposed *approximation* of the first phase of the analysis in a new version of Atomer using *sets of function calls*

the second phase looked for violations of atomic sequences obtained from Phase 1. Now, **atomic sets** are obtained from Phase 1; hence, the detection of atomicity violations needs to work with sets too. Again, the analysis looks for *pairs of functions that should be called atomically*, while this is not the case on some path through the CFG. This algorithm is identical to the algorithm before the approximation.

Nevertheless, it is needed to propose a new algorithm that derives the pairs of function calls (from the atomic sets) to be checked for atomicity (i.e., the set $\Omega \in 2^{\Sigma \times \Sigma}$). Intuitively, the second phase of the analysis now looks for non-atomic execution of any pair of functions \mathbf{f} , \mathbf{g} such that $\{\mathbf{f}, \mathbf{g}\}$ is a *subset* of some set of functions that were found to be executed atomically. In order to obtain the pairs, all possible pairs of functions are taken from atomic sets from Phase 1, i.e., all *2-element variations*. Formally, let S be an analysed program, and let $X_S \in 2^{2^{\Sigma} \times 2^{\Sigma}}$ be a set of all summaries of the program S . Then, all the atomic pairs (the first item of a pair may be empty if an atomic set consists of a single function) are obtained as follows:

$$\Omega = \{(x, y) \in \Sigma \times \Sigma \mid \exists (\mathcal{B}, C) \in X_S : \exists B \in \mathcal{B} : [|B| = 1 \wedge (x, y) \in \{\varepsilon\} \times B] \vee [|B| > 1 \wedge (x, y) \in B \times B \wedge x \neq y]\}$$

Example 4.1.2. For example, assume that Phase 1 analysed a function \mathbf{f} , which produced the summary $\chi_{\mathbf{f}} = (\mathcal{B}, C)$. Assume that before the approximation, a set of sequences of functions that should be called atomically was as follows: $\mathcal{B} = \{\mathbf{a} \cdot \mathbf{b} \cdot \mathbf{c}\}$. Then, the analysis looked for the following pairs of functions that are not called atomically: $\Omega = \{(\mathbf{a}, \mathbf{b}), (\mathbf{b}, \mathbf{c})\}$. Since the result of the first component of the summary was changed to the set \mathcal{B} of sets of functions that should be called atomically as follows: $\mathcal{B} = \{\{\mathbf{a}, \mathbf{b}, \mathbf{c}\}\}$, the analysis now looks for the following pairs of functions that are not called atomically (all 2-element variations): $\Omega = \{(\mathbf{a}, \mathbf{b}), (\mathbf{a}, \mathbf{c}), (\mathbf{b}, \mathbf{a}), (\mathbf{b}, \mathbf{c}), (\mathbf{c}, \mathbf{a}), (\mathbf{c}, \mathbf{b})\}$.

4.2 Advanced Manipulation with Locks

The original version of Atomer does not distinguish *different lock instances* in a program. Only calls of locks/unlocks are identified, and the parameters of these calls — *lock objects* — are not considered. Thus, if there are several lock objects used, the analysis does not work correctly.

In order to consider lock objects, it was proposed to distinguish between them using Facebook Infer’s built-in mechanism called *access paths* [25], explained below. The analyser does not perform a classical *alias analysis*, i.e., it does not perform a precise analysis for saying when arbitrary pairs of accesses to lock objects may alias (such an analysis is considered too expensive).

Access Paths The *syntactic access paths* [25] represent *heap locations* via the paths used to access them, i.e., they have the form of an expression consisting of a base variable followed by a sequence of fields. More formally, let Var be a set of all variables that can occur in a given program. Let $Field$ be a set of all possible field names that can be used in the program (e.g., structure fields). An access path π from the set Π of all access paths is then defined as follows:

$$\pi \in \Pi ::= Var \times Field^*$$

Access paths are already implemented in Facebook Infer. For instance, the principle of using access paths is used in an existing analyser in Facebook Infer — RacerD [5, 6, 17] — for data race detection. In general, no sufficiently precise *alias analysis* works *compositionally* and *at scale*. That is the motivation for using access paths in Facebook Infer.

Given a pair of accesses to lock objects, to determine whether these locks are equal, it is needed to answer the following question: “Can the accesses touch the same address?”. Remarkably, according to the authors of [5], access paths alone *almost* convey enough semantic information to answer the above question on their own. If two access paths are syntactically equal, it is almost (but not quite) true that they must refer to the same address. Syntactically identical paths can refer to different addresses if (i) they refer to different instances of the same object, or (ii) a prefix of the path is reassigned along one execution trace but not the other. These conditions cannot hold if an access path is *stable*, i.e., if none of its proper prefixes appears in assignments during a given execution trace, then it touches the same memory as all other stable accesses to the syntactic path. Therefore, the access paths’ syntactic equality is a reasonably efficient way to say (in an *under-approximate fashion*) that heap access touches the same address. Also, by using access paths, RacerD detected many errors in real-world programs, proving that the use of access paths can reveal real errors. This is why it was decided to use this principle to represent locks in Atomer.

During the analysis performed by Atomer (in both phases), each atomic section is identified by an access path of the lock that guards the section; see Sections 4.2.1, 4.2.2. Because *syntactically identical access paths* are used as the means for distinguishing atomic sections, some atomicity violations could be missed (or some false alarms could be reported) due to distinct access paths that refer to the same memory. However, the analysis’s precision is still significantly improved this way while preserving its *scalability*, and the stress is anyway put on finding likely violations, not on being *sound*.

Another limitation of Atomer in its basic version is that it does not count with *re-entrant locks* when a process can lock the same object multiple times without blocking itself, and then it should unlock the lock object the same number of times. This approach is, in fact, widespread, e.g., in Java, where so-called *synchronised blocks* are used, as demonstrates Listing 4.2. These blocks are re-entrant by default. To consider re-entrant locks in the analysis, the number of locks of individual lock objects is tracked in the abstract states of both phases of the analysis. A lock is unlocked as soon as this number decreases to 0. Also, an input parameter $t \in \mathbb{N}$ was proposed to limit the *upper bound* to which the analysis tracks precisely the number of times a given lock is locked. When this bound is reached, the *widening operator* is used to abstract the number to any value bigger than the bound. This is to ensure *termination* of the analysis. The idea of this upper bound limit comes from the approach used in RacerD.

```

1 public synchronized void f() {
2     x();
3     synchronized (this) {
4         a();
5         synchronized (this) { b(); c(); }
6         d();
7     }
8     y();
9 }

```

Listing 4.2: An example of *re-entrant locks* in Java using the `synchronized` keyword, which is implemented as a *monitor*. In the example, there are three locks used simultaneously over the same object. The entire method `f` is synchronised, which implicitly uses `this` as a *lock object*. Furthermore, the two *synchronised blocks* explicitly use the lock object `this`

4.2.1 Advanced Manipulation with Locks in Phase 1

Recall that the *detection of sets of calls to be executed atomically* is based on generating the pairs $(A, B) \in 2^\Sigma \times 2^\Sigma$. Now, these pairs are to be extended to store the access paths and the number of locks of lock objects that guard calls executed atomically, i.e., the B sets. Therefore, the (A, B) pairs are extended to tuples $(A, B, \pi, l) \in 2^\Sigma \times 2^\Sigma \times \Pi \times \mathbb{N}^\top$ where π is an access path that identifies a lock object that locks the atomic section that contains the calls from B , and l is the number of locks of the lock identified by π . For the clarity of the below description, let π be just a base variable, i.e., $\pi \in \Pi ::= \text{Var} \cup \{\varepsilon\}$. Note that π could also be ε , which is a special case when there is no lock associated with the (A, B) pair so far, i.e., B is empty and a lock was not called yet. \mathbb{N}^\top denotes $\mathbb{N} \cup \{\top\}$, where \top represents a number larger than t . Thus, the *abstract states* are elements of the set $2^{2^\Sigma \times 2^\Sigma \times \Pi \times \mathbb{N}^\top}$.

The analysis works as follows. When a function is called, it is appended to the A set of the element without an associated lock, i.e., the element where $\pi = \varepsilon$. Besides, the function is appended to the B sets of all the elements that have associated some lock that is currently locked, i.e., $\pi \neq \varepsilon$ and $l > 0$. When a lock is called, its identifier π_i is associated with the element without any lock associated with it, and the counter l of this element is set to 1. Then, l is incremented in all the elements where $\pi = \pi_i$ and $l > 0$. Moreover, it is created

a new empty element without a lock. Finally, when an unlock with the identifier π_i is called, l is decremented in all the elements where $\pi = \pi_i$ and $l > 0$. To formalise this process, let $inc_{\top} : \mathbb{N}^{\top} \rightarrow \mathbb{N}^{\top}$ and $dec_{\top} : \mathbb{N}^{\top} \rightarrow \mathbb{N}^{\top}$ be functions for incrementing and decrementing the number of locks of some lock objects w.r.t. the upper bound \top , respectively. These functions are defined as follows:

$$inc_{\top}(l) = \begin{cases} l + 1 & \text{for } l \neq \top \text{ and } l + 1 < t \\ \top & \text{otherwise} \end{cases} \quad dec_{\top}(l) = \begin{cases} l - 1 & \text{for } l \neq \top \text{ and } l - 1 \geq 0 \\ 0 & \text{for } l = 0 \\ t - 1 & \text{otherwise, i.e., } l = \top \end{cases}$$

The *initial abstract state* of a function is changed to $s_{init} = \{(\emptyset, \emptyset, \varepsilon, 0)\}$. During the analysis of a function \mathbf{g} with an abstract state $s_{\mathbf{g}}$, when a leaf function \mathbf{f} is called, the abstract state's transformation is changed as follows:

$$s_{\mathbf{g}} = \{p' \in 2^{2^{\Sigma} \times 2^{\Sigma} \times \Pi \times \mathbb{N}^{\top}} \mid \exists p \in s_{\mathbf{g}} : p' = \{(A', B', \pi', l') \in 2^{\Sigma} \times 2^{\Sigma} \times \Pi \times \mathbb{N}^{\top} \mid \\ \exists (A, B, \pi, l) \in p : [\pi = \varepsilon \wedge l = 0 \wedge (A', B', \pi', l') = (A \cup \{\mathbf{f}\}, B, \pi, l)] \vee [\pi \neq \varepsilon \\ \wedge ([l > 0 \wedge (A', B', \pi', l') = (A, B \cup \{\mathbf{f}\}, \pi, l)] \vee [l = 0 \wedge (A', B', \pi', l') = (A, B, \pi, l)])]\}\}$$

Further, when a lock identified by an access path π_i is called, the abstract state changes as follows:

$$s_{\mathbf{g}} = \{p' \in 2^{2^{\Sigma} \times 2^{\Sigma} \times \Pi \times \mathbb{N}^{\top}} \mid \exists p \in s_{\mathbf{g}} : p' = \{(A, B, \pi, l) \in 2^{\Sigma} \times 2^{\Sigma} \times \Pi \times \mathbb{N}^{\top} \mid \\ (A, B, \pi, l) = (\emptyset, \emptyset, \varepsilon, 0) \vee [(A, B, \varepsilon, 0) \in p \wedge \pi = \pi_i \wedge l = inc_{\top}(0)] \vee [(A, B, \pi, l') \in p \\ \wedge \pi \neq \varepsilon \wedge ((l' = 0 \vee \pi \neq \pi_i) \wedge l = l') \vee (l' > 0 \wedge \pi = \pi_i \wedge l = inc_{\top}(l')))]\}\}$$

Furthermore, when an unlock identified by the access path π_i is called, the abstract state changes as follows:

$$s_{\mathbf{g}} = \{p' \in 2^{2^{\Sigma} \times 2^{\Sigma} \times \Pi \times \mathbb{N}^{\top}} \mid \exists p \in s_{\mathbf{g}} : p' = \{(A, B, \pi, l) \in 2^{\Sigma} \times 2^{\Sigma} \times \Pi \times \mathbb{N}^{\top} \mid \\ (A, B, \pi, l) = (\emptyset, \emptyset, \varepsilon, 0) \vee [(A, B, \pi, l') \in p \wedge \pi \neq \varepsilon \\ \wedge ((l' = 0 \vee \pi \neq \pi_i) \wedge l = l') \vee (l' > 0 \wedge \pi = \pi_i \wedge l = dec_{\top}(l')))]\}\}$$

Other algorithms (e.g., calling an already analysed *nested* function) are changed analogically.

A *summary* $\chi_{\mathbf{f}} \in 2^{2^{\Sigma}} \times 2^{\Sigma}$ of a function \mathbf{f} is the same as earlier. Only access paths and lock counters from abstract states are ignored. This is, $\chi_{\mathbf{f}} = (\mathcal{B}, C)$, where:

- $\mathcal{B} = \{B' \in 2^{\Sigma} \mid \exists p \in s_{\mathbf{f}} : \exists (A, B, \pi, l) \in p : B \neq \emptyset \wedge B' = B\}$, where $s_{\mathbf{f}}$ is the abstract state at the end of the abstract interpretation of \mathbf{f} .
- $C = \bigcup_{c \in C'} c$, where $C' = \{c \in 2^{\Sigma} \mid \exists p \in s_{\mathbf{f}} : \exists (A, B, \pi, l) \in p : c = A \cup B\}$.

Example 4.2.1. Consider the function \mathbf{f} from Listing 4.3. There are two lock objects **L1** and **L2**, which are used simultaneously. Moreover, **L2** is locked several times without unlocking in between. Further, assume that **a**, **b**, **c** are leaf nodes of the call graph. After the extension described above, the produced summary is as follows: $\chi_{\mathbf{f}} = (\{\{\mathbf{b}\}, \{\mathbf{a}, \mathbf{b}, \mathbf{c}\}\}, \{\mathbf{a}, \mathbf{b}, \mathbf{c}\})$. Without the extension, the summary would be $\chi'_{\mathbf{f}} = (\{\{\mathbf{a}\}\}, \{\mathbf{a}, \mathbf{b}, \mathbf{c}\})$. The reason is that only the first locks/unlocks were detected. Other locks inside atomic sections and other unlocks outside atomic sections were ignored. Moreover, the abstract state after the execution of line 7 is as follows: $s_{\mathbf{f}_7} = \{(\emptyset, \{\mathbf{a}, \mathbf{b}\}, \mathbf{L1}, 1), (\{\mathbf{a}\}, \{\mathbf{b}\}, \mathbf{L2}, 2), (\{\mathbf{b}\}, \emptyset, \varepsilon, 0)\}$.


```

1 void f()
2 {
3     lock(&L1); // {a, b, c}
4     a();
5     lock(&L2); lock(&L2); // {b}
6     lock(&L2); unlock(&L2);
7     b();
8     unlock(&L2); unlock(&L2);
9     c();
10    unlock(&L1);
11 }

```

Listing 4.3: A code snippet used to illustrate the *advanced manipulation with locks* during the first phase of the analysis

4.2.2 Advanced Manipulation with Locks in Phase 2

The pairs Ω of functions that should be called atomically are computed the same way as earlier during the *detection of atomicity violations* in Phase 2. However, dealing with access paths and re-entrant locks must, of course, be reflected in the second phase of the analysis as well. For that, while looking for *atomicity violations of pairs of function calls*, from now, the analysis stores (in addition to pairs of the most recent function calls (\mathbf{x}, \mathbf{y}) and the set δ of pairs that have so far been identified as violating atomicity) all the most recent pairs of function calls locked under individual locks. Hence, the *abstract state* element gets the form $(\mathbf{x}, \mathbf{y}, \delta, \lambda) \in \Sigma \times \Sigma \times 2^{\Sigma \times \Sigma \times \mathbb{N}} \times 2^{\Sigma \times \Sigma \times \Pi \times \mathbb{N}^T}$, where $(\mathbf{x}, \mathbf{y}), \delta$ are as before, and λ is the set of the most recent function calls with their lock access paths and the number of locks of lock objects of these locks. Thus, the abstract states are elements of the set $2^{\Sigma \times \Sigma \times 2^{\Sigma \times \Sigma \times \mathbb{N}} \times 2^{\Sigma \times \Sigma \times \Pi \times \mathbb{N}^T}}$.

The analysis works as follows. When a function \mathbf{f} is called on some path that led to an abstract state $(\mathbf{x}, \mathbf{y}, \delta, \lambda)$, a new pair $(\mathbf{x}', \mathbf{y}')$ of the most recent function calls is created from the previous pair (\mathbf{x}, \mathbf{y}) such that $(\mathbf{x}', \mathbf{y}') = (\mathbf{y}, \mathbf{f})$. This pair is also stored in the locked pairs λ if there are any locks currently locked. Further, it is checked whether the new pair (or just the last call) violates the atomicity, and at the same time, the pair is not locked by any of the stored locks (i.e., $((\mathbf{x}', \mathbf{y}') \in \Omega \wedge \nexists (\mathbf{x}_\pi, \mathbf{y}_\pi, \pi, l) \in \lambda : (\mathbf{x}_\pi, \mathbf{y}_\pi) = (\mathbf{x}', \mathbf{y}')) \vee ((\varepsilon, \mathbf{y}') \in \Omega \wedge \nexists (\mathbf{x}_\pi, \mathbf{y}_\pi, \pi, l) \in \lambda : (\mathbf{x}_\pi, \mathbf{y}_\pi) = (\varepsilon, \mathbf{y}'))$). When the condition holds, the pair is added to the set δ of pairs that violate atomicity. When a lock with an identifier π_i is called, it is created a new empty element of λ with this identifier, and the lock counter l of this element is set to 1. Furthermore, the lock counter l of an element from λ with the access path π_i is incremented/decremented when a lock/unlock with the identifier π_i is called, respectively.

More formally, the *initial abstract state* of a function is defined as $s_{init} = \{(\varepsilon, \varepsilon, \emptyset, \emptyset)\}$. To formalise the analysis of a function, let \mathbf{f} be a called leaf function on a line $c \in \mathbb{N}$. Further, let s_g be the abstract state of a function \mathbf{g} being analysed before the function \mathbf{f} is called.

After the call of \mathbf{f} , the abstract state will be changed as follows:

$$s_g = \{(\mathbf{x}', \mathbf{y}', \delta', \lambda') \in \Sigma \times \Sigma \times 2^{\Sigma \times \Sigma \times \mathbb{N}} \times 2^{\Sigma \times \Sigma \times \Pi \times \mathbb{N}^\top} \mid \exists (\mathbf{x}, \mathbf{y}, \delta, \lambda) \in s_g : (\mathbf{x}', \mathbf{y}') = (\mathbf{y}, \mathbf{f}) \\ \wedge \lambda' = \{(\mathbf{x}'_\pi, \mathbf{y}'_\pi, \pi', l') \in \Sigma \times \Sigma \times \Pi \times \mathbb{N}^\top \mid \exists (\mathbf{x}_\pi, \mathbf{y}_\pi, \pi, l) \in \lambda : (\mathbf{x}'_\pi, \mathbf{y}'_\pi, \pi', l') = \\ (\mathbf{y}_\pi, \mathbf{f}, \pi, l)\} \wedge \delta' = \{(\mathbf{x}'', \mathbf{y}'', c') \in \Sigma \times \Sigma \times \mathbb{N} \mid (\mathbf{x}'', \mathbf{y}'', c') \in \delta \vee [(\mathbf{x}'', \mathbf{y}'') = (\mathbf{x}', \mathbf{y}') \vee \\ (\mathbf{x}'', \mathbf{y}'') = (\varepsilon, \mathbf{y}')] \wedge (\mathbf{x}'', \mathbf{y}'') \in \Omega \wedge c' = c \wedge \nexists (\mathbf{x}_\pi, \mathbf{y}_\pi, \pi, l) \in \lambda' : (\mathbf{x}_\pi, \mathbf{y}_\pi) = (\mathbf{x}'', \mathbf{y}'')]\} \}$$

Further, when a lock identified by an access path π_i is called, the abstract state is changed as follows:

$$s_g = \{(\mathbf{x}', \mathbf{y}', \delta', \lambda') \in \Sigma \times \Sigma \times 2^{\Sigma \times \Sigma \times \mathbb{N}} \times 2^{\Sigma \times \Sigma \times \Pi \times \mathbb{N}^\top} \mid \exists (\mathbf{x}, \mathbf{y}, \delta, \lambda) \in s_g : (\mathbf{x}', \mathbf{y}', \delta') = \\ (\mathbf{x}, \mathbf{y}, \delta) \wedge \lambda' = \{(\mathbf{x}_\pi, \mathbf{y}_\pi, \pi, l) \in \Sigma \times \Sigma \times \Pi \times \mathbb{N}^\top \mid (\mathbf{x}_\pi, \mathbf{y}_\pi, \pi, l) = (\varepsilon, \varepsilon, \pi_i, inc_\top(0)) \\ \vee [(\mathbf{x}_\pi, \mathbf{y}_\pi, \pi, l') \in \lambda \wedge [(\pi \neq \pi_i \wedge l = l') \vee (\pi = \pi_i \wedge l = inc_\top(l'))]]]\} \}$$

Furthermore, when an unlock identified by the access path π_i is called, the abstract state is changed as follows:

$$s_g = \{(\mathbf{x}', \mathbf{y}', \delta', \lambda') \in \Sigma \times \Sigma \times 2^{\Sigma \times \Sigma \times \mathbb{N}} \times 2^{\Sigma \times \Sigma \times \Pi \times \mathbb{N}^\top} \mid \exists (\mathbf{x}, \mathbf{y}, \delta, \lambda) \in s_g : \\ (\mathbf{x}', \mathbf{y}', \delta') = (\mathbf{x}, \mathbf{y}, \delta) \wedge \lambda' = \{(\mathbf{x}_\pi, \mathbf{y}_\pi, \pi, l) \in \Sigma \times \Sigma \times \Pi \times \mathbb{N}^\top \mid l > 0 \\ \wedge [(\mathbf{x}_\pi, \mathbf{y}_\pi, \pi, l') \in \lambda \wedge [(\pi \neq \pi_i \wedge l = l') \vee (\pi = \pi_i \wedge l = dec_\top(l'))]]]\} \}$$

A *summary* of the second phase of the analysis — $\chi_{\mathbf{f}} \in 2^{\Sigma \times \Sigma \times \mathbb{N}}$ — of a function \mathbf{f} is the same as earlier. Only the locked pairs λ with their access paths from abstract states are ignored. This is, $\chi_{\mathbf{f}} = \Delta = \bigcup_{\delta' \in \Delta'} \delta'$, where $\Delta' = \{\delta' \in 2^{\Sigma \times \Sigma \times \mathbb{N}} \mid \exists p \in s_{\mathbf{f}} : \exists (\mathbf{x}, \mathbf{y}, \delta, \lambda) \in p : \delta \neq \emptyset \wedge \delta' = \delta\}$ and $s_{\mathbf{f}}$ is the abstract state at the end of the abstract interpretation of \mathbf{f} .

Example 4.2.2. Consider the function \mathbf{g} from Listing 4.4. There are two lock objects **L1** and **L2**, which are used simultaneously. Further, assume that \mathbf{a} , \mathbf{b} are leaf nodes of the call graph. Then, assume that the result of the first phase of the analysis is that the pair (\mathbf{a}, \mathbf{b}) should be called atomically, i.e., $\Omega = \{(\mathbf{a}, \mathbf{b})\}$. Before the extension distinguishing of multiple lock instances, the analysis would report an atomicity violation for these functions (line 6). This is because the locks are not distinguished, and the unlock of **L1** (line 5) would unlock everything. On the other hand, after the extension, there are not reported any atomicity violations because the pair is still locked using **L2** (i.e., in Phase 2, it is produced the following summary: $\chi_{\mathbf{g}} = \emptyset$). The abstract state of \mathbf{g} after the execution of line 6 looks as follows: $s_{g_6} = \{(\mathbf{a}, \mathbf{b}, \emptyset, \{(\mathbf{a}, \mathbf{b}, \mathbf{L2}, 1)\})\}$.

4.2.3 Lock Guards

Finally, support for so-called *lock guard objects* has been proposed. Lock guards are objects associated with lock objects. One lock guard can be associated with multiple lock objects, and one lock object can be associated with multiple lock guards. When a lock guard is created, all lock objects associated with it are locked. When a lock guard is destroyed (usually, when a *scope of variables* is left), all lock objects associated with it are automatically unlocked. Exceptionally, under certain circumstances, lock guards can be locked/unlocked manually. Lock guards are widely used, especially in C++, but they are used, e.g., in Java as well. To cope with them, the analysis has been extended such that they are also identified

```

1 void g()
2 {
3     lock(&L1); // {}
4     lock(&L2); // {a, b}
5     unlock(&L1);
6     a(); b();
7     unlock(&L2);
8 }

```

Listing 4.4: A code snippet used to illustrate the *advanced manipulation with locks* during the second phase of the analysis

by *access paths*. The association between a lock guard and lock objects is modelled as a pair $(\pi_g, L) \in \Pi \times 2^\Pi$, where π_g is an access path that identifies the lock guard and L is a set of access paths that identify the lock objects associated with the guard identified by π_g . In the *abstract states* of both phases of the analysis, the associations between lock guards and lock objects are maintained as a set that is an element of the set $2^{\Pi \times 2^\Pi}$, i.e., there is a set of associations between lock guards and lock objects. Subsequent locks/unlocks of lock guards are then interpreted as a sequence of locks/unlocks of lock objects associated with these lock guards, respectively.

Example 4.2.3. Using Listing 4.5, it will be demonstrated how the analysis works with lock guards. There are used two lock objects **L1**, **L2**. **L2** is a *re-entrant lock* while **L1** is not. Besides, there are also two lock guards **G1** and **G2**. On line 9, at the same time, the lock guard **G1** is associated with both **L1** and **L2**, and the guard **G2** is associated only with **L2**. Conversely, the lock object **L1** is associated only with **G1**, and **L2** is associated with both **G1** and **G2**. These associations are stored in the analysis’s abstract states as the (π_g, L) pairs in the following set: $\{(\mathbf{G1}, \{\mathbf{L1}, \mathbf{L2}\}), (\mathbf{G2}, \{\mathbf{L2}\})\}$. For instance, when the scope of **G1** is left, the analysis abstractly interprets it as a sequence **L1.unlock()**; **L2.unlock()**;

4.3 Analysis’s Parametrisation

This section proposes solutions to reduce the number of *false alarms* reported by the basic version of Atomer. In particular, it is done by *parameterising* the analysis by some inputs provided by the user, i.e., in this case, the analysis does not have to be fully automatic. Namely, the parametrisation aims to filter out functions and *critical sections* that cause many reportings where most of them are very likely false alarms.

4.3.1 Specification of Critical Functions

One of the main reasons why Atomer in its first version reports false alarms is that, in practice, critical sections often interleave calls of functions that need to be executed atomically with common functions that need not be executed atomically (such as functions for printing to the standard output, functions for recasting variables to different types, functions related to iterators, and various other “safe” functions). Often, to find real atomicity violations, it is sufficient to focus on specific “critical” functions only.

```

1 #include <mutex>
2 std::mutex L1; std::recursive_mutex L2;
3 void f()
4 {
5     std::scoped_lock G1(L1, L2); // {a, b, c}
6     a();
7     {
8         std::lock_guard<std::recursive_mutex> G2(L2); // {b}
9         b();
10    }
11    c();
12 }

```

Listing 4.5: An example of *lock guards* in C++. The entire function `f` is locked with lock objects `L1`, `L2` using `std::scoped_lock`, which is a type of lock guard that can lock multiple lock objects at once (note that `std::scoped_lock` is available since C++ 17). However, there is also an inner atomic section locked with `L2` using `std::lock_guard`. `L2` is, in effect, locked twice—it is a *re-entrant lock* of a type `std::recursive_mutex`

For example, calls of *constructor* and *destructor* methods of classes do not lead to atomicity violations. Therefore, these calls can usually be ignored. Unfortunately, in general, it is not easy to differentiate between functions that should be set aside and functions to focus on because this distinction is application-specific. Therefore, the author decided to rely on the user to provide this information to the analysis. (Indeed, providing information of this kind is not so exceptional, e.g., for developers of libraries. A similar approach has also been chosen, e.g., in the *ANaConDA dynamic analyser* for concurrency issues [14], where the user can use so-called *hierarchical filters* to specify functions that the analysis should not monitor.) For this reason, the following input parameters of the analysis are proposed:

- a list of functions that will not be analysed,
- a list of functions that will be analysed (and all other functions will not be),
- a list of functions whose calls will not be considered, and
- a list of functions whose calls will be considered (and all other function calls will not be).

In other words, it is possible to specify *black-lists* and *white-lists* of functions to analyse and function calls to consider. It is also possible to combine these parameters, and they can be enabled for individual phases of the analysis. In the implementation of this approach mentioned in the following chapter, these parameters' values are read from input text files that contain one function name per line. Moreover, the implementation allows the user to specify sets of functions using *regular expressions* (in that case, the line must start with the letter `R` followed by whitespace).

4.3.2 Limitation of a Size of Critical Sections

Another issue often causing false alarms is that some programs contain “large” critical sections or critical sections that include function calls with a *deep hierarchy of nested function calls*. Both cases can cause massive and “imprecise” atomic sets that are the source of false alarms. Indeed, such “large” and/or “deep” critical sections are likely to contain a number of calls of functions that are not critical.

To resolve the “large” critical sections’ problem, the author proposes to parametrise the analysis by a parameter $d \in \mathbb{N}$ that limits the maximum length of a critical section to be taken into account. During its first phase, the analysis then discards all (A, B) pairs where $|B| > d$, i.e., it removes pairs where the number of functions in the set B (functions called atomically) is greater than the limit d .

To get to the above proposal of dealing with deeply nested critical functions, recall that, during the first phase of the analysis, when calling an already analysed *nested function*, the C set (i.e., the set of all called functions within a function) from its summary is used. If there is a deep hierarchy of nested function calls, the top level of the hierarchy uses function calls from all lower-level functions, leading to “large” critical sections. To avoid this problem, the summary $\chi = (\mathcal{B}, C) \in 2^{2^\Sigma} \times 2^\Sigma$ in Phase 1 is redefined as $2^{2^\Sigma} \times 2^{\mathbb{N} \times 2^\Sigma}$, i.e., C is no longer a set of all functions called within an analysed function. It is a set of pairs where each pair represents functions called at a particular level in the hierarchy of a nested function (0 means the top-level). For instance, the summary $\chi_f = (\emptyset, \{(0, \{\mathbf{a}, \mathbf{b}\}), (1, \{\mathbf{x}, \mathbf{y}\})\})$ of a function \mathbf{f} means that there were called functions \mathbf{a} , \mathbf{b} in \mathbf{f} and that there were called functions \mathbf{x} , \mathbf{y} in functions one level lower in the call tree (i.e., in functions directly invoked from \mathbf{f}). During the analysis, the summaries are passed among functions in the call hierarchy. Furthermore, the analysis uses a parameter $r \in \mathbb{N}$ to limit the number of levels considered during analysing nested functions.

More particularly, to derive these extended summaries, it is also needed to store in the abstract states the information about in which levels appear which function calls. When analysing a function \mathbf{f} and a leaf function \mathbf{g} is called, \mathbf{g} is added to the level 0 of functions called within \mathbf{f} . Furthermore, when calling an already analysed function \mathbf{h} , the functions called at a level i of \mathbf{h} are appended into the appropriate A or B sets of \mathbf{f} if $i < r$. Finally, the information about in which levels appear which function calls is passed from a lower-level function to a higher-level function. In particular, if some function is called at a level i of \mathbf{h} , it is added to the level $i + 1$ of the function \mathbf{f} (which is one level higher) if $i + 1 < r$.

Example 4.3.1. Assume functions \mathbf{x} , \mathbf{y} , \mathbf{z} from Listing 4.6. Their summaries for the first phase of the analysis are as follows: $\chi_z = (\emptyset, \{(0, \{\mathbf{z1}, \mathbf{z2}\})\})$, $\chi_y = (\emptyset, \{(0, \{\mathbf{y1}, \mathbf{y2}, \mathbf{z}\}), (1, \{\mathbf{z1}, \mathbf{z2}\})\})$, $\chi_x = (\emptyset, \{(0, \{\mathbf{x1}, \mathbf{x2}, \mathbf{y}\}), (1, \{\mathbf{y1}, \mathbf{y2}, \mathbf{z}\}), (2, \{\mathbf{z1}, \mathbf{z2}\})\})$. When the value of the parameter r is set to 1, the summary of the function \mathbf{x} is reduced as follows: $\chi'_x = (\emptyset, \{(0, \{\mathbf{x1}, \mathbf{x2}, \mathbf{y}\}), (1, \{\mathbf{y1}, \mathbf{y2}, \mathbf{z}\})\})$.

```

1 void z() { z1(); z2(); }
2 void y() { y1(); y2(); z(); }
3 void x() { x1(); x2(); y(); }

```

Listing 4.6: A code snippet used to illustrate the *limitation of considered nested functions*

4.4 Local/Global Atomicity Violations

The author of the thesis uses two notions—*local* and *global atomicity violations*—to distinguish two classes of “atomicity violations”. A local atomicity violation is when there is an atomicity violation within the *scope of a function* being analysed. However, in effect, it does not have to be an error in the context of the whole analysed program because the entire function may *always be locked* when it is used. Thus, it is ensured that all instructions of the function are always executed atomically in the program. On the other hand, when there is at least a single case of the function not being called atomically, it entails that there exists a possibility of a real error, which is called a global atomicity violation.

In the first version of Atomer, local/global atomicity violations are not distinguished. All atomicity violations are considered global ones. This is, of course, a source of *false alarms*. In extensive real-world programs, local atomicity violations are quite common because calls of some smaller functions are often included in atomic sections. Unfortunately, it can be challenging to differentiate between local and global atomicity violations in such programs due to *complex call graphs*.

Example 4.4.1. Consider Listing 4.7, where the functions `x`, `y` are leaf nodes of the call graph. `main` is the *top-level function* that contains the calls of functions `f` and `g`. From the first phase of the analysis, it is derived that the pair (x, y) should be called atomically because these functions are called in `g`, which is locked in `main`. In particular, the first phase’s summary of `main` is $\chi_{\text{main}} = (\mathcal{B}, C) = (\{\{g, x, y\}\}, \{(0, \{f, g\}), (1, \{x, y\})\})$, and $\exists B \in \mathcal{B} : \{x, y\} \subseteq B$. During Phase 2, there are reported atomicity violations on lines 1 and 2 due to non-atomic executions of the pair (x, y) within functions `f` and `g`, respectively. Consequently, the second phase’s summaries of these functions are as follows: $\chi_f = \{(x, y, 1)\}$, $\chi_g = \{(x, y, 2)\}$. Nevertheless, while the violation in `f` is global, the violation in `g` is, as a matter of fact, local because `g` is called exclusively from `main` where it is locked. Therefore, a false alarm is reported.

```
1 void f() { x(); y(); } // GLOBAL_ATOMICITY_VIOLATION: (x, y)
2 void g() { x(); y(); } // LOCAL_ATOMICITY_VIOLATION: (x, y)
3 void main()
4 {
5     f();
6     lock(&L); // {g, x, y}
7     g();
8     unlock(&L);
9 }
```

Listing 4.7: An example of so-called *local* and *global atomicity violations*

To identify local atomicity violations during Phase 2, the author proposes the following approach in a new Atomer’s version. First, two *severities* of the atomicity violation reporting are defined:

- **Warning** for local atomicity violations and
- **Error** for (real) global atomicity violations.

In particular, the sets δ of pairs that violate atomicity in the abstract states and the appropriate sets Δ from the summaries are extended from $2^{\Sigma \times \Sigma \times \mathbb{N}}$ to $2^{\Sigma \times \Sigma \times \mathbb{N} \times \mathcal{S}}$, where $\mathcal{S} = \text{Warning} \mid \text{Error}$.

When an atomicity violation within an analysed function is detected as earlier, it is labelled as global, i.e., the **Error** severity is assigned to it. In the course of the analysis, the information about the violations from function summaries (i.e., the pairs from Δ) are passed from lower-level functions to higher-level functions along to the call tree. Once a nested function call is performed inside an atomic section, all its violations are labelled as local, i.e., the **Warning** severity is assigned to them. Finally, atomicity violations (both local as well as global) are reported only from *top-level functions*, i.e., functions not called from any other functions; thus, they cannot be called under a lock somewhere higher in the call hierarchy.

Example 4.4.2. Consider Listing 4.7 again. The first analysis's phase derived that the pair (\mathbf{x}, \mathbf{y}) should be called atomically, as explained in Example 4.4.1. After extending Atomer for distinguishing local and global atomicity violations, the second phase of the analysis works as follows. First, functions **f** and **g** are analysed. This phase produces the following summaries for these functions: $\chi_{\mathbf{f}} = \{(\mathbf{x}, \mathbf{y}, 1, \text{Error})\}$, $\chi_{\mathbf{g}} = \{(\mathbf{x}, \mathbf{y}, 2, \text{Error})\}$. These violations must not be reported yet. Further, **main** is analysed. It receives summaries from **f** and **g**. All violations from **g** are, however, labelled as local atomicity violations. As a result, the summary of **main** is $\chi_{\text{main}} = \{(\mathbf{x}, \mathbf{y}, 1, \text{Error}), (\mathbf{x}, \mathbf{y}, 2, \text{Warning})\}$. In the end, the reporting is made from **main** because it is the only top-level function. On line 1, there is reported a global atomicity violation. On the contrary, it is reported just a warning (local violation) on line 2 because the call of this function is locked in **main**.

Chapter 5

Implementation of a New Version of Atomer

This chapter discusses the implementation of a new version of *Atomer*. It summarises the implementation of the whole analyser (i.e., it includes the implementation of the Atomer’s basic version from [18], also described in Chapter 3) with a particular focus on the new enhancements proposed in Chapter 4. Atomer is implemented as a *module of Facebook Infer* introduced in Section 2.3. The implemented algorithms are illustrated using convenient *pseudocode* and listings of simplified code written in *OCaml*, which is an implementation language of Atomer and Facebook Infer. Sections 5.1 and 5.2 describe the implementation of *both phases* of the analysis in more detail. Moreover, Section 5.3 outlines the implementation of support for *programming languages* and *locking mechanisms* newly added to Atomer.

The implementation of both the basic version and the new version of Atomer is *publicly available* at GitHub¹. An installation and user manual are available in Appendix B. Together with some examples, these manuals are also available in the attached memory media of the thesis (see Appendix A) and in Atomer’s Wiki². As it was already said, Atomer is implemented in OCaml, which is a *functional* programming language. However, it also allows using the *imperative* and *object-oriented* paradigms exploited in the Atomer’s implementation as well. Since the implementation of the first version of Atomer, Facebook Infer was many times updated (in its Git repository, there are several new commits every day). Due to maintaining Atomer up-to-date with Facebook Infer, and due to some refactoring of Atomer itself, the code of the first version of Atomer since its first implementation is considerably different.

Atomer itself is implemented in files in a directory `infer/src/atomicity`. In particular, the first phase of the analysis is implemented in `AtomicSets.ml[i]`, and its *abstract domain* Q is implemented in `AtomicSetsDomain.ml[i]`. The second phase is implemented in `AtomicityViolations.ml[i]`, and its abstract domain Q is implemented

¹The implementation of a new version of **Atomer** at GitHub as an *open-source* repository: <https://github.com/harmim/infer/tree/atomer-v2.0.0>. The implementation of the first version is in a branch `atomicity`. At the following address, there are available releases of both versions, including source code and executable binaries (v1.0.0 is the basic version, and v2.0.0 corresponds to the new version): <https://github.com/harmim/infer/releases>.

²**Atomer’s Wiki** provides an installation and user manual together with some examples. It is available at GitHub: <https://github.com/harmim/infer/wiki>.

in `AtomicityViolationsDomain.ml[i]`. The phases are implemented as separate analysers in Facebook Infer. The output of the first phase is the input to the second one (as earlier shown in Figure 4.1). These analysers are defined as modules of Facebook Infer in `infer/src/base/Checker.ml[i]`, and they are registered to the framework in `infer/src/backend/registerCheckers.ml`. The analysers are enabled only by specific command-line arguments; see the mentioned manuals.

Note that in the listings in the below sections, `Domain` refers to an abstract domain Q of a particular phase of the analysis (the `AtomicSetsDomain/AtomicityViolationsDomain` module for the first/second phase of the analysis, respectively). Then, `Domain.t` is a type of *abstract state* of the particular phase. Furthermore, `Domain.Summary.t` is a type of *summary* χ in the particular phase.

5.0.1 The Main Analysis Function

For both phases of the analysis, the analyser is implemented as an *abstract interpreter* using the `LowerHil` module, which transforms *Smallfoot Intermediate Language (SIL) instructions* (mentioned in Section 2.3) into *High-level Intermediate Language (HIL) instructions*. HIL instructions wrap SIL instructions and simplify their utilisation, i.e., it is an abstraction over SIL. For representing functions, a filtered view of the *forward CFG* that skips *exceptional control-flow paths* is used. This type of CFG corresponds to the `ProcCfg.Normal` module in Facebook Infer. The final analyser is created as a module `Analyser` using the `LowerHil.MakeAbstractInterpreter` function.

```

1 let analyse_procedure data : Domain.Summary.t option =
2   (* Should be the analysis of this function skipped? *)
3   if f_is_ignored (Pdesc.get_pname data.pdesc) then None
4   else
5     let pre : Domain.t =
6       Domain.initial (* an initial abstract state (pre-condition) *)
7     in
8       (* Compute the final abstract state (post-condition). *)
9       match Analyser.compute_post data ~initial:pre with
10      | Some (post : Domain.t) ->
11        (* Convert the abstract state to a summary. *)
12        let summary : Domain.Summary.t = Domain.Summary.create post in
13        Some summary
14      | None -> Logging.die InternalError "Analysis failed."

```

Listing 5.1: The *main analysis function* for analysing individual functions in a given program

The *main analysis function* `analyse_procedure` is implemented the same way for both phases of the analysis (for Phase 1, it is implemented in the `AtomicSets` module, and for Phase 2, it is implemented in the `AtomicityViolations` module). A simplified implementation is given in Listing 5.1. Facebook Infer’s backend invokes this function for each function in an analysed program. It computes a *summary* of the analysed function. At first, the function checks whether the analysis of an analysed function should be ignored (line 3). It

has to do with the *specification of critical functions* proposed in Section 4.3. The implementation of that is covered in Section 5.0.4. Further, `analyse_procedure` computes the final abstract state (*post-condition*) for the analysed function using the `Analysers.compute_post` function (line 9). As a *pre-condition*, the *initial abstract state* `Domain.initial` from an abstract domain is used. On line 12, the abstract state is appropriately converted to an interprocedural function summary using `Domain.Summary.create`. Finally, this summary is returned. The crucial parts are here the types `Domain.t`, `Domain.Summary.t`, and the functions `Domain.initial`, `Domain.Summary.create`. These are different for individual phases of the analysis. They are described in Sections 5.1, 5.2.

```

1 let exec_instr (s : Domain.t) data : HilInstr.t -> Domain.t = function
2   | Call (Direct callee, params) (* the CALL instruction *)
3     (* Should be this function call ignored? *)
4     when f_is_ignored callee ~actuals:(Some params) -> s
5   | Call (Direct callee, params, loc) -> ( (* the CALL instruction *)
6     (* Does the call relate to the locking? *)
7     match ConcurrencyModels.get_lock_effect callee params with
8     (* Process locks/unlocks/lock guards. *)
9     | Lock locks ->
10      Domain.apply_locks (get_paths locks) s
11     | Unlock locks ->
12      Domain.apply_unlocks (get_paths locks) s
13     | GuardConstruct {guard; locks} ->
14      Domain.apply_guard_construct (get_path guard) (get_paths locks) s
15     | GuardDestroy guard ->
16      Domain.apply_guard_destroy (get_path guard) s
17     | NoEffect -> ( (* the call of a classical function *)
18       let s : Domain.t =
19         (* Process the called function. *)
20         Domain.apply_call ~fName:(Pname.to_string callee) loc s
21       in
22       (* Read a summary of an already analysed nested function. *)
23       match data.analyse_dependency callee with
24       | Some (summary : Domain.Summary.t) ->
25         (* Apply the summary to the abstract state. *)
26         Domain.apply_summary summary s
27       | None -> s (* leaf node *) ) )
28   | _ -> s (* Do nothing for other instructions. *)

```

Listing 5.2: The implementation of the *abstract transformers* (also called *transfer functions*)

5.0.2 Abstract Transformers

The *abstract transformers* $\tau : Instr \times Q \rightarrow Q$ over an abstract domain Q (also called *transfer functions*) are implemented almost the same way for both phases of the analysis. In general, the abstract transformer takes an abstract state as its input and produces an abstract state as an output while interpreting an instruction of an analysed program. The imple-

mentation is illustrated in Listing 5.2. It is implemented using the `exec_instr` function (in the `AtomicSets/AtomicityViolations` module for Phase 1/Phase 2, respectively). It considers only function calls, i.e., instruction `CALL`. The abstract state remains unchanged when other instructions are executed (line 28). On line 4, there is checked whether the called function should be ignored (see Section 4.3). Otherwise, line 7 identifies calls related to *locking*. When the called function is some lock, `Domain.apply_locks` is called to update the abstract state appropriately (line 10). When the call is a construction of a *lock guard*, `Domain.apply_guard_construct` is used (line 14), etc. Note that besides the cases used in the pattern matching shown in the listing, there are also processed `GuardRelease`, `GuardLock`, and `GuardUnlock` for advanced manipulation with lock guards in the actual implementation. The functions `get_path/get_paths` are used for acquiring *access paths* of *lock objects*. The functions `Domain.apply_locks` and `Domain.apply_unlocks` are implemented differently for both phases of the analysis. This implementation is described in Sections 5.1, 5.2. However, the functions `Domain.apply_guard_construct` and `Domain.apply_guard_destroy` are the same for both phases. Thus, they are explained below. If the called function is not related to locking (line 17), the function call is interpreted in the abstract domain using `Domain.apply_call`. Further, line 23 tries to read a summary of the called function. If it is not a leaf node of the call tree, its summary is used to update the abstract state using `Domain.apply_summary` on line 26. The functions `Domain.apply_call`, `Domain.apply_summary` are essential. Thus, they are introduced in Sections 5.1, 5.2. Moreover, note that on line 5, there is a parameter `loc` of the called function. It is the location of the call in the source program. It is used to identify the location of an atomicity violation. Therefore, it is needed only in the implementation of the abstract transformers for the second phase of the analysis.

When a *lock guard is constructed*, an abstract state is updated using the function `Domain.apply_guard_construct` illustrated in Algorithm 5.1. It first appends the association of the lock guard with the underlying locks to the *guards* field of each element (program path) of the abstract state. Then, all the underlying locks are immediately locked using `Domain.apply_locks`, which is defined later. Both the lock guard and the underlying locks are identified through access paths.

Algorithm 5.1: Updating an abstract state after the *construction* of a *lock guard*

Data: lock guard's access path $\pi_g \in \Pi$; access paths $L \in 2^\Pi$ of locks associated with π_g ; abstract state $s \in \text{Domain.t}$

```

1 def apply_guard_construct( $\pi_g$ ,  $L$ ,  $s$ ):
2   for  $p \in s$  do  $p.guards \leftarrow p.guards \cup \{(\pi_g, L)\}$ ;    /* create a new guard */
3   return apply_locks( $L$ ,  $s$ );
4 end
```

After the *destruction of a lock guard*, an abstract state is changed using the function `Domain.apply_guard_destroy` given in Algorithm 5.2. For each element of the abstract state, it obtains locks associated with the lock guard being destroyed. This association is then removed from the abstract state's field *guards*. Finally, the associated locks are unlocked using the function `Domain.apply_unlocks` extended to apply unlock calls for a particular program path only.

Algorithm 5.2: Updating an abstract state after the *destruction* of a lock guard

Data: lock guard's access path $\pi_g \in \Pi$; abstract state $s \in \text{Domain.t}$

```
1 def apply_guard_destroy( $\pi_g$ ,  $s$ ):
2   for  $p \in s$  do
3      $L \leftarrow p.\text{guards}(\pi_g)$ ;           /* get associated locks */
4      $p.\text{guards} \leftarrow p.\text{guards} \setminus \{(\pi_g, L)\}$ ;      /* remove the guard */
5      $s \leftarrow \text{apply_unlocks}(L, s, p)$ ;
6   end
7   return  $s$ ;
8 end
```

5.0.3 Abstract Interpretation Operators

The abstract domains of both phases of the analysis are quite different. However, the *abstract interpretation operators* of these domains are the same because for both phases, the abstract states $s \in Q$ are *sets* of elements that represent an abstraction of program paths. The implementation of the operators is shown in Listing 5.3. Note that \mathbf{t} is a type of abstract state. In fact, it is an alias for `TSet.t`, a module representing a set of structures where the fields of these structures are defined differently in each phase of the analysis. Thus, each phase defines its own `TSet`. Particular abstract interpretation operators (implemented in Listing 5.3) are defined as follows:

- The *ordering operator* \sqsubseteq (`leq`) is defined as follows. Let `lhs` be the left-hand side of the operator and `rhs` the right-hand side. Then, $\text{lhs} \sqsubseteq \text{rhs}$ iff `lhs` is a *subset* of `rhs`. In other words, $s \sqsubseteq s' \iff s \subseteq s'$, where $s, s' \in Q$.
- The *join operator* \sqcup (`join`) is defined simply as the *union* of two abstract states, i.e., $s \sqcup s' \iff s \cup s'$.
- The *widening operator* ∇ (`widen`) is defined as joining the previous and the next abstract state since the domains are *finite*. In particular, $s \nabla s' \iff s \sqcup s'$.

```
1 (* lhs <= rhs if lhs is a subset of rhs. *)
2 let leq ~(lhs : t) ~(rhs : t) : bool = TSet.subset lhs rhs
3 (* Union of abstract states. *)
4 let join (s1 : t) (s2 : t) : t = TSet.union s1 s2
5 (* Join the previous and the next abstract state. *)
6 let widen ~(prev : t) ~(next : t) ~(num_iters : int) : t = join prev next
```

Listing 5.3: The implementation of the *abstract interpretation operators*

5.0.4 Specification of Critical Functions

In Section 4.3, there are proposed *input parameters* for specifying *black-lists/white-lists* of functions to analyse or function calls to consider. The decision of whether a function is analysed or a function call is considered is made in `analyse_procedure` (Listing 5.1) and

`exec_instr` (Listing 5.2), respectively. In particular, the function `f_is_ignored` implements the decision making whether a particular function is considered/ignored according to the input parameters. These parameters may be combined. For example, one can allow only specific functions to be analysed and ignore several selected function calls considered during the analysis. Moreover, these parameters may be enabled for individual phases of the analysis or specified differently for each phase. `f_is_ignored` is implemented in the `AtomicityUtils` module.

The parameters' values are read from *input text files*. How to enable the parameters (i.e., how to specify the file names of the input files using command-line arguments) is described in the user manual. The input file should contain one function name to consider/ignore per line. There may be empty lines, and line comments can be used with the character `#`. Moreover, it is possible to specify sets of function names using *regular expressions* (in that case, the line must start with the letter `R` followed by any number of whitespaces). The syntax of the regular expressions corresponds to the OCaml's regular expressions syntax³. An example of such an input file is in Listing 5.4. Processing of input files is implemented in the `functions_from_file` class in the `AtomicityUtils` module.

```
intValue # this is a comment
R_[A-Z]+.*
foo
Logger.log
R_\\(String\\)?\\.format
#this is another comment
R_\\(Int\\|Float\\)\\.toString
```

Listing 5.4: An example of an *input text file* with a *black-list/white-list* of functions to analyse or function calls to consider

5.0.5 Representation of Locks in Abstract States

In module `AtomicityUtils`, there is implemented a sub-module `Lock` for *working with locks* in both phases of the analysis. The essential parts of the implementation are shown in Listing 5.5. Data of this module are of a type `Lock.t`, a pair of an *access path* (`AccessPath.t`) of a lock object and the number of locks of this lock object, i.e., it corresponds to the (π, l) pairs of the abstract states in the proposal.

The bottom value of l is 0; this is `bot`. The top value of l is `Config.atomicity_lock_level_limit`; this is `top`. This value comes from the user (the default value is 5). This corresponds to the \top element (parameter t) from the proposal. How the user can set this value is described in the user manual. Locking/unlocking of the lock is implemented such that l is incremented/decremented, respectively. It is crucial that it is implemented so that l remains within bounds $[0, t]$ to ensure termination. The module also implements other functions, e.g., a function that finds out whether the lock object is currently locked, etc. However, other functions are not too interesting; thus, they are not shown in the listing.

³The syntax of OCaml's regular expressions: <https://ocaml.org/api/Str.html>.

```

1 module Lock = struct
2   (* A lock's access path with the number of times it has been acquired. *)
3   type t = AccessPath.t * int (* (p, l) *)
4   (* The bottom value of l. *)
5   let bot : int = 0
6   (* The top value of l. *)
7   let top : int = Config.atomicity_lock_level_limit
8   (* Increases the value of l. *)
9   let lock ((p, l) : t) : t = if l = top then (p, l) else (p, l + 1)
10  (* Decreases the value of l. *)
11  let unlock ((p, l) : t) : t = if l = bot then (p, l) else (p, l - 1)
12 end

```

Listing 5.5: The implementation of a module that represents *locks in abstract states*

5.1 Phase 1 — Detection of Atomic Sets

The first phase of the analysis is enabled by a command-line argument `--atomic-sets-only` (see the user manual for more details). The detected *atomic sets* by this phase (i.e., sets of functions *expected to execute atomically*) are printed into a text file, which is then used as an input to the second phase. This section, in particular, explains the concrete implementation of the abstract states and summaries. Moreover, it is outlined how they are used during the analysis. The output of this phase is in more detail described in Section 5.1.1.

Abstract State The implementation of the *abstract state* is illustrated in Listing 5.6. The abstract state is of the type `TSet.t` (`t` is an alias for it). `TSet` itself is defined on line 13. It is a module representing a set of structures `p` (lines 10–12). Each element corresponds to a single program path. Fields of this structure are composed of the types and modules defined on the lines above. Lines 15–21 define the *initial abstract state*. The fields are as follows:

- **calls**: a *set of function calls* (a function call is a string), i.e., an element of 2^Σ . It captures function calls before a lock. When a lock is called, a new (A, B) pair is created, and the **calls** set became the A set.
- **callsPairs**: a *set of the (A, B) pairs with associated locks*, where a lock \mathcal{L} is of the type `Lock.t` defined earlier. In fact, it is a set of the tuples (A, B, π, l) , i.e., an element of $2^{2^\Sigma \times 2^\Sigma \times \Pi \times \mathbb{N}^\top}$. These elements are considered currently locked, i.e., new function calls are appended to the B set.
- **finalCallsPairs**: a *set of the (A, B) pairs that are “complete”*, i.e., an element of $2^{2^\Sigma \times 2^\Sigma}$. These are the final pairs that were already unlocked (thus, the association with a lock is not needed anymore). In the end, the final atomic sets are derived from this field.
- **allCalls**: a *list of sets of function calls*, i.e., an element of $2^{\mathbb{N} \times 2^\Sigma}$. It is used to capture which functions were called at which levels of the call hierarchy. Indices of

the list represent individual levels (functions called in an analysed function are stored at level 0, functions directly invoked from the analysed function are at level 1, etc.). This field is used to limit the number of nested functions in critical sections using the proposed parameter r .

- **guards**: a *map* $\Pi \rightarrow 2^{\mathbb{N} \times \Pi}$. It maps a lock guard's access path to a list of access paths of associated locks. It can be viewed as a set of the pairs $(\pi_g, L) \in \Pi \times 2^\Pi$ from the proposal. (An ordering in the list is ignored, and it is treated as a set. The list is used because it simplifies the implementation in OCaml.)

```

1 module S = Set.Make(String) (* a set of strings *)
2 type ab = S.t * S.t (* (A,B) pair *)
3 module S_ab = Set.Make(struct type t = ab end) (* a set of (A,B) pairs *)
4 type ab_lock = ab * Lock.t (* (A,B) pair with a lock, i.e., (A,B,π,l) *)
5 module S_ab_lock = (* a set of the above *)
6   Set.Make(struct type t = ab_lock end)
7 module M_AP = Map.Make(AccessPath) (* a map where AccessPath is a key *)
8 type g = (* maps lock guards to a list of associated locks *)
9   AccessPath.t list M_AP.t (* (πg, L) where πg is a key and L is a value *)
10 type p = (* an element of an abstract state *)
11   { calls: S.t; callsPairs: S_ab_lock.t; finalCallsPairs: S_ab.t
12     ; allCalls: S.t list; guards: g }
13 module TSet = Set.Make(struct type t = p end) (* a set of the above *)
14 type t = TSet.t (* a type of the abstract state (an alias for TSet.t) *)
15 let initial : t = (* the initial abstract state of a function f *)
16   TSet.singleton (* an "empty" element (initial program path of f) *)
17     { calls= S.empty (* an empty set *)
18       ; callsPairs= S_ab_lock.empty (* an empty set *)
19       ; finalCallsPairs= S_ab.empty (* an empty set *)
20       ; allCalls= [S.empty] (* an empty set at level 0 of f *)
21       ; guards= M_AP.empty (* an empty map *) }

```

Listing 5.6: A definition of the *abstract state* in Phase 1

The primary functions of the abstract domain that work with abstract states are `apply_call`, `apply_locks`, and `apply_unlocks`. In earlier Listing 5.2, there is shown how these functions are utilised.

The function `apply_call` from Algorithm 5.3 demonstrates how an abstract state is updated when some classical function is invoked. For example, note that on line 4, the called function is appended to level 0 of the `allCalls` field. This is because the function is directly invoked from the analysed function, which is level 0. Moreover, e.g., notice that on line 8, the entire pair (A, B) is discarded when B is greater than the parameter d . This is to rid of “large” critical sections, as described in the proposal.

Furthermore, the implementation of the functions `apply_locks` and `apply_unlocks` is illustrated using the algorithms from Appendix C, namely, Algorithms C.2 and C.1, respectively.

Algorithm 5.3: Updating an abstract state after a *call of a function*

Data: called function $f \in \Sigma$; abstract state $s \in \text{Domain.t}$; maximum length $d \in \mathbb{N}$ of a critical section

```
1 def apply_call(f, s):
2   for p in s do
3     p.calls ← p.calls ∪ {f};
4     p.allCalls[0] ← p.allCalls[0] ∪ {f};          /* 0 is the top-level */
5     for (A, B, L) in p.callsPairs do
6       B ← B ∪ {f};
7       // rid of "large" critical sections
8       if |B| > d then p.callsPairs ← p.callsPairs \ {(A, B, L)};
9     end
10  end
11  return s;
12 end
```

Summary The implementation of the *summary* is given in Listing 5.7. It is of the type `Summary.t` (lines 3–4). It is composed of the modules defined on the lines above and earlier in Listing 5.6. The sub-module **Summary** also implements other functions not shown in the listing. Some of these functions are defined later. The summary is a structure that consists of the following fields:

- **atomicFunctions**: a set \mathcal{B} of sets of function calls, i.e., an element of 2^{2^Σ} . It captures sets of functions called atomically within an analysed function.
- **allCalls**: a list of sets of function calls, i.e., an element of $2^{\mathbb{N} \times 2^\Sigma}$. The purpose of this field is the same as the field **allCalls** of the abstract state, i.e., it is used to capture which functions were called at which levels of the call hierarchy. It is recorded to facilitate the derivation of atomic sets that show up higher in the call hierarchy. This field can be view as the C set from the proposal. However, it is extended to store the number of the level in the call hierarchy in which the functions were called, as proposed in Section 4.3.

```
1 module SS = Set.Make (S) (* a set of sets of strings *)
2 module Summary = struct
3   type t = (* a type of a summary *)
4     {atomicFunctions: SS.t; allCalls: S.t list}
5 end
```

Listing 5.7: A definition of the *summary* in Phase 1

The **Summary** module also implements the function **create**, which converts an *intraprocedural* abstract state to an *interprocedural* summary. This function is given in Algorithm 5.4. Its application can be seen in Listing 5.1 provided formerly.

Moreover, when calling an already analysed function, its summary is used to update an abstract state of an analysed function. This is done using the function **apply_summary**

Algorithm 5.4: Converting an abstract state to a *function summary*

Data: abstract state $s \in \text{Domain.t}$

```
1 def create(s):
2    $\mathcal{B} \leftarrow \emptyset$ ;
3    $C \leftarrow []$ ;
4   for  $p \in s$  do
5     for  $(A, B) \in p.\text{finalCallsPairs}$  do  $\mathcal{B} \leftarrow \mathcal{B} \cup \{B\}$ ;
6     for  $(i, C') \in p.\text{allCalls}$  do
7       if  $\exists C[i]$  then  $C[i] \leftarrow C[i] \cup C'$  else  $C[i] \leftarrow C'$ ;
8     end
9   end
10  return  $\{\text{atomicFunctions} = \mathcal{B}; \text{allCalls} = C\}$ ;
11 end
```

from Algorithm C.3 in Appendix C. On line 15, there is again used the parameter d to discard “large” critical sections. Another interesting part of the algorithm is the derivation of nested functions executed within the called function (lines 4–10). Here, the parameter r is used to limit the number of levels in the call hierarchy considered during analysing nested functions. How the function `apply_summary` is used in the analysis is shown in Listing 5.2 discussed previously.

5.1.1 Output of Phase 1

An output of the first phase of the analysis are *sets of functions* that were identified as those that are likely to be always executed atomically. The sets are given separately for each function in which they were identified. They are derived from the summaries of all analysed functions of a given program. At the end of the entire analysis, the sets are printed into a *text file* `atomic-sets` in an exact format. Each line of the file contains a sequence of detected atomic sets (separated by whitespace) within a particular function. A line starts with the name of the function followed by a colon and whitespace. Thus, whitespaces are significant. The set of atomic functions must be inside curly braces—a comma and whitespace separate individual functions inside the set. Empty lines and line comments (started with the character `#`) are ignored. Moreover, some statistical information about the number of atomic functions is printed on the last line. Listing 5.8 includes an example of such an output file.

```
functionA: {f1, f2, f3}
foo: {y, z}
bar_f: {a} {b, c}

# Number of (analysed functions; atomic sets; atomic functions): (3; 4; 8)
```

Listing 5.8: An example of an *output of the first phase* of the analysis (`atomic-sets` file)

This file is an *input file for the second phase* of the analysis. However, it may be manually edited by the user or created from scratch, but the format must be adhered to. The

derivation of the atomic sets and their printing to the output file is implemented in the `print_atomic_sets` function introduced in Algorithm 5.5.

Algorithm 5.5: *Printing atomic sets* from the summaries of all analysed functions

```

Data: all analysed functions  $F \in 2^\Sigma$ 
1 def print_atomic_sets( $F$ ):
2   for  $f \in F$  do
3      $\chi \leftarrow \text{read\_summary}(f)$ ;
4      $\text{printf}(\text{"\%s: "}, f)$ ;
5     // 'concat' concatenates elements of a set using a delimiter
6     for  $B \in \chi.\text{atomicFunctions}$  do  $\text{printf}(\text{"\%s} \text{"}, \text{concat}(B, ", "));$ 
7      $\text{print\_newline}()$ ;
8   end
9    $\text{print\_stats}(F)$ ;
10 end

```

5.2 Phase 2 — Detection of Atomicity Violations

The second phase of the analysis is enabled by a command-line argument `--atomicity-violations-only` (see the user manual for further details). This phase detects *atomicity violations*, i.e., violations of the atomic sets obtained from the first phase. The atomic sets are read from the file `atomic-sets`, described in the previous section. Therefore, this file must exist to run the second phase of the analysis successfully. This section covers the implementation of the abstract states and summaries and their usage in this phase. The output of this phase is a file containing reported atomicity violations, as discussed in Section 5.2.1.

Algorithm 5.6: Deriving *atomic pairs* from atomic sets

```

Data: atomic sets  $\mathcal{B} \in 2^{2^\Sigma}$ 
1 def create_atomic_pairs( $\mathcal{B}$ ):
2    $\Omega \leftarrow \emptyset$ ;
3   for  $B \in \mathcal{B}$  do
4      $B_{\text{list}} \leftarrow \text{Set.elements}(B)$ ;          /* convert a set to a list */
5     if  $|B| = 1$  then  $\Omega \leftarrow \Omega \cup \{(\varepsilon, B_{\text{list}}[0])\}$ ;
6     else
7       for  $i \leftarrow 0$  to  $|B| - 1$  do      /* create all 2-element variations */
8         for  $j \leftarrow i + 1$  to  $|B| - 1$  do
9            $\Omega \leftarrow \Omega \cup \{(B_{\text{list}}[i], B_{\text{list}}[j])\} \cup \{(B_{\text{list}}[j], B_{\text{list}}[i])\}$ ;
10        end
11      end
12    end
13  end
14  return  $\Omega$ ;
15 end

```

Atomic Pairs At the beginning of the second phase of the analysis, the input file is processed. All atomic sets from this file are transformed into *pairs* stored in an internal data structure $\Omega \in 2^{\Sigma \times \Sigma}$, which is a set of pairs of function calls that stores pairs that should be called atomically. Ω is globally accessible throughout the analysis.

The transformation from sets to pairs is done using the function `create_atomic_pairs` as demonstrated in Algorithm 5.6. In particular, to obtain the pairs, all possible pairs of functions are taken from the atomic sets, i.e., all *2-element variations* (lines 7–11). Besides, notice that single functions may also be considered when an atomic set contains only one function call (line 5).

```

1 type xy = string * string (* (x,y) pair *)
2 type xy_lock = xy * Lock.t (* (x,y) pair with a lock, i.e., (x,y, $\pi$ ,l)  $\in \lambda$  *)
3 module S_xy_lock = (* a set of the above *)
4   Set.Make(struct type t = xy_lock end)
5 type severity = Warning | Error (* severity  $\mathcal{S}$  of an atomicity violation *)
6 type v = (* a pair that violates atomicity, i.e., (x,y,c,S)  $\in \delta$  *)
7   xy * int * severity
8 module S_v = Set.Make(struct type t = v end) (* a set of the above *)
9 type p = (* an element of an abstract state *)
10   { lastPair: xy; violations: S_v.t; lockedLastPairs: S_xy_lock.t
11     ; guards: g }
12 module TSet = Set.Make(struct type t = p end) (* a set of the above *)
13 type t = TSet.t (* a type of the abstract state (an alias for TSet.t) *)
14 let initial : t = (* the initial abstract state of a function f *)
15   TSet.singleton (* an "empty" element (initial program path of f) *)
16     { lastPair= ("", "") (* an empty pair *)
17       ; violations= S_v.empty (* an empty set *)
18       ; lockedLastPairs= S_xy_lock.empty (* an empty set *)
19       ; guards= M_AP.empty (* an empty map *) }

```

Listing 5.9: A definition of the *abstract state* in Phase 2

Abstract State The implementation of the *abstract state* is illustrated in Listing 5.9. The abstract state is of the type `TSet.t` (`t` is an alias for it). `TSet` itself is defined on line 12. It is a module representing a set of structures `p` (lines 9–11). Each element corresponds to a single program path. Fields of this structure are composed of types and modules defined on the lines above (some types are defined already in the previous section). Lines 14–19 define the *initial abstract state*. The fields are as follows:

- **lastPair**: a *pair of function calls*, i.e., $\Sigma \times \Sigma$. It captures the (x, y) pair of the most recent calls of functions performed on a program path being explored. It is used for detecting whether this pair violates atomicity.
- **violations**: the *set δ of so far detected pairs that violate atomicity* on particular lines of code with the *severity* of the violation, i.e., an element of $2^{\Sigma \times \Sigma \times \mathbb{N} \times \mathcal{S}}$. In the end, this set is used for reporting atomicity violations.

- **lockedLastPairs**: the set λ of the most recent function calls with their locks, where a lock \mathcal{L} is of the type `Lock.t`. In fact, it is a set of tuples (x, y, π, l) , i.e., an element of $2^{\Sigma \times \Sigma \times \Pi \times \mathbb{N}^\Gamma}$. This set is used for checking whether the last pair is locked.
- **guards**: the same object as in the first phase of the analysis, i.e., a map $\Pi \rightarrow 2^{\mathbb{N} \times \Pi}$. It maps a lock guard's access path to a list of access paths of associated locks. It can be viewed as a set of the pairs $(\pi_g, L) \in \Pi \times 2^\Pi$ from the proposal.

The main functions of this phase's abstract domain that work with abstract states are `apply_call`, `apply_locks`, and `apply_unlocks`. Earlier, in Listing 5.2, there is shown how these functions are applied.

The function `apply_call` from Algorithm 5.7 demonstrates how an abstract state is updated when some classical function is invoked. At first, this function pushes the called function to the `lastPair` field and all locked pairs in the field `lockedLastPairs` (lines 4–9). Then, it is checked whether there are any new atomicity violations (lines 12–18). Note that on lines 13 and 17, the created violations are labelled as **Error**. Thus, they are considered as *global atomicity violations* so far. Later in the analysis—when analysing functions higher in the call hierarchy—these violations might be changed to **Warning** if this function is invoked under a lock. This is further described later on.

Algorithm 5.7: Updating an abstract state after a *call of a function*

Data: called function $f \in \Sigma$; actual line of code $c \in \mathbb{N}$; abstract state $s \in \text{Domain.t}$; atomic pairs $\Omega \in 2^{\Sigma \times \Sigma}$

```

1 def apply_call(f, c, s):
2   for p in s do
3     // 'snd' returns the second element of a pair
4     p.lastPair ← (snd(p.lastPair), f);    /* push f to the last pair */
5     for (xℒ, yℒ, ℒ) in p.lockedLastPairs do
6       // push f to all the locked pairs
7       xℒ ← yℒ;
8       yℒ ← f;
9     end
10    (x, y) ← p.lastPair;
11    // check atomicity violation
12    if (x, y) ∈ Ω and ∄(x, y, _) ∈ p.lockedLastPairs then
13      | p.violations ← p.violations ∪ {(x, y, c, Error)};
14    end
15    // check violation also for a single function
16    if (ε, y) ∈ Ω and ∄(_, y, _) ∈ p.lockedLastPairs then
17      | p.violations ← p.violations ∪ {(ε, y, c, Error)};
18    end
19  end
20  return s;
21 end

```

Furthermore, the `apply_locks` and `apply_unlocks` functions' implementation is illustrated using the algorithms from Appendix C, namely, Algorithms C.5 and C.4, respectively.

Summary The implementation of the *summary* in this phase of the analysis is given in Listing 5.10. It is of the type `Summary.t` (line 2). It is composed of the modules defined earlier in Listing 5.9. The sub-module `Summary` also implements other functions not shown in the listing. Some of these functions are defined later. The summary is a structure that contains the following fields:

- **violations:** the set Δ of all detected pairs that violate atomicity on particular lines of code with the severity of the violation, i.e., an element of $2^{\Sigma \times \Sigma \times \mathbb{N} \times \mathcal{S}}$. It is constructed from the sets δ of an abstract state of a given function.

```

1 module Summary = struct
2   type t = {violations: S_v.t} (* a type of a summary *)
3 end

```

Listing 5.10: A definition of the *summary* in Phase 2

The `Summary` module also implements the function `create`, which converts an *intraprocedural* abstract state to an *interprocedural* summary. This function is given in Algorithm 5.8. It simply takes atomicity violations from an abstract state and puts them into a single set. An application of this function can be seen in Listing 5.1 provided at the beginning of this chapter.

Algorithm 5.8: Converting an abstract state to a *function summary*

Data: abstract state $s \in \text{Domain.t}$

```

1 def create(s):
2    $\Delta \leftarrow \emptyset$ ;
3   for  $p \in s$  do  $\Delta \leftarrow \Delta \cup p.violations$ ;
4   return {violations =  $\Delta$ };
5 end

```

Furthermore, when calling an already analysed function, its summary is used to update an abstract state of an analysed function. This is done using the function `apply_summary` from Algorithm C.6 in Appendix C. The algorithm works as follows. When the abstract state is considered inside an atomic block, all atomicity violations from the summary of the called function are passed to atomicity violations of the abstract state as warnings (lines 8–10). Otherwise, they are passed unchanged (line 5). This means that *global* atomicity violations may be changed to *local* atomicity violations of the called functions. As a consequence, these local atomicity violations will be reported as warnings, as discussed in the next section. How the function `apply_summary` is used in the second phase of the analysis is shown in Listing 5.2.

5.2.1 Output of Phase 2

An output of the second phase of the analysis is the final reported *atomicity violations*. The violations are printed to the standard output. Moreover, they are stored in both plain text format (`infer-out/report.txt`) and JSON format (`infer-out/report.json`). An

example of text output is provided in Listing 5.11. Notice that in JSON format, there are available more information about reported violations, e.g., an *error trace*.

```
#0
analysed_program.c:36: error: Atomicity Violation
  Atomicity Violation! - Functions 'foo' and 'bar' should be called atomically.
  34. void do_something()
  35. {
  36.   foo(); bar();
      ^
  37. }
  38.

#1
analysed_program.c:58: warning: Local Atomicity Violation
  Local Atomicity Violation! - Functions 'x' and 'y' should be called atomically.
  56. if (x() > 6)
  57. {
  58.   y(args);
      ^
  59. }
  60.

Found 2 issues

                                Issue Type(ISSUED_TYPE_ID): #
      Local Atomicity Violation(ATOMICITY_VIOLATION_WARNING): 1
      Atomicity Violation(ATOMICITY_VIOLATION_ERROR): 1
```

Listing 5.11: An example of an *output of the second phase* of the analysis (`report.txt` file)

The process of reporting atomicity violations is implemented using the function from Algorithm 5.9. As line 3 indicates, atomicity violations are reported only from *top-level functions* of the call tree. These functions also contain atomicity violations from lower levels, because the summaries are passed *upwards along the call tree* during the analysis. An error is reported a little differently for local and global violations (lines 6–9). The reporting itself is implemented using the **Reporting** module provided by the Facebook Infer framework.

5.3 Support for New Languages and Locks

Besides all the proposed and implemented extensions described in the previous sections/chapters, Atomer was extended to support analysis of *C++* and *Java* programs that the first version of Atomer did not support. As already said in Section 2.3, in general, Facebook Infer can analyse programs written in C, C++, Java, Objective-C (and C# [40]). The Facebook Infer’s *frontend* compiles input program into the *SIL language* and represents them as a CFG. Individual analyses are then performed on SIL. It is thus possible to extend the frontend for supporting other languages. However, in practice, there are not negligible differences among the ways programs from different languages look like in SIL. Therefore, individual non-trivial analysers have to be adapted for specific languages.

Algorithm 5.9: *Reporting atomicity violations* from the summaries of all analysed functions

```

Data: all analysed functions  $F \in 2^\Sigma$ 
1 def report_atomicity_violations( $F$ ):
2   for  $f \in F$  do
3     if is_top_level( $f$ ) then
4        $\chi \leftarrow \text{read\_summary}(f)$ ;
5       for  $(x, y, c, S) \in \chi.\text{violations}$  do
6         switch  $S$  do
7           case Error do report_global_violation( $x, y, c$ );
8           case Warning do report_local_violation( $x, y, c$ );
9         end
10      end
11    end
12  end
13 end

```

The basic version of Atomer supported only C language using *PThread mutex locks*. The new version of Atomer was adjusted to support also other C locks and even more complicated *mutual exclusion* C++ and Java mechanisms. Facebook Infer’s built-in modules implemented in `infer/src/absint/ConcurrencyModels.ml[i]` are used to identify locks for different languages. Mainly, the function `get_lock_effect` is used. The particular usage of the function can be seen in Listing 5.2 at the beginning of the chapter.

In the new version of Atomer, the `ConcurrencyModels` module has been significantly extended to support a broader range of locks. Most importantly, there is significantly better and more accurate work with *lock guards* (and their *locking strategies*) and support for *re-entrant locks*. Further, there is new/extended support, e.g., for *read/write locks*, *timed locks*, *shared locks*, *spinlocks*, locking multiple locks *at once* using special functions, etc. In particular, the following locking mechanisms are supported in the new version of Atomer:

- all C/C++ locks from the `pthread.h`⁴ library, including, e.g., spinlocks;
- C++ locks (guards) from standard libraries `std::mutex` and `std::shared_mutex`⁵, including, e.g., C++ 17 `scoped_lock` or `shared_mutex`;
- C++ locks (guards) from libraries of Apache Thrift⁶, Boost⁷, and Facebook Folly⁸;
- Java locks from the standard package `java.util.concurrent.locks`⁹; and
- Java lock guards (*monitors*), i.e., *synchronised* blocks and methods using the `synchronized` keyword.

⁴POSIX thread C/C++ library `pthread.h`: <https://pubs.opengroup.org/onlinepubs/9699919799/basedefs/pthread.h.html>.

⁵Standard C++ *multi-threading* libraries: <https://en.cppreference.com/w/cpp/thread>.

⁶C++ framework **Apache Thrift**: <https://thrift.apache.org>.

⁷Third party C++ library **Boost**: <https://www.boost.org>.

⁸Third party C++ library **Facebook Folly**: <https://github.com/facebook/folly>.

⁹Java (SE 11) package `java.util.concurrent.locks`: <https://docs.oracle.com/en/java/javase/11/docs/api/java.base/java/util/concurrent/locks/package-summary.html>.

Chapter 6

Experimental Evaluation of the New Version of Atomer

This chapter is devoted to *testing* and *experimental evaluating* the new version of *Atomer* as proposed and implemented in Chapters 4 and 5, respectively. Each Atomer’s enhancement/extension was tested independently as soon as it had been implemented. It was tested on suitable programs created for testing purposes and got inspired by test programs from [18]. Section 6.1 includes the evaluation of analysis *scalability*. Further, Section 6.2 shows experiments performed on test programs derived from the static analyser *Gluon*. Moreover, Section 6.3 is about the experimental evaluation of the new version of Atomer on publicly available *real-life complex* programs. Finally, Section 6.4 concludes the experimental evaluation and discusses future work.

Testing on Hand-Crafted Examples

The correct behaviour of the first version of Atomer was already tested in [18] on a suitable test suite containing C programs that use *PThread locks*. Test programs from this test suite include sequences of function calls inside and outside *atomic blocks*, *not paired lock/unlock calls*, iteration, selection, and *nested function calls* inside atomic sections. The programs were designed to check all essential aspects of Atomer’s analysis, i.e., to check whether various parts of *abstract domains* work well w.r.t. to the proposal.

The implementation of the new version of Atomer was, at first, tested on the test suite from [18]. However, the reference outputs for these test programs were modified to adapt the newly implemented features, e.g., the approximation of working with *sequences of function calls* by working with *sets of function calls* has to be reflected in the reference outputs. Further, new test programs were derived from those in [18], and they were changed in order to check all aspects of the extensions and enhancements implemented in the new version of Atomer. During testing, various kinds of different lock mechanisms for C/C++/Java were used.

This way, the correct functioning of the analysis of Atomer’s new version has been validated (w.r.t. the proposal). All of these testing programs are available in the attached memory media; see Appendix A.

6.1 Scalability Benchmark

The *scalability* of the analysis has been tested on a subset of a *publicly available benchmark* from [21]. It consists of *real-life low-level complex concurrent C* programs derived from the Debian GNU/Linux distribution. The entire benchmark was initially used for an experimental evaluation of Daniel Kroening’s static deadlock analyser for C/PThreads implemented in the CPROVER framework. However, it can also be used for the evaluation of Atomer’s scalability. For the evaluation, 7 programs with deadlocks and 54 deadlock-free programs (806,431 lines of code in total) were used.

The experiments¹ were run in a *Docker container* on Windows 10 (WSL 2) with a 3.3 GHz Intel® Core™ i5-2500K 64-bit processor and 6 GB RAM running the Debian GNU/Linux x86-64 10.9 (Buster) operating system. Table 6.1 shows aggregated results of the evaluation. There are average and total times of analyses for both phases of the analysis for both the basic version of Atomer (Atomer v1.0.0, i.e., the version before the *approximation with sets* proposed in the thesis) and the new version (Atomer v2.0.0, i.e., the version after the approximation). Also, there is shown the number of *timeouts* of individual analyses where the timeout \top was set to 10 minutes. It is evident that, on average, the new version of Atomer is about *two times faster*.

Table 6.1: *Aggregated results of the evaluation of Atomer’s analysis scalability*

	Atomer v1.0.0		Atomer v2.0.0	
	Phase 1	Phase 2	Phase 1	Phase 2
Average Time (s)	70.98	109.11	37.96	50.93
Total Time (s)	4,117	5,892	2,164	2,750
Timeouts	4	9	3	4

For a closer look, Table 6.2 provides an overview of analyses’ times of selected programs from the benchmark. Above that, there are **tgrep** and **sort** from GNU Core Utilities that also use PThreads. It can be seen that the approximation decreased the analysis time a lot for some of these programs. Nonetheless, in one case, the timeout was also exceeded in the new version of Atomer. This shows that the analysis may still take a quite long time in some cases even after the approximation. The reason is that in these low-level programs, there are vast and complicated functions (thousands lines of code) that contain *in-lined code*.

6.2 Evaluation on Validation Programs Derived from Gluon

Gluon is a prototype tool for static analysis of *contracts for concurrency* in Java programs developed in [13, 37] (also discussed in Section 3.1). The functionality of Gluon was validated on a set of small benchmarking programs with known atomicity violations, which can be seen as contract violations. These programs were adapted from [2, 3, 41], where they are typically used to evaluate atomicity violation detection methods. In [13, 37], these

¹Performed **scalability experiments** are available at GitHub: <https://github.com/harmim/vut-dip/tree/master/benchmarks/scalability>.

Table 6.2: *Individual results* of the evaluation of Atomer’s *analysis scalability* ($\top > 600$ s)

Program	LOC	Atomer v1.0.0		Atomer v2.0.0	
		Phase 1 Time (s)	Phase 2 Time (s)	Phase 1 Time (s)	Phase 2 Time (s)
btscanner 2.1	9,142	\top	N/A	47.33	69.04
daemon (libslack) 0.6.4	42,857	280.75	40.85	11.40	10.43
crossfire-client-gtk2 1.71	41,185	1.74	\top	1.71	14.58
c-icap 0.4.2	39,254	154.86	20.70	2.37	6.14
hmmer2 (hmmcalibrate) 2.3.2	38,301	109.79	\top	10.07	\top
towitoko 2.0.7	12,339	38.23	77.50	2.27	10.69
mesa-demos 8.3.0	2,424	397.10	\top	46.50	5.83
freecell-solver 3.26.0	15,408	37.28	0.78	3.45	9.81
tgrep 1	3,190	240.36	\top	0.92	0.74
sort (GNU Coreutils)	7,523	5.33	0.39	1.49	0.91

programs were appropriately redesigned, and the necessary contracts for each program were created.

Moreover, in [13], to validate a *dynamic approach* for contract violations, a subset of the benchmarking programs was rewritten to a C++ version. This subset was also used to validate the new version of Atomer. Validation results² are given in Table 6.3. The table shows the number of atomicity violations that can really occur in a given program (*Real Atomicity Violations*), the number of the real violations that were found by the new version of Atomer (*True Positives*), and the number of false positives (*False Positives*).

Table 6.3: Atomer’s *validation results* on test programs derived from *Gluon* [13, 37]

Benchmark	LOC	Real Atomicity Violations	True Positives	False Positives
Account [41]	56	1	1	0
Coord03 [2]	119	1	1	0
Coord04 [3]	55	1	1	0
Local [2]	30	1	1	0
NASA [2]	97	1	1	0

The experiments were performed such that the contracts from [13] were manually given to the second phase of Atomer’s analysis as *atomic sets*. As the table shows, Atomer detected *all known atomicity violations* with the *absence of false positives*. Times of the analysis are not included in the table because all the programs were analysed in less than a second.

Furthermore, the author of this thesis created a *correct version* of each benchmark, i.e., a version where atomicity violations are fixed by adding appropriate locks. Then, it was shown that the first phase of the new version of Atomer was able to derive suitable contracts *automatically*. Subsequently, Atomer found no atomicity violations in such corrected programs.

²Experiments on validation programs derived from *Gluon*: <https://github.com/harmim/vut-pp1/tree/master/testing-programs/contracts-first-experiments-atomer>.

6.3 Experiments with Real-Life Programs

Since the new version of Atomer supports the analysis of Java programs, two *open-source real-life extensive* (both $\sim 250,000$ lines of code) Java programs were analysed — *Apache Cassandra*³ 3.11 and *Apache Tomcat*⁴ 8.5. In [13], there were reported several *atomicity-related bugs* in these programs. It turned out that the reported bugs were real atomicity violation errors, and they were later fixed.

When Atomer first analysed these programs (at that time, without most of the extensions presented in this thesis), the bugs were successfully *rediscovered*, but quite some *false alarms* were reported. However, after all the improvements proposed in the thesis were implemented, the number of false alarms was *significantly reduced*. In particular, the support for *re-entrant locks* increased the analysis’ precision a lot because these types of locks appear pretty often in these programs. Moreover, many “large” atomic sections in these programs dramatically increase the number of reported false alarms. Therefore, the *parameters of the analysis* presented in the thesis were used. In particular, the maximum length of critical sections was limited with the parameter d set to 20, and the number of levels considered during analysing nested functions was limited with the parameter r set to 10. Further, also using the parametrisation of the analysis, several “non-critical” functions were ignored during the analysis (e.g., `String.format`, `*.toString`, `*.toArray`, `Log.debug`, `Integer.valueOf`, etc.).

With all these Atomer’s extensions, when analysing one of the source files of Tomcat where a real atomicity violation was reported before (i.e., `org.apache.catalina.core.StandardContext` — this package itself contains 6,684 lines of code, and the packages it depends on, that must have been analysed too, contain tens of thousands lines of code), the number of reported errors decreased from ~ 800 to 228. Then, when analysing one of the source files of Cassandra where a real atomicity violation was reported before (i.e., `org.apache.cassandra.streaming.StreamSession` — this package itself contains 1,014 lines of code, and the packages it depends on, that must have been analysed too, contain tens of thousands lines of code), the number of reported errors decreased from ~ 700 to 112.

Obviously, most of the previously reported errors were false alarms. However, it is still challenging to say which of these errors are real atomicity violations. Indeed, the author suspects some of the warnings correspond to real errors, but so far, he has not managed to confirm that.

6.4 Summary of the Evaluation and Future Work

Using small *hand-crafted programs*, it was *successfully validated* that the implemented new version of Atomer works correctly w.r.t. the proposal. This validation includes all *scalability/precision* enhancements as well as support for other programming languages and *advanced locking mechanisms* that were not supported in the first version of Atomer.

³Open-source NoSQL database system **Apache Cassandra**: <https://cassandra.apache.org>.

⁴Open-source Java HTTP web server **Apache Tomcat**: <https://tomcat.apache.org>.

Scalability testing in Section 6.1 shows that the new version of Atomer can analyse quite *extensive real-life low-level* programs in a reasonable time. Moreover, it was proven that the new version of Atomer is *significantly faster* than the first version implemented in [18]. This improvement is caused by the *approximation* of working with *sequences of function calls* by working with *sets of function calls*. Regarding performance, Atomer may be directly usable for large real-life projects. However, the performance of Atomer strongly depends on the size of *critical sections* (i.e., the number of function calls inside critical sections). Thus, it is needed to filter out “large” critical sections using the proposed parameters in some cases.

Sections 6.2 and 6.3 demonstrate that the new version of Atomer can also reveal *real atomicity-related bugs* in C++ and Java programs. In the case of the complicated and large Java programs from Section 6.3, it was shown that the accuracy of the new version of Atomer was *significantly improved*. However, unfortunately, quite some false alarms are still reported along with real errors. It can be demanding to properly classify which of the reported errors are real errors and false alarms. This can make the analysis unusable on some kinds of real-world programs. But still, it is possible to find appropriate values of the proposed parameters for a particular program to obtain good results. For instance, in this way, one can specify several specific “critical functions” that Atomer should focus on. This should often lead to much better results.

The above implies that future work will focus mainly on further improvements of the accuracy of the analysis. Besides, the analysis results may be used as an input for *dynamic analysis*, checking whether the atomicity violations are real errors. For example, one could use the ANaConDA [14] dynamic analyser, which uses *noise-based testing* with *extrapolated checking* for violations of contracts for concurrency. ANaConDA could be instructed to concentrate its analysis and *noise injection* to those sets whose atomicity was found broken. Furthermore, it seems promising to consider *formal parameters* and *distinguish the context* of called functions during the analysis to reduce false alarms.

Another exciting idea is to use *machine learning* to learn appropriate values of the analysis’ parameters (introduced in this thesis) for particular programs. *Automation* of the analysis can be preserved using this approach, and the precision of the analysis may be high enough (however, the precision would highly depend on *training data*). Another solution to reducing the number of false alarms is to do some *statistical ranking* of reported errors and display only the most critical errors to the user.

During experimenting with Tomcat, the author of this thesis found out that several reported errors were due to calling functions with parameters that were *local* variables (i.e., *non-shared* variables). These cases cannot lead to real atomicity violations. Therefore, it would be helpful to statically identify whether function parameters are shared variables. If they are not, calls of such functions should be ignored.

Chapter 7

Conclusion

The thesis started by describing the principles of program analysis with a focus on *static analysis* and *abstract interpretation*. Further, *Facebook Infer* was described—a concrete static analysis framework that uses abstract interpretation. Next, there were described *contracts for concurrency*. The major part of the thesis then aimed to describe static analyser *Atomer*¹—proposed and implemented within the author’s bachelor’s thesis [18]—implemented as a Facebook Infer’s module for detection of *atomicity violations* in *multi-threaded programs*. It was described its limitations, and thereafter, it was introduced the proposal and implementation of its extensions and improvements. Lastly, the experimental evaluation of the new features and improvements was discussed together with possible future work.

The first version of *Atomer* works on the level of *sequences of function calls*. It is based on the assumption that sequences of function calls executed *atomically once* should probably be executed *always atomically*, and it naturally works with sequences of calls. In the thesis, to improve *scalability*, the use of sequences was *approximated* by *sets*. Furthermore, several new features were implemented, e.g., support for *C++* and *Java*, including various *advanced kinds of locks* these languages offer (such as *re-entrant locks* or *lock guards*); or a more precise way of *distinguishing between different lock instances*. Moreover, the analysis has been *parametrised*, so the user can provide some input values to the analysis to make the analysis more accurate.

The introduced enhancements were *successfully tested* and *experimentally evaluated* also on *extensive real-life software* where *real bugs* were *successfully rediscovered*. It turned out that such innovations improved the *accuracy* and *scalability* of the analysis a lot. However, *Atomer*’s accuracy can be further increased. There are still some other improvements and ideas to work on, for instance, considering *formal parameters* and *distinguishing the context* of called functions; or combinations with a *dynamic analysis*. Another interesting idea is to use *machine learning* to learn appropriate values of the analysis’ parameters for particular programs. Further, it is needed to perform more experiments on real-life programs to find and report *new bugs*.

It is expected that the work on this project will continue within the VeriFIT group at FIT BUT. The preliminary results of the thesis were presented at the Excel@FIT’21 conference, where it *won two awards*.

¹The implementation of a new version of **Atomer** is available at GitHub as an *open-source* repository: <https://github.com/harmim/infer/tree/atomer-v2.0.0>.

Bibliography

- [1] ALLEN, F. E. Control Flow Analysis. In: *Proceedings of a Symposium on Compiler Optimization*. Urbana-Champaign, Illinois: ACM, New York, NY, USA, July 1970, p. 1–19. DOI: 10.1145/800028.808479. ISBN 9781450373869.
- [2] ARTHO, C., HAVELUND, K. and BIERE, A. High-level data races. *Software Testing, Verification and Reliability*. November 2003, vol. 13, no. 4, p. 207–227. DOI: 10.1002/stvr.281. ISSN 0960-0833.
- [3] ARTHO, C., HAVELUND, K. and BIERE, A. Using Block-Local Atomicity to Detect Stale-Value Concurrency Errors. In: WANG, F., ed. *Automated Technology for Verification and Analysis*. Taipei, Taiwan: Springer, Berlin, Heidelberg, October 2004, vol. 3299, p. 150–164. ATVA’04. Lecture Notes in Computer Science. DOI: 10.1007/978-3-540-30476-0_16. ISBN 978-3-540-23610-8.
- [4] BLACKSHEAR, S. *Getting the most out of static analyzers* [online]. San Jose Convention Center: The @Scale Conference, 2. September 2016 [cit. 2021-04-21]. Available at: <https://atscaleconference.com/videos/getting-the-most-out-of-static-analyzers>.
- [5] BLACKSHEAR, S., GOROGIANNIS, N., O’HEARN, P. W. and SERGEY, I. RacerD: Compositional Static Race Detection. *Proceedings of the ACM on Programming Languages*. New York, NY, USA: Association for Computing Machinery. October 2018, vol. 2, OOPSLA’18, p. 144:1–144:28. DOI: 10.1145/3276514. ISSN 2475-1421.
- [6] BLACKSHEAR, S. and O’HEARN, P. W. Facebook Engineering. *Open-sourcing RacerD: Fast static race detection at scale* [online]. 19. October 2017 [cit. 2021-04-21]. Available at: <https://code.fb.com/android/open-sourcing-racerd-fast-static-race-detection-at-scale>.
- [7] CALCAGNO, C., DISTEFANO, D., O’HEARN, P. W. and YANG, H. Compositional Shape Analysis by Means of Bi-Abduction. In: *Proceedings of the 36th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. Savannah, GA, USA: ACM, New York, NY, USA, January 2009, p. 289–300. POPL’09. DOI: 10.1145/1480881.1480917. ISBN 978-1-60558-379-2.
- [8] COUSOT, P. *Abstract Interpretation* [online]. Revised 5. August 2008 [cit. 2021-04-21]. Available at: <https://www.di.ens.fr/~cousot/AI>.
- [9] COUSOT, P. *Abstract Interpretation in a Nutshell* [online]. [cit. 2021-04-21]. Available at: <https://www.di.ens.fr/~cousot/AI/IntroAbsInt.html>.

- [10] COUSOT, P. Abstract Interpretation Based Formal Methods and Future Challenges, invited paper. In: WILHELM, R., ed. « *Informatics — 10 Years Back, 10 Years Ahead* ». Berlin, Heidelberg: Springer Berlin Heidelberg, March 2001, vol. 2000, p. 138–156. Lecture Notes in Computer Science. DOI: 10.1007/3-540-44577-3_10. ISBN 978-3-540-44577-7.
- [11] COUSOT, P. and COUSOT, R. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In: *Conference Record of the Fourth Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. Los Angeles, California: ACM Press, New York, NY, January 1977, p. 238–252. POPL’77. DOI: 10.1145/512950.512973. ISBN 9781450373500.
- [12] COUSOT, P. and COUSOT, R. Comparing the Galois Connection and Widening/Narrowing Approaches to Abstract Interpretation, invited paper. In: BRUYNOOGHE, M. and WIRSING, M., ed. *Proceedings of the International Workshop Programming Language Implementation and Logic Programming*. Springer-Verlag, Berlin, Germany, January 1992, p. 269–295. PLILP’92. Leuven, Belgium, 13–17 August 1992, Lecture Notes in Computer Science 631. DOI: 10.1007/3-540-55844-6_101. ISBN 978-3-540-47297-1.
- [13] DIAS, R. J., FERREIRA, C., FIEDOR, J., LOURENÇO, J. M., SMRČKA, A., SOUSA, D. G. and VOJNAR, T. Verifying Concurrent Programs Using Contracts. In: *10th IEEE International Conference on Software Testing, Verification and Validation*. Tokyo, Japan: IEEE Computer Society, Los Alamitos, CA, USA, March 2017, p. 196–206. ICST’17. DOI: 10.1109/ICST.2017.25. ISBN 9781509060313.
- [14] FIEDOR, J., MUŽIKOVSKÁ, M., SMRČKA, A., VAŠÍČEK, O. and VOJNAR, T. Advances in the ANaConDA Framework for Dynamic Analysis and Testing of Concurrent C/C++ Programs. In: *Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis*. Amsterdam, Netherlands: Association for Computing Machinery, New York, NY, USA, July 2018, p. 356–359. ISSTA’18. DOI: 10.1145/3213846.3229505. ISBN 978-1-4503-5699-2.
- [15] FLANAGAN, C. and FREUND, S. N. Atomizer: A Dynamic Atomicity Checker for Multithreaded Programs. In: *Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. Venice, Italy: Association for Computing Machinery, New York, NY, USA, January 2004, p. 256–267. POPL’04. DOI: 10.1145/964001.964023. ISBN 158113729X.
- [16] FLANAGAN, C. and FREUND, S. N. FastTrack: Efficient and Precise Dynamic Race Detection. In: *Proceedings of the 30th ACM SIGPLAN Conference on Programming Language Design and Implementation*. Dublin, Ireland: Association for Computing Machinery, New York, NY, USA, June 2009, p. 121–133. PLDI’09. DOI: 10.1145/1542476.1542490. ISBN 9781605583921.
- [17] GOROGIANNIS, N., O’HEARN, P. W. and SERGEY, I. A True Positives Theorem for a Static Race Detector. *Proceedings of ACM Programming Languages*. New York, NY, USA: Association for Computing Machinery. January 2019, vol. 3, POPL’19, p. 57:1–57:29. DOI: 10.1145/3290370. ISSN 2475-1421.

- [18] HARMIM, D. *Static Analysis Using Facebook Infer to Find Atomicity Violations*. Brno, CZ, 2019. Bachelor’s thesis. Brno University of Technology, Faculty of Information Technology. Department of Intelligent Systems. Supervisor VOJNAR, T. Available at: <https://www.fit.vut.cz/study/thesis/21689>.
- [19] HARMIM, D. Advanced Static Analysis of Atomicity in Concurrent Programs through Facebook Infer. In: *Proceedings of the Excel@FIT’21*. Brno, CZ: Brno University of Technology, Faculty of Information Technology, 2021. Available at: <http://excel.fit.vutbr.cz/submissions/2021/057/57.pdf>.
- [20] HOARE, C. A. R. An Axiomatic Basis for Computer Programming. *Commun. ACM*. New York, NY, USA: Association for Computing Machinery. October 1969, vol. 12, no. 10, p. 576–580. DOI: 10.1145/363235.363259. ISSN 0001-0782.
- [21] KROENING, D., POETZL, D., SCHRAMEL, P. and WACHTER, B. Sound Static Deadlock Analysis for C/Pthreads. In: *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*. Singapore, Singapore: ACM, New York, NY, USA, August 2016, p. 379–390. ASE’16. DOI: 10.1145/2970276.2970309. ISBN 978-1-4503-3845-5.
- [22] KŘENA, B. and VOJNAR, T. Automated Formal Analysis and Verification: An Overview. *International Journal of General Systems*. Taylor & Francis. November 2012, vol. 42, no. 4, p. 335–365. DOI: 10.1080/03081079.2012.757437. ISSN 0308-1079.
- [23] LAMPORT, L. Time, Clocks, and the Ordering of Events in a Distributed System. *Commun. ACM*. Association for Computing Machinery, New York, NY, USA. July 1978, vol. 21, no. 7, p. 558–565. DOI: 10.1145/359545.359563. ISSN 0001-0782.
- [24] LENGÁL, O. and VOJNAR, T. *Abstract Interpretation. Lecture Notes in Static Analysis and Verification*. Brno, CZ: Brno University of Technology, Faculty of Information Technology, 2020. Available at: <https://www.fit.vutbr.cz/study/courses/SAV/public/Lectures/sav-lecture-05-ai.pdf>.
- [25] LERCH, J., SPÄTH, J., BODDEN, E. and MEZINI, M. Access-Path Abstraction: Scaling Field-Sensitive Data-Flow Analysis with Unbounded Access Paths (T). In: *30th IEEE/ACM International Conference on Automated Software Engineering*. Lincoln, NE, USA: IEEE Computer Society, Los Alamitos, CA, USA, November 2015, p. 619–629. ASE’15. DOI: 10.1109/ASE.2015.9. ISBN 978-1-5090-0025-8.
- [26] MARCIN, V. *Static Analysis Using Facebook Infer Focused on Deadlock Detection*. Brno, CZ, 2019. Bachelor’s thesis. Brno University of Technology, Faculty of Information Technology. Department of Intelligent Systems. Supervisor VOJNAR, T. Available at: <https://www.fit.vut.cz/study/thesis/21920>.
- [27] MEYER, B. Applying “Design by Contract”. *Computer*. Washington, DC, USA: IEEE Computer Society Press. October 1992, vol. 25, no. 10, p. 40–51. DOI: 10.1109/2.161279. ISSN 0018-9162.
- [28] MINSKY, Y., MADHAVAPEDDY, A. and HICKEY, J. *Real World OCaml: Functional Programming for the Masses*. 1st ed. Sebastopol, CA: O’Reilly Media, 2013. ISBN 978-1-449-32391-2.

- [29] MUŽIKOVSKÁ, M. *Towards Parameterized Contract Validator in ANaConDA Framework*. Brno, CZ, 2018. Bachelor's thesis. Brno University of Technology, Faculty of Information Technology. Department of Intelligent Systems. Supervisor SMRČKA, A. Available at: <https://www.fit.vut.cz/study/thesis/20642>.
- [30] MØLLER, A. and SCHWARTZBACH, I. M. *Static Program Analysis*. Department of Computer Science, Aarhus University, November 2020. Available at: <https://cs.au.dk/~amoeller/spa>.
- [31] NIELSON, F., NIELSON, H. R. and HANKIN, C. *Principles of Program Analysis*. 2nd ed. Berlin, Heidelberg: Springer Berlin Heidelberg, January 2005. ISBN 978-3-642-08474-4.
- [32] REPS, T., HORWITZ, S. and SAGIV, M. Precise Interprocedural Dataflow Analysis via Graph Reachability. In: *Proceedings of the 22nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. San Francisco, California, USA: ACM, New York, NY, USA, January 1995, p. 49–61. POPL'95. DOI: 10.1145/199448.199462. ISBN 0-89791-692-1.
- [33] RICE, H. G. Classes of Recursively Enumerable Sets and Their Decision Problems. *Transactions of the American Mathematical Society*. 1953, vol. 74, no. 2, p. 358–366. DOI: 10.1090/s0002-9947-1953-0053041-6. ISSN 0002-9947.
- [34] RIVAL, X. and KWANGKEUN, Y. *Introduction to Static Analysis: An Abstract Interpretation Perspective*. 1st ed. Cambridge: The MIT Press, February 2020. ISBN 978-0-262-04341-0.
- [35] SAVAGE, S., BURROWS, M., NELSON, G., SOBALVARRO, P. and ANDERSON, T. Eraser: A Dynamic Data Race Detector for Multi-Threaded Programs. In: *Proceedings of the Sixteenth ACM Symposium on Operating Systems Principles*. Saint Malo, France: Association for Computing Machinery, New York, NY, USA, October 1997, p. 27–37. SOSP'97. DOI: 10.1145/268998.266641. ISBN 0897919165.
- [36] SHARIR, M. and PNUELI, A. Two Approaches to Interprocedural Data Flow Analysis. In: MUCHNICK, S. S. and JONES, N. D., ed. *Program Flow Analysis: Theory and Applications*. Prentice Hall Professional Technical Reference, January 1981, chap. 7, p. 189–211. ISBN 0137296819.
- [37] SOUSA, D. G., DIAS, R. J., FERREIRA, C. and LOURENÇO, J. M. Preventing Atomicity Violations with Contracts. *CoRR*. Ithaca, New York, USA: Cornell University Library, arXiv.org. May 2015, abs/1505.02951. ISSN 2331-8422.
- [38] TARSKI, A. A Lattice-Theoretical Fixpoint Theorem and Its Applications. *Pacific Journal of Mathematics*. June 1955, vol. 5, no. 2, p. 285–309. DOI: 10.2140/pjm.1955.5.285. ISSN 00308730.
- [39] VALMARI, A. The State Explosion Problem. In: REISIG, W. and ROZENBERG, G., ed. *Lectures on Petri Nets I: Basic Models*. Springer, Berlin, Heidelberg, 1998, vol. 1491, p. 429–528. Lecture Notes in Computer Science. DOI: 10.1007/3-540-65306-6_21. ISBN 978-3-540-65306-6.

- [40] VILLARD, J. Facebook Engineering. *Infer powering Microsoft's Infer#, a new static analyzer for C#* [online]. 14. December 2020 [cit. 2021-04-21]. Available at: <https://engineering.fb.com/2020/12/14/open-source/infer>.
- [41] VON PRAUN, C. and GROSS, T. Static Detection of Atomicity Violations in Object-Oriented Programs. *Journal of Object Technology*. Zürich, Switzerland: Laboratory for Software Technology, ETH Zürich. June 2004, vol. 3, no. 6, p. 103–122. DOI: 10.5381/jot.2004.3.6.a5. ISSN 16601769.
- [42] YI, K. Facebook Research. *Inferbo: Infer-based buffer overrun analyzer* [online]. 6. February 2017 [cit. 2021-04-21]. Available at: <https://research.fb.com/inferbo-infer-based-buffer-overrun-analyzer>.

Appendices

List of Appendices

A	Contents of the Attached Memory Media	79
B	Installation and User Manual	80
C	Atomer's Implementation Algorithms	84

Appendix A

Contents of the Attached Memory Media

This appendix lists the contents of the attached memory media. In particular, the attached memory media contains the following:

- `/xharmi00-thesis-2021.pdf`
 - This thesis in PDF format.
- `/thesis-latex/`
 - \LaTeX source files of the thesis.
 - It can be compiled into PDF format using `pdflatex` and `bibtex` with the `make` command (`Makefile` is included). However, the attached PDF file has been compiled in Overleaf.
- `/README.md`, `/README.html`
 - Files containing (among other things) an installation and user manual. It is available in both Markdown and HTML formats.
- `/src/`
 - The source code of the Facebook Infer framework with the Atomer plugin.
 - Atomer is “hooked” in several files of the framework. Nevertheless, the majority of the implementation is located in the sub-directory `infer/src/atomicity/`.
- `/infer-linux-x86-64-atomer-v2.0.0/`
 - Compiled Facebook Infer with the new version of Atomer (v2.0.0). It has been compiled on the Debian GNU/Linux x86-64 10.9 (Buster) operating system.
 - Executable binaries are available in the sub-directory `bin/`.
- `/examples/`
 - Some example programs with reference outputs of the analysis. These programs have been used for experimental testing of the correctness of the analysis.

Appendix B

Installation and User Manual

This appendix serves as an installation and user manual. An extended manual (together with some examples) can be found at *Atomer's Wiki* (<https://github.com/harmim/infer/wiki>) and on the attached memory media (see Appendix A).

Further, it is assumed that the *working directory* contains all the files from the attached memory media, as described in Appendix A. If not, the `/src/` directory may be cloned using `git` through SSH or HTTPS using one of the following commands:

```
git clone git@github.com:harmim/infer.git --branch atomer-v2.0.0
--single-branch --depth 1 --recurse-submodules src
git clone https://github.com/harmim/infer.git --branch atomer-v2.0.0
--single-branch --depth 1 --recurse-submodules src
```

Alternatively, it can be downloaded using `curl` as follows:

```
curl -sL
https://github.com/harmim/infer/archive/refs/tags/atomer-v2.0.0.tar.gz
| tar -xz
mv infer-atomer-v2.0.0 src
```

If the `/infer-linux-x86-64-atomer-v2.0.0/` directory containing executable binaries is not present, it can be downloaded using `curl` as follows:

```
mkdir infer-linux-x86-64-atomer-v2.0.0
curl -sL
https://github.com/harmim/infer/releases/download/atomer-v2.0.0/infer-atomer-v2.0.0.tar.xz
| tar -xJ -C infer-linux-x86-64-atomer-v2.0.0
```

Installation Manual

In the `/infer-linux-x86-64-atomer-v2.0.0/` directory, there is Facebook Infer with Atomer already compiled. Thus, it can be directly used. In the sub-directory `bin/`, there are executable binaries. Manual pages can be found in the sub-directory `share/man/`. The binaries have been compiled for the Debian GNU/Linux x86-64 10.9 (Buster) operating system. It can be installed system-wide using the following commands:

```
PATH="infer-linux-x86-64-atomer-v2.0.0/bin:$PATH"
MANPATH="infer-linux-x86-64-atomer-v2.0.0/share/man:$MANPATH"
```

Another option is to build Atomer from the source or to run it in Docker. These options are discussed in the following sections. **Running pre-compiled binaries in Docker is a recommended and quick method.**

Building Atomer from the Source

Building Atomer from the source may be very time-consuming.

At first, it is required to install Facebook Infer's dependencies and then to compile Facebook Infer with Atomer. Here are the prerequisites to be able to compile Facebook Infer with Atomer on Linux:

- `opam` \geq 2.0.0,
- `pkg-config`,
- Java (only needed for the Java analysis),
- `gcc` \geq 5.X or `clang` \geq 3.4 (only needed for the C/C++ analysis),
- `autoconf` \geq 2.63,
- `automake` \geq 1.11.1,
- `cmake` (only needed for the C/C++ analysis).

See `/src/docker/atomer-v2.0.0/Dockerfile` for inspiration on how to install the dependencies on Linux. It includes the dependencies needed to build Facebook Infer with Atomer on Debian 10 (Buster).

The installation of Facebook Infer with Atomer can be done using the following commands:

```
cd src
./build-infer.sh
sudo make install
```

The official Facebook Infer's installation manual (which also includes instructions for other operating systems) can be found at <https://github.com/harmim/infer/blob/atomer-v2.0.0/INSTALL.md>.

Facebook Infer with Atomer now should be installed system-wide. Executable binaries are located in `/src/infer/bin/`.

Running Atomer in Docker

Alternatively, Facebook Infer with Atomer can be easily built and run in Docker. The installation is much simpler and faster. At first, it is needed to install *Docker Engine*. It is available on a variety of Linux platforms, macOS, and Windows 10. Get it and follow the instructions from <https://docs.docker.com/engine/install>.

The first option is to build Facebook Infer with Atomer from the source in Docker. There is prepared a Docker image for it. *Building the image can take a few hours.* The image is based on Debian Linux. After running a container built from this image, Facebook Infer with Atomer will be compiled and installed system-wide in the container. The image can be built and the container run using the command below:

```
./src/docker/atomer-v2.0.0/run.sh # there is also run.bat for Windows
```

The second option is to run *pre-compiled binaries* of Facebook Infer with Atomer in Docker. There is prepared another Docker image for it. *Building the image can take, at most, a few minutes.* The image is based on Debian Linux as well. After running a container built from this image, binaries of Facebook Infer with Atomer will be downloaded and installed system-wide in the container. The image can be built and the container run using the command below:

```
./src/docker/atomer-v2.0.0-bin/run.sh # there is also run.bat for Windows
```

User Manual

This section assumes that Facebook Infer with Atomer is installed system-wide and executable by the command `infer`.

In general, to analyse a C/C++ program with Facebook Infer, it can be done using the following command (for a single source file):

```
infer -- gcc -c source_file.c
```

Java programs can be analysed using the following command:

```
infer -- javac source_file.java
```

Another option is to analyse the entire project with Makefile using the following:

```
infer -- make <target>
```

For advanced usage, see <https://fbinfer.com/docs/infer-workflow>. Many other build systems may be used; see <https://fbinfer.com/docs/analyzing-apps-or-projects>.

Atomer is *deactivated* by default. The analysis has to be executed twice (it has two phases). Each phase runs with a different command-line argument. The first phase derives *sets of functions that should be executed atomically* into a file `atomic-sets`. The second phase then *detects atomicity violations* according to this file. So, the analysis can be triggered using the following commands:

```
infer capture -- gcc -c source_file.c
# Without the '-only' suffix, it is triggered with other Infer's analyses.
infer analyze --atomic-sets-only # --atomic-sets
infer analyze --atomicity-violations-only # --atomicity-violations
```

To analyse Java programs, use the following commands:

```
infer capture -- javac source_file.java
infer analyze --atomic-sets-only # --atomic-sets
infer analyze --atomicity-violations-only # --atomicity-violations
```


Input Parameters

The following parameters may be used when running the `infer analyze` command:

`--atomicity-ignored-function-calls path`

- Specify a file with function names (one function name per line; considered as a regular expression if the line starts with `R` followed by whitespace, an exact match otherwise) whose *calls should be ignored* during the analysis.

`--atomicity-ignored-function-analyses path`

- Specify a file with function names (one function name per line; considered as a regular expression if the line starts with `R` followed by whitespace, an exact match otherwise) whose *analysis should be ignored* during the analysis.

`--atomicity-allowed-function-calls path`

- Specify a file with function names (one function name per line; considered as a regular expression if the line starts with `R` followed by whitespace, an exact match otherwise) whose *calls should be allowed* during the analysis. Other functions will be ignored.

`--atomicity-allowed-function-analyses path`

- Specify a file with function names (one function name per line; considered as a regular expression if the line starts with `R` followed by whitespace, an exact match otherwise) whose *analysis should be allowed* during the analysis. Other functions will be ignored.

`--atomic-sets-locked-functions-limit int`

- Specify the *maximum number of function calls* that could appear in a *critical section* in the first phase of the analysis. Critical sections with more function calls will be ignored. This is the parameter *d*. The default value is 20.

`--atomic-sets-functions-depth-limit int`

- Specify the *maximum depth in the hierarchy of function calls* to which function calls will be considered during the first phase of the analysis. This is the parameter *r*. The default value is 10.

`--atomicity-lock-level-limit int`

- Specify the *maximum expected level of ownership over the same lock object*. An over-approximation of the number of times the lock has been acquired. This is the parameter *t*. The default value is 5.

For more information about the input parameters and other helpful parameters, see `man infer` or `man infer-analyze`.

Appendix C

Atomer's Implementation Algorithms

This appendix provides an overview of some of the *implementation algorithms* of the new version of Atomer. These algorithms are either too long or not that important to be placed in the main content of the thesis.

Phase 1 — Detection of Atomic Sets

This section deals with algorithms related to the *first phase* of the analysis.

Algorithm C.1: Updating an abstract state after an *unlock*

```
Data: access paths  $L \in 2^{\Pi}$  of locks being unlocked; abstract state  $s \in \text{Domain.t}$ 
1 def apply_unlocks( $L, s$ ):
2   for  $p \in s$  do
3     for  $(A, B, \mathcal{L}) \in p.\text{callsPairs}$  do
4       if  $\mathcal{L}.\pi \in L$  then                                     /* unlock an existing lock */
5          $\mathcal{L} \leftarrow \text{Lock.unlock}(\mathcal{L})$ ;
6         if  $\neg \text{Lock.is\_locked}(\mathcal{L})$  then
7           // the lock is "completely" unlocked
8            $p.\text{finalCallsPairs} \leftarrow p.\text{finalCallsPairs} \cup \{(A, B)\}$ ;
9            $p.\text{callsPairs} \leftarrow p.\text{callsPairs} \setminus \{(A, B, \mathcal{L})\}$ ;
10           $p.\text{calls} \leftarrow \emptyset$ ;
11        end
12      end
13    end
14  end
15  return  $s$ ;
16 end
```

Algorithm C.2: Updating an abstract state after a *lock*

Data: access paths $L \in 2^\Pi$ of locks being locked; abstract state $s \in \text{Domain.t}$

```
1 def apply_locks( $L, s$ ):
2   for  $p \in s$  do
3     lockCreated  $\leftarrow$  false;
4     for  $\pi \in L$  do
5       locked  $\leftarrow$  false;
6       for  $(A, B, \mathcal{L}) \in p.\text{callsPairs}$  do
7         if  $\pi = \mathcal{L}.\pi$  then                                /* lock an existing lock */
8           locked  $\leftarrow$  true;
9            $\mathcal{L} \leftarrow \text{Lock.lock}(\mathcal{L})$ ;
10        end
11      end
12      if  $\neg \text{locked}$  then                                    /* create a new lock */
13        lockCreated  $\leftarrow$  true;
14         $p.\text{callsPairs} \leftarrow p.\text{callsPairs} \cup \{(p.\text{calls}, \emptyset, \text{Lock.create}(\pi))\}$ ;
15      end
16    end
17    if lockCreated then  $p.\text{calls} \leftarrow \emptyset$ ;
18  end
19  return  $s$ ;
20 end
```

Algorithm C.3: Updating an abstract state with a *summary of a called function*

Data: summary $\chi \in \text{Domain.Summary.t}$ of a called function; abstract state $s \in \text{Domain.t}$; maximum length $d \in \mathbb{N}$ of a critical section; the number $r \in \mathbb{N}$ of levels for the consideration of nested calls

```
1 def apply_summary( $\chi, s$ ):
2   for  $p \in s$  do
3     nestedCalls  $\leftarrow \emptyset$ ;
4     for  $(i, C) \in \chi.\text{allCalls}$  do /* get nested calls from lower-levels */
5       if  $i + 1 < r$  then /* move the calls one level lower */
6         if  $\exists p.\text{allCalls}[i + 1]$  then  $p.\text{allCalls}[i + 1] \leftarrow p.\text{allCalls}[i + 1] \cup C$ ;
7         else  $p.\text{allCalls}[i + 1] \leftarrow C$ ;
8       end
9       if  $i < r$  then nestedCalls  $\leftarrow$  nestedCalls  $\cup C$ ;
10    end
11     $p.\text{calls} \leftarrow p.\text{calls} \cup \text{nestedCalls}$ ;
12    for  $(A, B, \mathcal{L}) \in p.\text{callsPairs}$  do
13       $B \leftarrow B \cup \text{nestedCalls}$ ;
14      // rid of "large" critical sections
15      if  $|B| > d$  then  $p.\text{callsPairs} \leftarrow p.\text{callsPairs} \setminus \{(A, B, \mathcal{L})\}$ ;
16    end
17  end
18  return  $s$ ;
19 end
```

Phase 2 — Detection of Atomicity Violations

This section deals with algorithms related to the *second phase* of the analysis.

Algorithm C.4: Updating an abstract state after an *unlock*

Data: access paths $L \in 2^\Pi$ of locks being unlocked; abstract state $s \in \text{Domain.t}$

```

1 def apply_unlocks( $L, s$ ):
2   for  $p \in s$  do
3     for  $(x, y, \mathcal{L}) \in p.\text{lockedLastPairs}$  do
4       if  $\mathcal{L}.\pi \in L$  then                                     /* unlock an existing lock */
5          $\mathcal{L} \leftarrow \text{Lock.unlock}(\mathcal{L})$ ;
6         if  $\neg \text{Lock.is\_locked}(\mathcal{L})$  then
7           // the lock is "completely" unlocked
8            $p.\text{lockedLastPairs} \leftarrow p.\text{lockedLastPairs} \setminus \{(x, y, \mathcal{L})\}$ ;
9         end
10      end
11    end
12  end
13  return  $s$ ;
14 end

```

Algorithm C.5: Updating an abstract state after a *lock*

Data: access paths $L \in 2^\Pi$ of locks being locked; abstract state $s \in \text{Domain.t}$

```

1 def apply_locks( $L, s$ ):
2   for  $p \in s$  do
3     for  $\pi \in L$  do
4        $\text{locked} \leftarrow \text{false}$ ;
5       for  $(x, y, \mathcal{L}) \in p.\text{lockedLastPairs}$  do
6         if  $\pi = \mathcal{L}.\pi$  then                                     /* lock an existing lock */
7            $\text{locked} \leftarrow \text{true}$ ;
8            $\mathcal{L} \leftarrow \text{Lock.lock}(\mathcal{L})$ ;
9         end
10      end
11      if  $\neg \text{locked}$  then                                         /* create a new lock */
12         $p.\text{lockedLastPairs} \leftarrow p.\text{lockedLastPairs} \cup \{(\varepsilon, \varepsilon, \text{Lock.create}(\pi))\}$ ;
13      end
14    end
15  end
16  return  $s$ ;
17 end

```

Algorithm C.6: Updating an abstract state with a *summary of a called function*

Data: summary $\chi \in \text{Domain.Summary.t}$ of a called function; abstract state

$s \in \text{Domain.t}$

```
1 def apply_summary( $\chi$ ,  $s$ ):
2   for  $p \in s$  do
3     if  $p.\text{lockedLastPairs} = \emptyset$  then      /* outside an atomic section */
4       // take summary violations unchanged
5        $p.\text{violations} \leftarrow p.\text{violations} \cup \chi.\text{violations}$ ;
6     else                                     /* inside an atomic section */
7       // take summary violations as warnings
8       for  $(x, y, c, S) \in \chi.\text{violations}$  do
9         |  $p.\text{violations} \leftarrow p.\text{violations} \cup \{(x, y, c, \text{Warning})\}$ ;
10      end
11    end
12  end
13  return  $s$ ;
14 end
```
