

Parallelism Detection using Dynamic Instrumentation in a Virtual Machine

Christopher Edward Atkin-Granville



Master of Science
Computer Science
School of Informatics
University of Edinburgh

2013

Abstract

Contemporary computers are highly based around highly parallel architectures through chip multiprocessors, instruction-level parallelism and graphics processing units with potentially thousands of cores. Despite this, many popular programs are based around a sequential programming paradigm. This project investigates the use of the Graal compiler infrastructure in order to dynamically profile running programs within the Java Virtual Machine and determine which hot loops are good candidates for automatic parallelism transformations, possibly JIT recompilation to an OpenCL target.

We consider two main approaches to trace collection: exact approaches and probabilistic approaches.

Acknowledgements

I wish to thank my supervisors, Dr. Christophe Dubach and Dr. Björn Franke for their insightful and valuable contributions to the project.

Also in need of thanks are my parents, Sandra and Ian who have supported me throughout my entire University career.

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

(Christopher Edward Atkin-Granville)

To my grandfather, Leslie.

Table of Contents

1	Introduction	1
1.1	Golden Age	1
1.2	Cheating the System	3
1.3	Hello, Parallelism	3
1.4	Parallelising Compilers	4
1.5	Contributions	6
1.6	Outline	6
2	Background	9
2.1	Parallelist Approaches	9
3	Related Work	15
3.1	Dynamic Parallelism and Parallelism Extraction	15
3.2	Instrumentation	18
3.2.1	Bytecode Instrumentation	18
3.2.1.1	Java Agents	18
3.2.1.2	ASM	19
3.2.1.3	Javassist	20
3.2.2	Aspect-Oriented Programming	20
3.2.3	AspectJ/ABC	22
3.2.3.1	Array and Loop Pointcuts	22
3.2.4	Hybrid Models	23
3.2.4.1	DiSL	23
3.2.4.2	Turbo DiSL	24
4	The Graal Compiler Infrastructure	27
4.1	Background	27
4.2	Introduction	29
4.3	Intermediate Representations	29

4.4	Graph Transformations	31
4.4.1	The <code>.class</code> File Format - Constant Pools	32
4.5	Snippets	32
4.6	Replacements	33
4.7	Optimisations and Deoptimisations	34
4.8	Summary	37
5	Dynamic Instrumentation	39
5.1	Introduction	39
5.2	Automatic Approaches	39
5.2.1	Graal	39
5.3	Manual Approaches	43
5.4	Summary	44
6	The Runtime Library	45
6.1	Introduction	45
6.2	Trace Storage	45
6.2.1	Exact Approaches: Hash Tables and Sets	46
6.2.2	Probabilistic Approaches: Bloom Filters	47
6.3	Dependency Analysis Algorithms	48
6.3.1	Why are dependencies important?	49
6.3.2	Dependency Theory	49
6.3.3	Offline Algorithms	51
6.3.4	Online Algorithms	51
6.4	Implementation Details	53
6.4.1	Entry Point	53
6.4.2	Instrument Implementation	53
6.4.2.1	Java Auto(un)boxing	54
6.4.3	Trace Storage and Configuration	55
6.4.3.1	Exact - Hash Set	55
6.4.3.2	Inexact - Bloom Filters	56
6.5	Summary	57
7	Methodology	59
7.1	Introduction	59
7.2	Experimental Setup	59
7.3	Repeats	60
7.4	Measurement Methodology	60

7.4.1	Execution Time	60
7.4.2	Dependencies	60
7.4.3	Memory Usage	60
7.5	Parametric Benchmarks	62
7.5.1	Motivation	63
7.5.2	Problem Statement	63
7.5.3	Algorithm Design	63
7.5.4	Implementation Details	65
7.6	Test Harness Design and Implementation	65
7.7	Result Collection	66
7.8	Summary	66
8	Results	67
8.1	Introduction	67
8.2	Basic Testing	67
8.2.1	No Dependencies	67
8.2.2	All Dependent	67
8.3	Parametric	69
8.4	N-Body	70
8.5	Overhead	70
8.5.1	Execution Time	70
8.5.2	Memory Usage	70
8.6	Analysis	70
8.7	Summary	70
9	Conclusion	71
9.1	Concluding Remarks	71
9.2	Contributions	71
9.3	Unsolved Problems	71
9.4	Future Work	71
	Bibliography	73

List of Figures

4.1	The compiler pipeline in Microsoft Roslyn, a similar project to Graal. Roslyn operates at a higher level than Graal, instead manipulating abstract syntax trees. [Figure Ros]	28
4.2	Graal HIR created for a vector addition using two array literals	30
4.3	Relationships between IR levels and lowering types in Graal	32
4.4	Method substitution in Graal	34
4.5	Relationship between optimisations and deoptimisations [adapted from Schwaighofer, 2009, p. 24]	36
5.1	A high-level graph (with inlining disabled) for a simple method taking an array actual parameter, returning index 0	41
5.2	Information available at compile-time for the <code>LoadIndexedNode</code> shown in figure 5.1	41
6.1	A simple hash table for a phone book [Wikipedia, 2013a]	46
6.2	Bloom filter operation with $m = 18$ and $k = 3$ [Wikipedia, 2013b]	47
6.3	Finite state machine for the online algorithm	52
6.4	Class diagram for <code>InstrumentSupport</code>	53
6.5	Trace class hierarchy	56
8.1	Overview of the memory usage when all iterations are dependent	68

Chapter 1

Introduction

In this chapter, some background to the problem of parallelism and parallel expression is presented. Some alternative approaches to dynamic detection are covered (i.e., alternative parallel expression techniques), with a critical analysis of their strengths and merits. Lastly, a brief overview of the outline of this dissertation is presented.

Ever since the introduction of the first microprocessors in the early 1970s, there has been a trend within the microprocessor industry affectionately called Moore's Law. Although not a law in the proper scientific sense (rather, it is more of an observation of the behaviour of the industry), it does accurately describe the trend of the number of transistors which can be placed in a given area¹. The trend so far has been that this number double every 18 months. Altogether, Moore's Law successfully described the behaviour of the semiconductor industry until roughly five years ago.

1.1 Golden Age

During this time of rapid advancement, programmers had to expend very little effort in order to improve performance of their programs on newer hardware. In the best-case scenario, literally no change was required whatsoever - not even a recompilation of the program. The underlying hardware was improving, and when one combines this fact with the separation of concern between user-level applications and the underlying hardware (i.e., the abstraction layer that compilers introduce, with a special focus on ISA abstraction) meant that developers could simply urge their users to buy new hardware for a performance improvement.

¹Which is emphatically *not* the idea that processor speed doubles every 18 months - which is a common misconception. Moore's Law is also applicable to other VLSI products, such as memory.

In a slightly less-than-idea scenario, the higher transistor counts being allowed would allow semiconductor designers to add new features to the ‘base’ instruction set of their choice - for example on x86 there have been several additions over the years (MMX, 3DNow!, SSE, PAE, x86-64 etc). In these cases, programmers would simply need to recompile their programs with a compiler that would take advantage of the new extensions. Platforms supporting just-in-time (JIT) compilation such as Java, C# etc would need to replace existing virtual machines (VMs) with ones capable of using the new instruction sets.

In many ways, this time could be seen as a golden age of computer architecture. Transistors were cheap and plentiful and the promise was always there that next year transistors would be even cheaper and more plentiful. Semiconductor manufacturers started experimenting with radical new designs (not all of which were successful, for example Intel’s NetBurst which promised speeds of up to 10GHz by, amongst other techniques, involved utilising an extremely long pipeline). Consumers were confident that a new machine would be significantly faster than the machine they purchased a mere 12 months prior. Enabled by the new-found performance of processors, application developers would start to introduce many new layers of abstraction (and indirection), which would allow for safer, stabler programs to be written using high-level languages such as Ruby, Python, Perl and PHP. These extremely-high-level (EHL) languages (sometimes called scripting languages) commonly sacrificed execution speed for programmer ease of use, safely, new features and other such advantages. Indeed, this phenomenon even became widespread in lower-level languages via Java and C#, both of which introduced a virtual machine between the application and the hardware. In many cases, these virtual machines were specifically designed (at least initially) for the languages for which they were designed (in that they were not initially designed to be ‘language agnostic’), meaning they may have allowed features that are difficult to implement lower in the stack. For example, the Java Virtual Machine (JVM) includes opcodes such as `invokespecial` (which calls a special class method), `instanceof` and other such codes specifically designed for an object-oriented language². These features are enabled via high-performance processors, and would likely not exist (or certainly, not be mainstream) without these processors.

²There is currently an effort to add new instructions to the JVM designed to ease execution of languages with non-object-oriented paradigms

1.2 Cheating the System

However, these increases cannot occur indefinitely. There exists not only a fundamental lower-bound on the size of an individual transistor (as a result of quantum tunnelling), but also the extent to which contemporary techniques can provide performance improvements. For example, many common processors exploit instruction-level parallelism (ILP) by executing several instructions at the same time - pipelining. This is achieved by effectively duplicating many stages of the pipeline and the supporting infrastructure. Besides the standard issues with pipelining (data, control and structural hazards spoiling issue flow, multi-cycle instructions spoiling commit flow and the like which can be solved via trace caches, as done in the Pentium 4), there exists a larger problem. As the degree to which ILP is exploited in a processor increases, the complexity of the supporting infrastructure increases combinatorially. Hence, this is clearly not the 'silver bullet' which ILP was once thought to be. The extent to which current processors exploit ILP are not likely to increase significantly in the next several years, barring a revolutionary breakthrough in processor manufacturing, ILP detection/exploitation etc.

About a decade ago, it was a commonly held belief that the path to improving processor performance was to make a single-core processor increasingly powerful, through a combination of higher clock speeds (which manifested itself as the so-called 'Megahertz War') and architectural improvements. Although this did come true to an extent (eventually culminating in the 3.8GHz Intel Pentium 4), this period did not yield the kind of performance that was expected (see above). The main reason for this was a simple one - transistors with higher switching frequencies produce more heat. This, when combined with the fact that Moore's Law would allow higher transistor counts per unit area meant that around 2006 to 2007 manufacturers were unable to improve performance much more simply through increasing the clockspeed.

1.3 Hello, Parallelism

There existed no simple solution to this problem. For decades developers were used to having to expend little to no effort to realise potentially significant performance improvements. The solution that industry converged upon was that of parallelism - to improve performance not by increasing the performance of a single processor, but to provide many processors each of which are slightly slower when taken individually. When combined together (with a multi-threaded program), the culmination of these

processors would be more performant than a single processor could ever be.

Parallelism (and concurrency) was not a new idea. For decades parallelism had been used for the most compute-intensive problems (such as ray tracing and scientific computing). These kinds of problems are usually ‘embarrassingly parallel’ - each unit of work is totally independent from all other pieces of work. Examples of this include ray-tracing, where each ray can be simulated independently; rasterisation, where each pixel can be computed in parallel and distributed scientific problems such as SETI@Home. Indeed, concurrency has been part of developers’ standard toolkit for many years since the advent of GUIs. In Java developers commonly use helper classes such as `SwingWorker` to run compute-intensive GUI tasks in a thread independent from UI event processing in order to prevent the UI ‘hanging’ when performing long-running computations.

However the level of parallelism present in most applications is fairly superficial. Even using tools such as `SwingWorker` does not introduce a significant level of parallelism. For example, imagine a button that invokes a `SwingWorker` which executes a loop for many iterations. Although that loop is running on a different thread, that loop is *still* executing sequentially. A significant performance improvement could be realised if the developer had introduced structures and processes that allow the loop to be executed in parallel; unfortunately these transformations are non-trivial and hence are usually not performed.

Regardless of the main reason that parallelism hasn’t been introduced to any significant degree in programs (i.e., there was not a pressing need to), there are still many barriers to introducing parallelism. The main problem is likely that most developers simply do not have the required education or experience to do so. Parallelism and concurrency introduces many subtle timing errors that appear transiently. Scheduling algorithms are usually non-deterministic, which makes reasoning about them (either formally or informally) difficult. The behaviour of multi-threaded programs can change with varying number of processors.

1.4 Parallelising Compilers

The result of increasing level of parallelism of programs going from simply computing UI events in parallel to performing actual results in parallel is that it becomes substantially harder to write the parallel algorithms. Indeed, many universities do not cover parallel algorithms until final year undergraduate, postgraduate level or in

some cases not at all. This is because designing, analysing and implementing parallel algorithms is much harder to do than writing sequential ones.

In order to overcome this problem, research into the possibility of automatically parallelising compilers was started. Unfortunately, designing a compiler that is capable of parallelising arbitrary loops is arguably an intractable problem; the semantics of imperative/procedural languages do not convey enough semantic information. Additionally, the problem of data dependence became important - if an iteration is dependent on the iterations previous to it, then the iterations cannot be executed out-of-order.

Dependency detection is of critical importance to this field, as a compiler can only parallelise when iterations can be executed out-of-order. The first works in this field were in the area of static analysis. Static analysis is a compile-time process, where the compiler attempts to introspect the dependencies between iterations, parallelising if it is able to determine that there is no dependency. Although this approach would appear to be adequate, in-fact it has a significant disadvantage. Static analysis can only reason about dependencies which are present at compile-time – run-time dependencies are not possible. Although for simple programs this is likely not an issue, for more complex programs – which are more likely to benefit from parallelisation – this becomes a significant disadvantage.

As an alternative to this approach is *dynamic* analysis – analysis which occurs at run-time. This approach has the advantage of being able to detect dynamic dependencies, but at the cost of increased execution time and memory usage. This dissertation investigates this dynamic analysis through *instrumenting* programs dynamically. This instrumentation collects information about each loop and iteration, and builds data structures for suitable dependency analysis.

However, dynamic dependency analysis requires storing the address of each access per iteration, in order to ensure that no two iterations access the same location. Keeping a store of these addresses can require a significant amount of space. This dissertation investigates the use of exact approaches, in the form of sets based on hash tables, against probabilistic approaches based on bloom filters. The use of bloom filters is advantageous because unlike hash tables, they have a space complexity of $S = O(1)$, as opposed to $S = O(n)$ (where n is the number of accesses) for hash tables. However, bloom filters are probabilistic in nature, and this is taken into consideration in our analysis.

1.5 Contributions

- **Comparison between hash sets and bloom filters for dependency analysis**

The first contribution is the comparison between exact approaches (hash sets) and inexact approaches (bloom filters). The use of bloom filters was not found in a the literature review for this dissertation (section 3). We compare and contrast the relative overheads of using exact and inexact approaches, as well as the impact on the systems ability to detect parallelisable loops.

- **Implementation of parametric benchmark**

Also introduced is a parametric benchmark that allows the statistic generation of hazards. This benchmark has been used to not only measure the correctness of the framework, but also measure the relative overhead of the instrumentation. The parametric benchmark will be useful for testing the performance of similar techniques in the future.

- **Dynamic dependency analysis framework**

Lastly, a framework for dynamically proving dependencies between loop iterations is presented. Although the current framework is ‘stand-alone’, the framework has been designed to ‘drop-into’ the just-in-time compiler in the Java Runtime Environment (JRE). An analysis of possible techniques for JRE-hosted runtime detection techniques is also presented.

1.6 Outline

This dissertation is split into several chapters.

Chapter 1 outlines the problem background and context, including an overview of the current most common approaches to parallelism expression. Chapter 3 describes the previous work both the areas of dynamic parallelism detection and parallelising compilers/runtime systems. Chapter 4 is an outline of the Graal compiler infrastucture, the main tool used in the project. Chapter 5 provides an overview of the possible approaches to instrumenting Java bytecode where appropriate. Chapter 6 introduces the approaches to trace storage from both theoretical and practical perspectives. A software engineering-based overview of the approaches used is also included. Chapter 7 describes the experimental design, configuration and other parameters. Chapter 8 presents the findings and a critical analysis of the work. The last chapter, chapter 9

draws final conclusions about the work, and suggests possible areas of future work in this area.

Chapter 2

Background

2.1 Parallelist Approaches

One solution to this problem that has been advocated by many theoretical computer scientists and language designers is to switch language paradigm to functional languages such as Haskell and Erlang. Languages and systems of this class do have significant advantages - both theoretical and applied - over more conventional paradigms. Purely functional languages have first-class, referentially-transparent functions which can be easily reasoned about by both the user and the compiler. However, here are two main disadvantages to this approach. Although it is the most optimal from a theoretical perspective (and mainstream switching to a functional language would bring many benefits, not just increased parallelism), it would require rewriting existing programs in a functional language. Moreover, most developers do not have experience with functional languages - and are not, unfortunately, willing to learn. It is not a problem of technology - more one of education.

Moving more towards the practical side of things, there are two main alternatives. The first is to use language-level constructs which are built into the semantics of the language. There has been some work in this area, and they usually focus on extending existing languages with parallel semantics. These additions usually provide language-level constructs for message passing, parallel loop construction, barriers and other such low-level parallelism primitives. Although there are many such research languages (Click-5, Chapel, X10 etc), there exists no commonly used language supporting these features. There are also some hybrid languages, such as Lime which are backwards compatible with existing infrastructures [Dubach et al., 2012]. The advantages of these languages is that because they support parallelism primitives

intrinsically, reasoning and inference (both by humans and in an automated fashion) is much easier than it is for other languages. Adding support for parallelism at the language level also, depending on the level of abstraction used, implies tying a language to a particular parallel paradigm (or set of paradigms). If a language implements low-level constructs instead of higher-level constructs (similar to algorithmic skeletons except at the language level), then the user has to implement many commonly used operations manually, leading to bugs and inconsistencies between implementations. The education issue is also present with language-level constructs. Despite these disadvantages, some progress has been made in this area by Microsoft with PLINQ (Parallel Language Integrated Query), which is a declarative extension to CLI languages. PLINQ is successful because it requires little to no changes from standard LINQ - this is the ideal scenario in many cases (although note that PLINQ is not a general-purpose solution).

A more pragmatic approach to language-level constructs is a library-based approach such as POSIX Threads (Pthreads), OpenMP and OpenMPI. Pthreads is an implementation of the POSIX (*Portable Operating System Interface*) standard regarding threads and is therefore compatible with a wide variety of hardware and operating systems (including some non-POSIX compliant systems such as Windows). OpenMP is a library for C, C++ and Fortran for shared-memory programming, although it is commonly used to implement parallel loops (i.e., work sharing). OpenMPI is similar to OpenMP with the exception that it is designed for distributed-memory programming instead. The advantage of a library-based approach is that users can ‘mix-and-match’ between different libraries - for example, a user can use both OpenMP and OpenMPI within the same program. However, many of the libraries require significant low-level knowledge of parallelism, and are not particularly user friendly.

To illustrate the substantial difficulties of adding parallel semantics to existing programs with sequential semantics, let us consider the challenges programmers face when using parallel features in programming languages, such as threads. Mutual exclusion is, currently, the most widely adopted parallel programming paradigm, as it is the easiest to use (although it does not necessarily offer the greatest performance; nevertheless it can offer performance greater than that of sequential programs). Some common issues include [Fraser, 2004; Herlihy et al., 1993; Cole, 2013]:

- **Deadlock:** where two or more threads are waiting on each other to finish before continuing
- **Livelock:** similar to deadlock, except that the threads are changing state yet not achieving work

- **Priority inversion:** when a thread with high priority is pre-empted by a thread with a lower priority
- **Race conditions:** the result of two or more threads attempt to access the same resource(s) at the same time, leading to non-deterministic behaviour

There have been some attempts to mitigate these issues (e.g. Liu et al. [2011]), but they have not yet been widely adopted.

Many web-scale companies (Facebook, Google, Yahoo! etc) have large datasets that require (offline) processing and are now using parallel infrastructures such as Hadoop (MapReduce) to support this computation. Such infrastructures are similar to library-based solutions except that they also provide management solutions and other such features. Hadoop, along with its sister projects HDFS, Hive and HBase, provide an entirely framework for programming, executing, distributing and managing parallel applications. An ‘all-in-one’ solution such as Hadoop is extremely attractive simply because it provides everything one needs to set up distributed/parallel applications. The disadvantage is that, in general, such frameworks are only suited to a single kind of parallel framework. In the case of Hadoop, it can only support MapReduce-based computations. In other words, all computations that are not based around a MapReduce model are incompatible with Hadoop. Additionally, such frameworks typically require a large number of resources in order to be effective (Hadoop may actually *reduce* performance of computations if the number of processors is small), although that disadvantage is not inherent to MapReduce-style computations.

Perhaps the ‘holy grail’ of parallelism is the idea of an auto-parallelising compiler. In other words, a compiler that can infer enough information about the semantics of the program in order to apply transformations that convert a program with sequential semantics into one with parallel semantics. This approach sounds extremely attractive, as it would not require (ideally) any effort on the behalf of the programmer; despite this there are significant disadvantages. Firstly, given the scope and context of contemporary common languages (C, C++, Objective-C, C#, Java), applying compile-time transformations to add highly context-sensitive parallel semantics is likely an intractable problem. The reason for this is that, with current compiler technology, it is not possible to infer enough information regarding the semantics and syntax of the languages to reason about them. To illustrate this point, even when additional parallel semantics are added to existing languages (with sequential semantics), loops may still not be easily parallelisable because of pointer aliasing. For example, imagine a function that returns an array of values which cannot be determined at compile-time. If that array is then used to index another array (say in a vector addition), reasoning

about the parallel semantics of that vector addition are intractable at compile-time.

The last kind of parallelism is hardware-supported parallelism through constructs such as transactional memory (TM). TM solutions have the advantage of being easily programmable (to the extent that they are as easy as marking a method as *synchronized*, similar to how Java's monitors operate i.e., coarse-grained parallelism) whilst retaining the performance of fine-grained parallelism. The disadvantage of this approach is that high-performance TM architectures are reliant on hardware support (although software-based approaches do exist, they typically exhibit significantly lower performance characteristics than a similar hardware-based approach). Additionally, simply adding TM into a program does not add parallel semantics to a program with sequential semantics - constructs that allow e.g. threads to be created still need to be added. TM can rather be seen as of a way to make parallel programming easier, not a complete solution.

Transactional memory is, in effect, an optimistic memory model in that multiple threads attempt several transactions at the same time, and any conflicting writes are *rolled fine*. The final operation to be performed in a transaction is a *commit* operation - where the changes are recognised permanently in global state. However, unlike database management systems which are concerned with the ACID properties [Garcia-Molina et al., 2002, p. 14]:

- **Atomicity:** the 'all-or-nothing' execution of transactions
- **Consistency:** no transaction can move global state from a consistent state to an inconsistent state
- **Isolation:** the fact that each transaction must appear to be executed as if no other transaction is executing at the same time
- **Durability:** the condition that the effect on global state of the transaction must never be lost once the transaction is complete

However, transactional memory when used in the context of computer systems (TM can be implemented in both hardware and software, but as the semantics are the same we shall consider only transactional memory 'in the large'), we are principally concerned with only atomicity and isolation, as we assume that changes to memory do not need to be durable (memory operations are transient) [Marshall, 2005].

Another way of utilising hardware support in parallelism is to make use of low-level constructs such as compare-and-swap, test-and-set (and test-and-test-and-set) and so on. These mechanisms are very simple, but they allow programmers to implement

more complicated and sophisticated parallelism constructs using them. However, such approaches are equivalent to using a low-level language such as assembly or C to write a modern application - they are simply too low-level for programmers and are highly error prone.

It is interesting to note that languages with theoretical, rather than pragmatic, roots display excellent characteristics for automatic parallelisation. For example, functional languages display characteristics such as referential transparency. Referential transparency is the property that a function is composed entirely of *pure functions* - functions that do not modify global state. The advantage of this, along with the other properties of functional languages, is that it allows the compiler to reason about the program. Functions that display referential transparency lend themselves to easy automatic parallelisation.

Indeed, this idea of higher-level languages displaying good parallelisation characteristics extends to languages other than just functional languages. Many declarative languages (where the programmer describes *what* he/she wants to happen, in a sense without expressing control flow) lend themselves to automatic parallelisation, because they allow the compiler (or interpreter) to reason about the language.

Chapter 3

Related Work

The idea of an automatic parallelising compiler is not a particularly new one, and indeed has been the focus of much research since the dawn of structured programming with Fortran [Backus, 1979].

In this chapter, some background of both parallelising compilers and parallelism detection will be presented. The areas have a rich and full history spanning many decades (indeed, the objective of automatic parallelising compilers has been sought for many years), so this is a somewhat brief introduction; only the major results are considered.

3.1 Dynamic Parallelism and Parallelism Extraction

Programs and/or runtime systems are attractive because they require no access to the source code. In 2011, Yang et al. [2011] introduced one of the main advances on the field, *Dynamic Binary Parallelisation*. It requires no access to the source code, and instead operates only using object code. The mechanism through which it operates is by detecting hot loops – loops where the program spends a majority of execution time – and parallelises them, executing the parallel versions speculatively. This speculation is likely the cause of the inefficiencies of their approach - using 256 cores they achieved a somewhat negligible performance improvement of just 4.5 times, and the reason for this is likely hazards spoiling the speculation frequently, which requires expensive rollbacks. However, one of the advantages of the approach presented in this dissertation is that it does not require the use of speculative execution - a loop is only executed if it can be proven to contain no inter-iteration dependencies.

The main difference between Yang et al.'s work and the work outlined here is the nature of the dynamic detection. Yang et al. used dynamic trace analysis, which can only identify the hot loops within a program, and not whether the iterations within those loops have dependencies. The work presented in this dissertation *does* determine whether there are inter-loop dependencies, therefore the use of speculative execution is not required. Although we have not yet done so, the addition of hot-loop detection would be a trivial addition to our framework.

Wang et al. [2009] used a technique called *backwards slicing* [Weiser] in order to preserve essential dependence and data flow. The advantage of an approach based on slicing is that it can detect parallelism regardless of the granularity. This is in contrast to the work presented in this dissertation which can currently only detect parallelism at a loop-level. Wang et al. [2009] represents the first work in the area of dynamic parallelisation using a splicing-based approach. An approach based on slicing can determine different 'sections' of a program (such as blocks, if-then-else statements and so on) which can be executed in parallel, allowing the runtime system to execute in parallel. The algorithms in this paper are speculative - if a dependency is detected at runtime, the parallelisation for both blocks must be rolled back and instead executed in parallel. Wang et al. present a new set of algorithms for extraction of parallelism. These algorithms can:

- Determine loop-unrolling factors to expose maximum parallelism
- Identify backwards slices for hot regions
- Exploit speculation in program slicing
- Group the large number of slices

The approach to slicing is based on static analysis, with an additional speculative execution model. The speculations performed are:

- **Memory speculation:** memory dependencies on operations that occur infrequently ($<0.1\%$ according to profiling) are ignored
- **Control speculation:** control flow edges that are rarely executed ($<0.5\%$) are cut
- **Branch prediction:** highly identifiable branches are cut

Performance of the system was adequate, with a roughly 3x performance improvement with unlimited threads, and 1.8x parallelism with four threads. The major disadvantage to this approach is that the complexity of the slicing algorithms increases exponentially with the length of each slice, meaning that only relatively small slices

can be detected and hence parallelised.

An alternative approach was taken by Ketterlin and Clauss. Their technique ‘raises’ binary code into an intermediate form (similar to how Graal converts bytecode, see chapter ?? for more details), and then applies a well-known parallelising transformation.

Dong et al. present a framework that is capable of dynamic parallelisation for general-purpose graphics processing cards (GPGPU technology). GPGPU technology allows the large number of processors in GPGPUs to be used by the user for executing large programs. These processors are simpler than the general-purpose cores found in normal CPUs, and so there are some constraints placed on programs.

The framework presented takes a binary loop body and generates IR from it, similar to both Graal and Ketterlin and Clauss. Then, the IR is combined with statically-determined loop information to perform static analysis on the code. A translator is then used to transform the sequential C code into C code that can be executed in parallel using CUDA. The framework is advantageous in that it can parallelise multiple nested loops, increasing performance for many applications. However, it relies on static analysis for dependency analysis, and it does not use speculation to increase performance. However, despite these limitations, the results are acceptable, showing

Another approach taken is the one taken by Tournavitis et al. [2009]. Instead of applying more traditional algorithms to the problem, Tournavitis et al. instead use machine learning, specifically Support Vector Machines. The aims and philosophy behind the paper are similar to this project: Tournavitis et al. argue that static parallelism detection has failed because there is not enough information available at compile-time - the same reasoning that this project assumes.

The paper presents an approach to auto-parallelisation based around a combined profile of statically and dynamically generated metadata, which is then fed into a pre-trained support vector machine (SVM). This SVM provides information about how the loop should be parallelised. This approach is significant because it compromised three approaches - static analysis, dynamic analysis and machine learning. A significant performance

3.2 Instrumentation

Although in this dissertation we consider dynamic instrumentation, a significant amount of work is available in the literature describing *static instrumentation*. Static instrumentation modifies bytecode directly at load-time and is ‘dumb’ in the sense that it is not possible to do any kind of analysis on control flow graph etc. This is unlike dynamic instrumentation which has access to the full control flow graph and other such information.

Here, we consider some of the various different approaches to static instrumentation.

3.2.1 Bytecode Instrumentation

The first approach uses so-called *bytecode instrumentation* (BCI) - that is, modifying the bytecode directly, either at compile-time or runtime. This approach is advantageous in that arbitrary commands can be inserted, and is only subject to the limitations of the bytecode format. In this sense, arbitrary functionality can be inserted into `.class` files. However, it is complex (requiring advanced knowledge of the JVM and bytecode formats), as well as being difficult to use. Graal already performs a lot of the work that would be required with this kind of approach, in that it detects control flow and memory dependencies from bytecode. Such systems require a large degree of programmer effort.

3.2.1.1 Java Agents

At the heart of BCI is the idea of Java Agents [Javabeats.com, 2012]. In order to understand them, however, one must first understand some details of the Java platform.

Unlike some other languages (for example, C and C++¹), Java is a dynamically linked language. That is to say that the various different libraries (JARs) that Java programs used are linked at run-time, rather than compile-time. The advantage to this approach is that it allows distributables to be smaller in size (recall Java Network Launching Protocol, a method for launching (and therefore, distributing) Java applications over the Internet). The disadvantage to this approach is that it can lead to ‘dependency hell’, although through the combination of versioning metadata in JARs and the extreme backwards compatability mantra in Java, this is not currently a significant issue.

There are three main class loaders in Java:

¹Although note that C and C++ *also* support dynamic linking

- The system class loader loads the classes found in the `java.lang` package
- Any extensions to Java are loaded via the extension loader
- JARs found within the class path are loaded with the lowest precedence, these include the majority of user-level libraries

Java Agents manipulate class files at load-time, through the `java.lang.instrument` package. The package defines the `ClassFileTransformer` interface, which provides implementations of class transformers. An advantage of this approach is that since Java Agents are included in the core Java package, developers wishing to use *Locomotion* would not need to download and install Graal.

There are several different libraries which provide an abstraction layer for such bytecode transformations. The following sections describe various different libraries that could be used (or use themselves) for agent-based bytecode instrumentation.

3.2.1.2 ASM

Despite Java Agents providing the capability to manipulate raw Java bytecode (indeed, the bytecode is made available as a `byte` array), performing such transformations are difficult and awkward. For this reason, there exists many different libraries for manipulating Java bytecode, ObjectWeb ASM being one of them.

ASM [Bruneton et al., 2002] is a simple-to-use bytecode manipulation library, itself written in Java. It uses a high-level abstraction for the bytecode, which is advantageous because it allows developers to remain unconcerned with the specifics of control flow analysis, dependency analysis and other such concerns.

Although now common, when ASM was first developed it was considered particularly innovative because it allowed the use of the visitor pattern [Gamma et al., 2012, p. 331] for traversing bytecode. The visitor pattern ‘allows for one or more operations to be applied to a set of objects at runtime, decoupling the operations from the object structure’ [McDonald, 2008]. The advantage to this approach is that it allows a user to walk a serialized object graph *without* de-serialising it or defining large numbers of classes (for reference, an alternative to ASM for bytecode generation, BECL [Apache Foundation, 2013] contains 270 classes for representing each bytecode). Additionally, it allows users to also reconstruct a modified version of the graph (in ASM, graphs are immutable).

Several well-known existing projects use ASM already for bytecode generation, in-

cluding the Groovy programming language [Strachan and The Groovy Project, 2013]. I also have some experience of ASM through the *Compiling Techniques* coursework [Franke, 2013].

ASM was selected as a possible alternative for this project for several reasons. Firstly, its visitor pattern-based approach to bytecode generation/modification is high-level and easily understandable. It also has high performance, being a factor of 12 more performance than BECL for serialisation/deserialisation, and a factor of 35 times more performance than BECL for computing maximum stack frame sizes. ASM is also superior to BECL for performing modifications, although with a significantly lower margin of just a factor of four.

The reason for this performance improvement is likely the way ASM and BECL are designed. BECL follows a strict, classical interpretation of object-oriented design principles. Although ‘good’ software design, it is well known that object models have considerable overhead.

3.2.1.3 Javassist

Javassist [Chiba, 1998] is similar to ASM in that it is also a library for manipulating bytecode, but its method of operation is significantly different. It allows for run-time polymorphism, by dynamically switching implementation of classes at run-time.

There are also additional libraries for manipulating Java bytecode, but due to their features (or lack of), they were not considered for this project.

3.2.2 Aspect-Oriented Programming

Aspect-Oriented Programming (AOP) [Kiczales, 1996] is another dialect of object-oriented programming that aims to significantly increase separation of concern within programs, so that programs are more loosely coupled. AOP is a direct descendent from object-oriented programming as well as reflection. Reflection allows programmers to dynamically introspect classes at run-time; changing values and so on. AOP takes this to another level, by allowing so-called *advice* to be specified (essentially the additional behaviour to be added) and added to *join points*, which are arbitrary points of control flow within the program.

When combined, aspect-oriented systems add these behaviours to the program in question at compile-time through a process called *weaving*.

To illustrate this concept, consider the problem of logging method calls. In traditional systems, at each function/method definition the programmer would need to add specific logging code:

```
1 def function name(...) {  
2     if (DEBUG)  
3         println("function called at " + time());  
4  
5     // other statements  
6 }
```

Listing 3.1: Traditional use of advice in programs

This behaviour is called an *aspect* (an area of a program which may be repeated several times which is unrelated to the purpose of the program). If the behaviour is to change (e.g. for example, changing the call to `date()` to a call to `time()` instead), each method declaration must be changed manually - a time consuming and potentially error-prone task.

Instead, the use of aspects allows the programmer to remove this functionality, and combine a pointcut and advice into an *aspect*:

```
1 def aspect TraceMethods {  
2     def pointcut method-call: execution.in(*)  
3         and not(flow.in(this));  
4  
5     before method-call {  
6         println("function called at " + time());  
7     }  
8 }
```

Listing 3.2: AOP-based advice equivalent to listing 3.1

Although that from a software engineering perspective this is clearly a superior solution (decreases coupling, increases reuse, increases separation of concern), the use of aspects has not been widely adopted. There are several likely causes for this, such as:

- **Lack of education:** like the other models that we have seen in section ??, widespread adoption of new programming construct requires that the average programmer can understand the feature without in-depth education in

the model. AOP is somewhat counter to intuitive definitions of imperative or procedural languages, which hampers their adoption.

- **Lack of language support:** no widely adopted programming language comes with AOP included, or with an AOP library included in the standard library. Standard licensing issues also apply to third-party additions (e.g. GPLv2/3 differences).
- **Unclear flow control:** perhaps the single largest issue with AOP. As noted by Constantinides et al. [2004], aspects introduce effectively unconditional branches into code, mimicking the use of `goto` which Dijkstra famously considered harmful [Dijkstra, 1968].
- **Unintended consequences:** defining aspects incorrectly can lead to incorrect (global) state, e.g. renaming methods and so on. If a team of developers are unaware of each other's modifications at weave-time, there may unintended consequences and subtle (or substantial) bugs introduced.

3.2.3 AspectJ/ABC

AspectJ [Kiczales et al., 2001] is an extension to the Java language that adds aspect-oriented features. It is a project of the Eclipse Foundation (of Eclipse IDE fame). The usage of AOP within Java is a somewhat natural extension as aspects can be seen as the modularisation of behaviour (concerns) over several classes - and not to forget that AOP was originally developed as an extension to object-oriented languages.

3.2.3.1 Array and Loop Pointcuts

However, the limitations of the AspectJ join-point model are somewhat obvious for this project. To be specific, 'vanilla' AspectJ cannot define point cuts for neither array accesses or loops - a combination of which would be required for this project. In addition, the vanilla AspectJ implementation is not particularly extensible, which means that defining new point-cuts is somewhat difficult.

There is, however, an implementation of AspectJ which *is* designed to be more extensible and compatible (mostly) with the original AspectJ implementation - abc, the AspectBench Compiler for AspectJ [Allan et al., 2005].

Although abc itself does not include point-cuts for either array access or loops, there exists two projects which, if combined, could offer the required features for this

project.

LoopsAJ [Harbulot and Gurd, 2005] is an extension to abc that adds a loop join point. This is not a trivial addition - when loops are compiled, they are compiled to forms that lose loop semantics (and instead use `goto` instructions). There are several forms that a loop can take, and a significant proportion of Harbulot and Gurd's work is in the identification of loops from the bytecode.

For array access, the ArrayPT project [Chen and Chien, 2007] adds additional array access capabilities to abc. Although the included point cut does include array access, it is somewhat limited and cannot determine either the index, nor the value to be stored. ArrayPT adds these capabilities to abc. ArrayPT defines two new point cuts, `arrayset(signature)` and `arrayget(signature)`. ArrayPT relies on the `invokevirtual` bytecode in the JVM.

It is anticipated that, if these projects are combined, it would present a feasible approach to instrumentation.

3.2.4 Hybrid Models

3.2.4.1 DiSL

Recently, there has been renewed interest in Java bytecode instrumentation. Clearly, the use of aspect-oriented techniques is advantageous, but the current implementations (AspectJ/abc) are deeply flawed. In a sense, they are *static* - they rely on predefined join and point-cuts before any aspect definitions can be constructed. DiSL is considered a hybrid approach because, unlike AspectJ which relies on access to source code, it uses an agents-based approach to aspect-oriented programming.

DiSL (*Domain Specific Language for Instrumentation*) [Marek et al., 2012] is a new approach to a domain-specific language (which incidentally, implies that DiSL is declarative) for bytecode instrumentation. It does rely on the use of aspects, but it instead uses an open join-point model where any area of bytecode can be instrumented.

- Lower overheads
- Greater expressibility of aspect and join-point definition
- Greater code coverage
- Efficient synthetic local variables for data exchange between join-points

As opposed to AspectJ, which requires compile-time definition of join-points, DiSL uses an open-ended join point format which can be evaluated at weave-time. This allows arbitrary regions of bytecodes to be used as join points. *Markers* are used to specify such bytecode regions (markers are included for common join points, such as method calls and, unusually, exception handling - a novel addition to aspect-systems in Java although control-flow analysis can be used to implement user-defined markers), while *guards* allow users to further restrict selected join-points. Guards are essentially predicates which have access to only static information which can be evaluated at weave-time.

DiSL implements advice in the form of *code snippets*. Note the distinction between DiSL snippets and Graal snippets - although they are similar, DiSL snippets allow arbitrary behaviour to be inserted whilst Graal snippets are used to mainly lower complex bytecodes into simpler ones. Unlike other aspect-systems, DiSL does not support ‘around’ advice. However, this is not usually regarded as a disadvantage per-se as synthetic local variables mitigate this.

The semantics of snippets and guards is novel in DiSL. Both have complete access to local static (i.e., weave-time) reflective join-point information, meaning they can make (theoretically) unbounded numbers of references to static contexts. In addition, snippets have access to dynamic (i.e., run-time) information, including local variables and the operand stack.

Marek et al. present benchmarks of overheads with DiSL versus AspectJ, and their results are promising - a factor of three lower overheads, yet DiSL manages greater code coverage than AspectJ (the number of join-points captured is greater).

In conclusion, DiSL represents a significant advancement in aspect-systems in general. DiSL allows many semantics of dynamically-typed languages to be expressed in the (statically-typed) Java language.

3.2.4.2 Turbo DiSL

An extension to DiSL, Turbo DiSL has been proposed by Furia and Nanz [2012, p. 353-368]. Turbo DiSL is essentially an optimiser for DiSL which processes the bytecode produced by ‘vanilla’ DiSL.

There are several advantages of Turbo DiSL over DiSL. For example, instead of requiring expressions to be placed into separate classes, Turbo DiSL allows these expressions to be placed in the same class, increasing maintainability. Turbo DiSL

also performs some standard compiler optimisations on DiSL-generated code, such as pattern-based code simplification, constant propagation and conditional reduction. These are supported by a novel partial evaluation algorithm.

Turbo DiSL implements conditional reduction using partial evaluation. Many conditional control-flow statement expressions can be evaluated at weave-time – Turbo DiSL removes these dead blocks. DiSL replaces these with `pop` commands², resulting in program correctness remaining unchanged.

In addition, an approach similar to peephole-based optimisation. For example, Turbo DiSL reduces unrequired instruction such as jumping to the next instruction, or optimising the conditional reduction effects. For each `pop` instruction found, the source bytecodes are found (i.e., which bytecodes push the to-be-popped operands). If those bytecodes are side-effect free, they both (the `pop` and the source) removed.

The authors present an analysis of Turbo DiSL performance characteristics. The benchmarks selected were from the DaCapo benchmarks [Blackburn et al., 2006]. There is a considerable increase in weave-time of a factor of 7.64 above the baseline, which clearly shows the drawbacks of partial evaluation. However, Turbo DiSL outperforms DiSL by a factor of 5.18 and 13 for startup and steady-state respectively - a considerable improvement.

The authors present several uses cases where TurboDiSL is superior to DiSL (dynamic instrument configuration, tracking monitor ownership, field access analysis and execution trace profiling). However, this author speculates that, in spite of the aforementioned increase in weave-time, Turbo DiSL will completely supersede DiSL in all situations.

²<http://homepages.inf.ed.ac.uk/kwxm/JVM/pop.html>

Chapter 4

The Graal Compiler Infrastructure

Graal is a new approach to Java compiler engineering. Here, Graal is introduced and technically assessed. It's strengths and weaknesses are also evaluated.

4.1 Background

The basis of the project is the Graal compiler infrastructure [Oracle and OpenJDK, 2012]. Graal is an experimental project developed mainly at Oracle Labs (although there are some additional collaborators at AMD) under the OpenJDK programme.

The aim of the Graal project is to develop a Java compiler written in Java itself - *'a quest for the JVM to leverage its own J'*. Graal is, in essence, a Java compiler written in Java. However, this doesn't fully explain the Graal project.

Virtually all languages used commonly in industry have had their compilers go through a so-called 'bootstrap' process. This bootstrapping process involves writing the compiler for a language in the language it is intended to compile. In many cases the first version of a compiler is written in a different language - commonly used languages include the standard C and C++ due to their performance. There are many examples of this in the real-world - GCC is written in a combination of C and C++¹, LLVM/Clang is written in C++ etc. There are several advantages to this approach - in essence, this process is a kind of informal proof that the language has matured to a level that is capable of supporting a program as complicated as a compiler. In effect, bootstrapping a compiler displays shows that a language (and associated platform) has a certain level of 'maturity'; that it is now ready for large software projects (or at

¹Although the project is currently converting all C code to C++

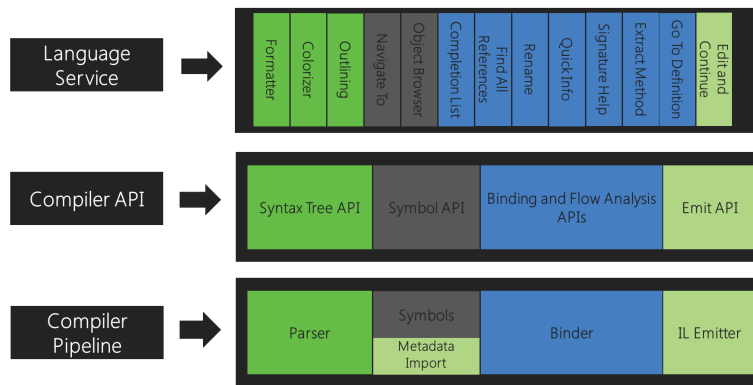


Figure 4.1: The compiler pipeline in Microsoft Roslyn, a similar project to Graal. Roslyn operates at a higher level than Graal, instead manipulating abstract syntax trees. [Figure Ros]

least not totally unprepared for them).

Graal is an attempt to bring this approach to the Java Virtual Machine. Note that there are other projects with attempts to bootstrap parts of the Java platform - for example, the Maxine VM is a Java virtual machine written in Java. However, because Java is unlike most other platforms (in that it not only requires a compiler, but also a virtual machine, class library and such), the compiler has, until the advent of the Graal project, remained written in C++.

Another feature of Graal is that it allows users (of the compiler) to interface directly with the compiler. Common compilers (GCC, ICC etc) are seen as ‘black-boxes’, where a user invokes the compiler, waits for a while, and then a resulting object file or binary is produced. Graal is a part of a new generation of compilers that expose APIs to users, which means users can change parts of the compilation process to suit their needs, ease debugging and other such advantages. With modern languages and platforms being required to target multiple different machine classes (module, desktop, laptop and server/cloud), this is a crucial advantage over more conventional languages and platforms. There are only a few examples of this new generation of compiler, but another - somewhat more mainstream example - is Microsoft’s Roslyn project [Ros] for their .NET platform. The new Windows Runtime include deep metadata integration into the platform (which is the basis for, amongst other things, the Common Object System in .NET languages²).

²Which allows inter-language types to be considered equivalent - a C int is semantically equivalent to a C++ int, a C# int, a JavaScript int and so on

4.2 Introduction

Graal is somewhat different than other compilers. As opposed to other compilers, which use a combination of parsers and lexers to produce their IRs from source code, Graal builds the IR from Java bytecode instead. This approach has several advantages for this project, the main being that we cannot assume that the source code is available to many legacy programs. Another advantage to this approach is that it would allow the detection mechanism to not only be performed upon user-provided programs, but also system-level libraries as well (for example, the Java Collections Framework). Lastly, it would allow the instrumentation of all JVM languages, not just Java. There are many languages which have been designed for the JVM, such as Scala, Groovy and Clojure, and there are also some languages which have been ported to the JVM such as Python (through Jython), Ruby (through JRuby) and PHP (through Quercus).

4.3 Intermediate Representations

Like many compilers, Graal uses several different internal representations at different stages of compilation. Each of the representations has a distinct, non-overlapping use case; despite this the graphs are somewhat similar in structure.

As is common in many compilers, Graal uses graphs for intermediate representations. These graphs combine several different kinds of graph together into a single form, such as control flow and memory dependency (data flow).

To illustrate this, consider the figure 4.2, created using the Ideal Graph Visualiser (IGV).

The format for graph visualisations is the following:

- Red edges represent control flow
- Blue edges represent memory dependence
- Black edges are defined as edges which are not control flow or memory dependence. In reality, they are mainly used for associating `FrameState` nodes where appropriate

The text inside nodes uses the following format:

```
<node-id> NodeName <additional>
```

In some cases, `<additional>` contains the value associated with the code (for subtypes

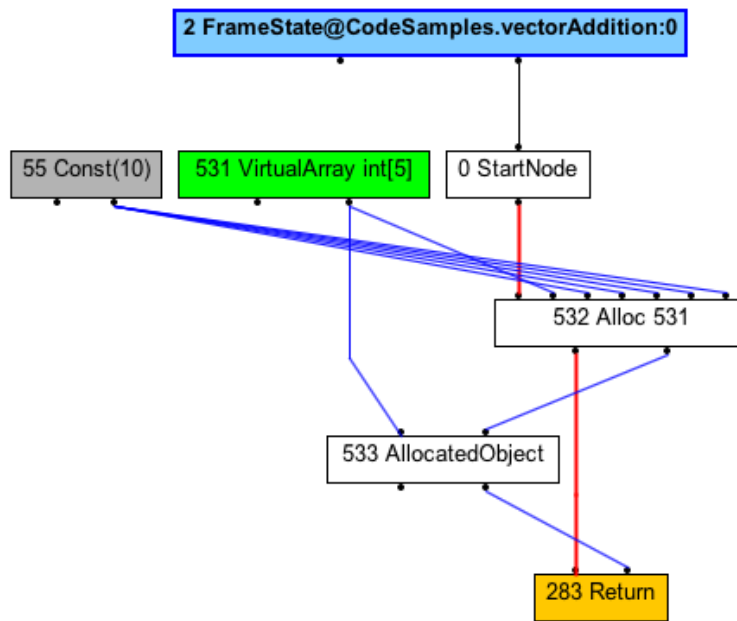


Figure 4.2: Graal HIR created for a vector addition using two array literals

of `ValueNode`). Node IDs are unrelated to the ordering within the graph. Time is represented on the y-axis of the graph, flowing down the page (so in the case of figure 4.2, node 0 (`StartNode`) is executed first, followed by node 532 and 283 in that sequence).

Figure 4.2 is a representation of a vector addition using two `int` array literals as arguments. The original source code contains a loop, which itself contains two array load operations, an addition and an array store operation. From it, we can see clearly the different kinds of node relationships in Graal’s IR. The `StartNode` is the start of the method, which is followed by an allocation (notice that the loop was been optimized away by Graal), the result of which is then returned from the method.

The allocate has two actual memory dependencies:

- A `VirtualArray` node is used to create an array. We can see that the type of this array is `int[]` with length five.
- A `ConstantNode` is referenced five times, one for each position of the array.

Without both of these dependencies, the compiler cannot create the array, and assign the correct values. Note that, because the JVM allocates default values of zero to declared-yet-unassigned variables, if the `ConstantNode` was not used, the array would consist of zeros.

4.4 Graph Transformations

One of the ways the abilities of Graal are manifested is through its capacity to apply transformations to the various intermediate representations. The main use for this is to allow users of the compiler to add custom behaviour at the various stages of the compilation process, but the mechanism extends to additional uses. For example, custom behaviour can be inserted into programs, as well as the graphs dumped for inspection in external tools.

As of the time of writing, Graal uses a hybrid approach to transforming graphs between the IRs - suites and phases. They are, in effect, essentially the same thing - they both consist of a sequence of transformations (in Graal terminology, a *phase*) which are applied to a graph in a well-defined order (although the user can change the ordering if required, as well as disable certain phases³. The result of a sequence of phases is the representation used at the next lower-level of abstraction. The three phases in Graal are high-level, mid-level, and low-level.

Todo: clear up the difference between runtime-specific lowering and target-specific lowering.

The high-level phase is used mainly for the graph-building phase. Graal uses the concept of `ResolvedJavaMethods`, which are internal 'linkages' to constant pools found within `.class` files. It is also used for runtime-specific lowering⁴, for example to handle multiple machine instruction set architectures - although the Java language is unconcerned with JVM implementation details such as the endianness of the machine's CPU, the runtime must, by definition, be aware of this fact. Examples of runtime-specific lowering are found throughout Graal, but examples include multiple implementations of the `GraalRuntime` interface - one for each ISA that Graal supports. Single Static Assignment (SSA) form is used at the high-level in order to allow for optimisations to be performed.

Mid-level phases remove the SSA form, and includes target-specific lowering. The low-level phases are analogous to the backends of more traditional compilers, and deals with low-level issues such as register allocation and code generation.

It is important to note that lowering is an iterative process, because a given lowering phase may result in further 'lowerable' nodes being added to the graph.

There are two main classes of lowering in Graal, and they are split into nodes which implement `Lowerable` and those which implement `LIRLowerable`. `Lowerable` nodes

³A common usage for this is to disable the inlining optimisation.

⁴*Lowering* is a mechanism to convert a complex bytecode into a simpler one, much like the difference between RISC and CISC architectures



Figure 4.3: Relationships between IR levels and lowering types in Graal

are transformed between the high-level and mid-level, and `LIRLowerable` nodes are transformed between the mid-level and low-level – figure 4.3 shows this diagrammatically.

4.4.1 The .class File Format - Constant Pools

In order to properly understand `ResolvedJavaMethods`, one needs to understand some parts of the `.class` file format. The source for this section is the Java Virtual Machine specification [Lindholm et al., 2013, p. 69].

`.class` files are structured containers for Java bytecode streams. However, they are not ‘plain old data structures’, as would be indicated by that description. Instead, they are laid out in such a way to increase performance for the JVM.

Unlike some other executable file formats, the JVM does not rely on the (relative) positions of the various kinds of definition permitted. In this context, constants refer to *all* immutable identifiers - and not simply to the language-level construct of constants (e.g., `public static final int THOMAS.KLAUS = 9;`). Each constant has an entry in the *constant pool* - a table containing `cp_info` structures. A `ResolvedJavaMethod` contains a link to the offset of a `cp_info` structure for a method.

4.5 Snippets

In order to engineer the lowering and optimisations mechanisms in a way that follows ‘good’⁵ manner that follows solid software engineering principles (low code reproduction, cohesion and linkage etc.), Graal uses the snippets mechanism.

Snippets are Graal graphs represented as Java source methods. They are used for lowering bytecodes with runtime-dependent semantics, such as `CHECKCAST`⁶. Snippets allow certain complex bytecodes to be replaced with simpler ones. Furthermore, snippets have the additional constraint that they must be deoptimisation (see section 4.7) free.

⁵For some definition of ‘good’

⁶<http://homepages.inf.ed.ac.uk/kwxm/JVM/checkcast.html>

Each class of snippets has an associated template, which provides the mechanism through which snippets are instantiated, and the graph modified. These are usually made available through a static nested class of the snippet class. Every time a snippet is activated, the template class creates a template for the snippet, adds arguments (via the `Arguments` class), and then modifies the graph, replacing the target bytecodes with the snippet.

Snippets are highly related to lowering, as they are used to lower nodes between the different IRs.

4.6 Replacements

In this feature, which is perhaps unique amongst compiler implementations, Graal allows the user to dynamically change implementations of methods at compile-time – in other words, Graal allows compile time method polymorphism. In a sense, replacements are similar to snippets (section 4.5), in that they both represent graph replacements. They differ in that snippets are a low-level abstraction, used for replacing bytecodes with other bytecodes, whereas replacements are for replacing method implementations.

The idea is simple: to replace the body of a method dynamically, at compile-time. This requires performing a compilation on both the replacer and replacee. The name and signatures of the methods must be identical.

This mechanism is exposed via both the `Replacements` interface, and the `@MethodSubstitution` annotation. This annotation includes various metadata, such as whether the replaced method should be substituted in all cases, and the identifier and signature of the method that will be replaced. The class which this is defined is then passed to an implementation of the `Replacements` interface, which performs a runtime-appropriate transformation.

In all cases, however, the basic idea of the transformation is the same. Both methods are compiled, and the surrounding pre-call and post-return of the replacee⁷. The nodes for the method starts and returns (or end nodes) are replaced until the transformation is complete. Figure 4.4 displays this mechanism.

⁷Of course, at this stage this compilation should have *already* been performed

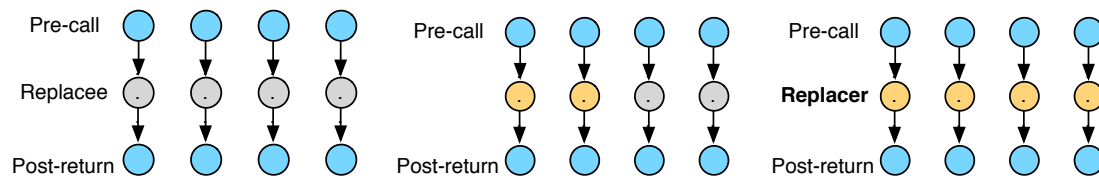


Figure 4.4: Method substitution in Graal

4.7 Optimisations and Deoptimisations

For maximum performance, Graal applies optimisations aggressively and speculatively. The potential performance advantages of speculative optimisation outweighs the possible downside of occasionally requiring deoptimisation. In addition, probability of certain blocks being executed can be computed at compile-time, which in-essence provides Graal with additional metadata regarding which blocks it may prove fruitful to optimise. If a loop is to be executed 100 times, with a low probability of exception-causing errors being thrown per iteration, it makes sense to optimise the loop for the cases where no exception will be thrown. The performance advantage of doing so will be greater than the performance decrease as a result of deoptimisation in the few(er) cases where it will be required.

Some of the optimisations that Graal applies are described here.

- **Dead and unreachable code elimination**

Dead code is code that may be executed, but has no effect on program correctness or semantics (i.e., it has no meaning [Aho et al., 2007, p. 533]). Code is said to be unreachable if there are no possible inputs which would result in execution of the block the statements exist in.

- **Type-checked method inlining**

Method inlining (also called *procedure integration*) replaces calls to procedures with copies of their bodies [Muchnick, 1997, p. 465]. However, in some cases the result of these transformations may not be type-safe, even in statically-typed languages such as Java [Glew and Palsberg, 2002]. However, transformations are possible that ensure these method inlinings are type-safe.

- **Probability of exception throwing**

Many nodes in Graal have an associated *probability* - i.e., the probability of their execution, and hence, the basic block within which they lie.

To illustrate, imagine the following loop:

```
for (int i = 1 to 100) {  
    if (i == 100)  
        doSomethingStrange();  
    else  
        doSomethingNormal();  
}
```

In 99% of cases, `doSomethingNormal()` will be executed, and in 1% of cases `doSomethingStrange()` will be executed.

This behaviour can be used to speculatively merge blocks into *extended basic blocks* - a maximal sequence of instructions beginning with a leader that contains no join nodes other than its first node [Muchnick, 1997, p. 175]. These extended basic blocks can then be optimised in the same way that 'normal' blocks can be.

If the above example were used, it could optimise for the `doSomethingNormal()` calls, and use a deoptimisation for the call to `doSomethingStrange()`.

Probabilities are computed by post-order graph traversal. When it encounters control flow nodes, Graal uses information about the split to divide probability between probability between the successors. At merge nodes, Graal sums the probability of all predecessors. Lastly, loop frequencies are propagated and each fixed node within loops have their probabilities multiplied by the loop's frequency.

Note: Graal also includes references to 'UseExceptionProbabilityForOperations', but this feature has not yet been completed.

- **Loop limit checks**

This refers to the practice of "inserting a predicate on the entry path of a loop and raising a common trap if the check or condition fails" [Peng, 2010]. Loops are split from a single block into several blocks containing a pre-, main and post-loop.

The advantages of approaches based around this technique are three-fold:

- Loop predication can be applied to outer loops without code size increment
- Checks can be eliminated from the entire iteration space
- The approaches can be applied to loops with calls

However, there are some cases where an optimisation has been applied, but it leaves programs in an inconsistent state. In this case, a deoptimisation occurs. More specifi-

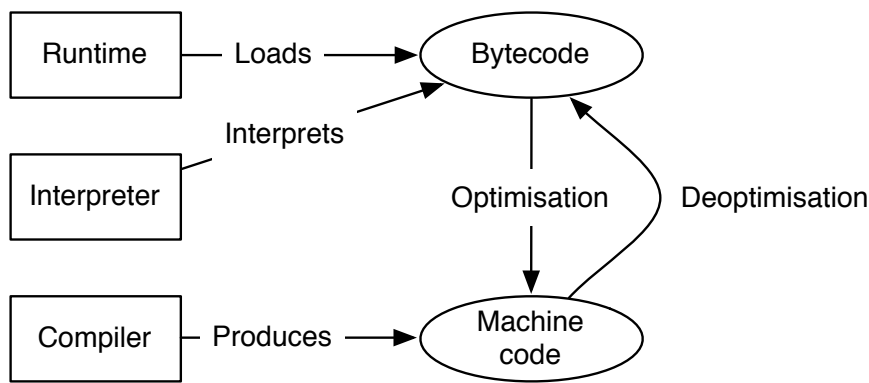


Figure 4.5: Relationship between optimisations and deoptimisations [adapted from Schwaighofer, 2009, p. 24]

cally, deoptimisation [Hölzle et al., 1992] is the process of converting an optimised stack frame into an unoptimised one. In HotSpot, an optimised frame is the result of the JIT compiler and unoptimised frames run using the interpreter. Figure 4.5 shows the interplay between the optimisations and deoptimisations. Optimisation-deoptimisation is controlled through an assumption mechanism - in order to be (or remain) optimised, a set of assumptions must remain true. If any of the assumptions becomes false, the optimisation is broken and deoptimisation occurs. When deoptimisation occurs, the runtime returns to the bytecode index associated with the optimisation, and executes using the interpreter instead. This technique is used because Graal currently relies on the deoptimisation framework in HotSpot™, which uses the approach described here.

Deoptimisations can occur for several reasons:

- **Exceptions:** certain kinds of exceptions can trigger deoptimisations as they disrupt control flow. These exceptions are `NullCheckException`, `BoundsCheckException`, `ClassCastException`, `ArrayStoreException` and `ArithmeticException`. Additionally, if an exception handler has not been compiled, a deoptimisation will occur.
- **Violated assumptions:** if any of the assumptions made during optimisation are violated (i.e., their state is modified), the optimisation is marked as invalidated.
- **Unreached code:** if there is unreachable code that could not be determined statically, this code is not optimised.

4.8 Summary

In this section, we have introduced the Graal compiler infrastructure and some of the advantages that it brings over more traditional compilers. Graal is at the forefront of compiler technology, and is part of a new generation of compilers that provide user-facing APIs. Graal also optimises speculatively and aggressively through a deoptimisation-based model; some of the optimisations and possible reasons for deoptimisations have been introduced.

Graal is distinguished from other compilers by its flexibility. It is designed to allow users to use the compiler-as-a-service, a new paradigm in compiler technology. Compilers-as-a-service promise a new level of abstraction for compiler technology, with use cases such as code analysis, online refactoring and code generation. Until now, compilers have been a black box which users have little say over the internal workings of. In the future, this will be no longer the case - compilers will adapt to suit the users and programs requirements, and Graal is the first step in this area for the Java world.

Chapter 5

Dynamic Instrumentation

5.1 Introduction

Flexible and efficient instrumentation is the foundation of this project. Here, possible approaches to this instrumentation are discussed, along with critical analyses.

There are two main approaches to the instrumentation: automatic, and manual. We consider both of them, using Graal for an automatic, dynamic approach. We conclude by showing that, currently, it is not yet possible to use Graal for automatic instrumentation – the reason for this is explained. As an alternative, manual instrumentation is used, and we explain the disadvantages and impact of this approach.

5.2 Automatic Approaches

In this context, an automatic approach to instrumentation is a technique that can be applied directly at compile-time (or just before run-time in an agent-like fashion). Such approaches require no human intervention whatsoever, and have the additional advantage of having access to additional compile-time meta-data which is not possible to gain at run-time or using manual instrumentation.

5.2.1 Graal

In the first instance, an automatic approach based on Graal was investigated. Graal was chosen in the first instance because it is the most flexible approach: it allows

the combination of compile-time evaluation, as well as run-time evaluation of the dependency algorithms.

As mentioned, Graal uses various different forms of intermediate representation, based on graphs (see section 4.3). In principle, it should be possible to modify these graphs in order to invoke the instrumentation (described in chapter 6).

Modifying graphs in Graal is made easy through a simple API. A user can define custom phases (see section 4.4). Pseudocode for this is as follows:

```

1      // create the phase
2      public class MyPhase extends Phase {
3          @Override
4          protected void run(StructuredGraph graph) {
5              // modify the graph
6              for (LoadIndexedNode node:
7                  graph.getNodes(LoadIndexedNode.class)) {
8                  graph.addAfterFixed(node,
9                      graph.add(new MyCustomNode()));
10             }
11         }
12     }
13
14     // add the phase
15     Suites s = Graal.getRequiredCapability(SuitesProvider.
16     class)
17         .createSuites();
18     s.getHighTier().appendPhase(new MyPhase());

```

Listing 5.1: Sample code for adding a phase and manipulating a graph

The call to `getHighTier()` can be replaced with equivalent methods for the other IRs (i.e., `getLowTier()` and `getMidTier()`).

As we can see in figure 5.1, which displays the high-level graph for a simple method that takes an array actual parameter and returns index 0 of that array, there are specialised nodes for array accesses: `LoadIndexedNode`. Note that there is also an equivalent for array store operations, `StoreIndexedNode`.

This is the basis for selecting the appropriate nodes to instrument. One additional

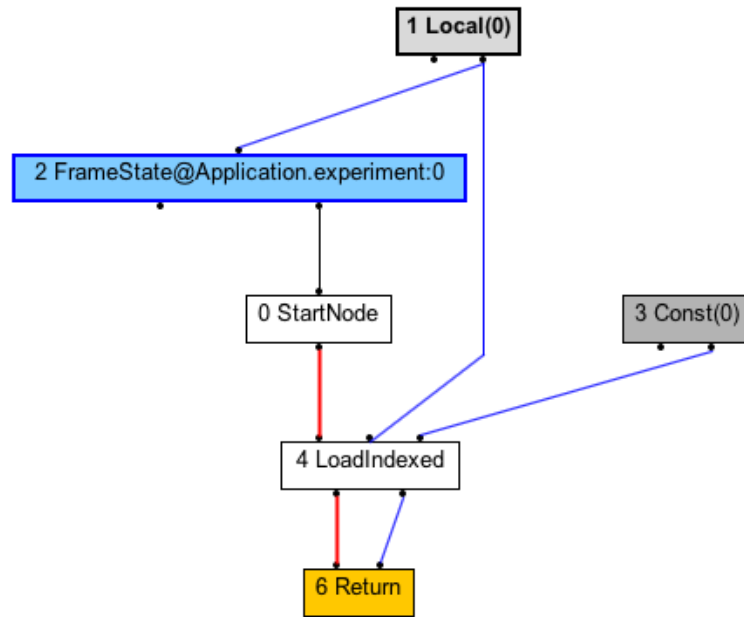


Figure 5.1: A high-level graph (with inlining disabled) for a simple method taking an array actual parameter, returning index 0

advantage of Graal is the richness of the information available at compile-time. Figure 5.2 shows all the available information, as seen in Eclipse’s debugger.

Name	Value
▶ this	MemoryOperationInstrumentationPhase (id=30)
▶ graph	StructuredGraph (id=31)
▼ node	LoadIndexedNode (id=34)
▶ array	LocalNode (id=57)
▶ elementKind	Kind (id=59)
▶ graph	StructuredGraph (id=31)
▶ id	4
▶ index	ConstantNode (id=69)
▶ modCount	0
▶ next	ReturnNode (id=72)
▶ nodeClass	NodeClass (id=74)
▶ predecessor	StartNode (id=77)
▶ stamp	ObjectStamp (id=79)
▶ typeCacheNext	null
▶ usages	NodeUsagesList (id=83)

3 | Const(0)

Figure 5.2: Information available at compile-time for the `LoadIndexedNode` shown in figure 5.1

The availability of predecessor and next nodes, as well as nodes representing indexes (in this case, a `ConstantNode`) highlights another advantage of using Graal for these transformations - static analysis is possible, which means that instrumentation can be disabled if dependencies can be proven statically.

There are several possibilities for adding instrumentation in Graal.

1. Add a custom node type, which is lowered to an `invoke` instruction after each

applicable node.

2. Replace the nodes in question with a custom node that implements the same operation in addition to the required behaviour.
3. Javassist includes a method for replacing all array accesses with method calls [Javassist Project, 2013]. It may be possible to use Graal to replace the method call targets where instrumentation is required. In cases where no instrumentation is required, the method would return the array value (or perform the store), which would then be inlined by the compiler.

In cases 1 and 2, then node would need to implement `Lowerable`, an interface in Graal that allows nodes to be lowered between IR levels.

However, there is currently a limitation in Graal which means that it is not currently possible to insert calls to (static) methods. This is because there is no *bytecode index* (or BCI) for the interpreter to return to if a deoptimisation occurs (details on Graal's deoptimisation mechanisms are available in section 4.7). Since methods calls cannot be guaranteed to be deoptimisation-free, they cannot be inserted into graphs. Any static methods that modify abstract data types (i.e., the storage formats described in section 6.2) cannot be assumed to be deoptimisation free. Modifying a data structure violates the optimisation assumptions, causing a deoptimisation. The interpreter then tries to resume from the BCI of the invoke. However, there is no BCI associated with inserted invocation, meaning that the interpreter cannot resume in the event of the inevitable deoptimisation. The end result is that insertion of arbitrary behaviour through invoke nodes to static methods is not currently possible in the Graal system.

To illustrate the concept of BCIs, consider the following example. In Java bytecode, each instruction has an associated index (this example was compiled using `javap -c` from a basic 'hello, world!' application):

Compiled from "Hello.java"

```
public class Hello {
    public Hello();
    Code:
        0: aload_0
        1: invokespecial #1    // Method java/lang/Object."<init>":()V
        4: return

    public static void main(java.lang.String[]);
    Code:
```

```

0: getstatic      #2      // Field java/lang/System.out:Ljava/io/PrintStream;
3: ldc           #3      // String Hello, world
5: invokevirtual #4      // Method java/io/PrintStream.println:(Ljava/lang/String;
8: return
}

```

We can clearly see the BCIs of the various operations: the `invokespecial #1` has a BCI of 1, and so on.

This is a limitation of the Graal platform. In the coming months, the Graal core developers are adding this required feature. Once the feature has been added, it is a simple modification (as the infrastructure has already been created) to add this into Graal. Indeed, the required infrastructure for this transformation has been created – once the support for it is available, only would need to enable the transformations.

As an alternative, we used manual instrumentation (section 5.3). Although this does have the disadvantage of requiring both source-code access as well as human effort, this will not affect the results or correctness of the evaluation as the semantics are equivalent. The results and conclusions of this report will not be affected by this shortcoming of Graal. The required framework and infrastructure has already been developed, the mechanism through which the framework is used cannot affect the results in this case.

5.3 Manual Approaches

The major alternative to automatic instrumentation is manual instrumentation, where the user manually performs the required transformations in the source code.

These transformations consist of replacing array access and store operations with method calls, whilst adding relevant metadata if required. For example, consider the following array access:

```
1 int c = a[b];
```

Listing 5.2: Standard array access in Java

Manual instrumentation refers to replacing this operation, and all such operations, with the following (or equivalent):

```
1 int c = access-array(a, b);
```

Listing 5.3: Instrumented array access

The `array-access` method (and the implied `array-store` method) performs the dependency checking algorithms using the techniques as described in section 6.3.4.

5.4 Summary

In this section, we have considered the use of Graal for automatic instrumentation. We have seen that, due to a limitation in the current version of Graal, it is not possible to add arbitrary static method invoke instructions. Regardless of this limitation, there are other possible approaches that could also (in principle) be used for automatic instrumentation.

Lastly, we have seen how manual instrumentation is possible, and the reasons why the alternative approach used (manual instrumentation) will not affect the results or conclusions of this dissertation.

Chapter 6

The Runtime Library

6.1 Introduction

In this chapter, the work performed towards implementing the runtime library, known as Locomotion, is presented. This library handles several core functions critical to the infrastructure of the application:

- Functions for collecting trace analyses
- Implementations of trace storage backends
- Algorithms for offline and online dependency analysis

The programs described within this chapter are open-source, released under The University of Edinburgh GPL license. They are available at <https://github.com/chrisatkin/locomotion>.

The fully-qualified package identifier for the runtime library is `uk.ac.ed.inf.icsa.locomotion.instrumentation`.

6.2 Trace Storage

Generating traces for large programs – the kind of programs which would benefit from hot-loop analysis – requires a large amount of storage as it scales linearly with the number of memory operations ($S = O(n)$). Although the number of storage operations conforming to the requirements in a program may be relatively small, this number is increased when the standard library is included.

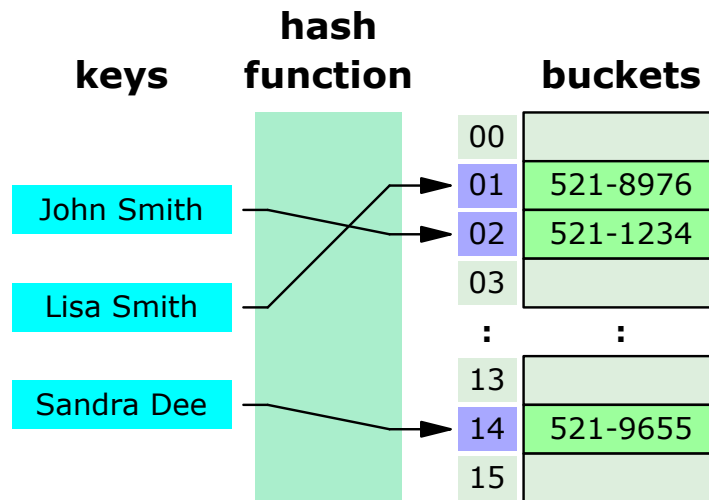


Figure 6.1: A simple hash table for a phone book [Wikipedia, 2013a]

The main problem that the storage format must be able to determine is this. Given trace T and access α , is $\alpha \in T$?

6.2.1 Exact Approaches: Hash Tables and Sets

Exact approaches provide an accurate deterministic response to this question. They are not probabilistic or statistical in nature.

A hash function maps keys to indexes, where buckets are stored. Buckets contain all items with the same key. Ideally, a lookup of $T = O(1)$ is possible (i.e., a hash function with ideal distribution). Worst-case cost is $T = O(n)$, in the case of multiple collisions or a low load factor, as the bucket will need to be traversed sequentially to find the item.

It is common to use hashing by division for the hash function:

$$f(k) = k \bmod D$$

Where k is the key and D is the size (i.e., number of positions of the hash table) [Sartaj, 1998].

The constant-time lookup assumes that the *load factor* L is bounded below some constant. The load factor is computed by $L = \frac{n}{k}$.

Figure 6.1 demonstrates a simple hash table without collisions.

6.2.2 Probabilistic Approaches: Bloom Filters

The main disadvantage to using exact approaches is that the storage required scales linearly such that $S = O(n)$. For anything but the most trivial programs, this means that it becomes infeasible to store traces for all memory operations.

One alternative is the use of Bloom Filters [Bloom, 1970]. A Bloom Filter is a randomised data structure which supports membership queries, with the possibility of false positives. In the context of parallelism detection, this means that we may conclude that a loop is not parallelisable when in reality, it is.

The operation of a Bloom Filter is simple: there exists a bit vector of size m and a number k of hash functions (which could use universal hashing [Carter and Wegman, 1979]). Upon insertion of an item i , for each hash k_n the value of $v_n = k_n(i)$ is computed, which is an integer in the range $0..m$. The corresponding index of the vector is then set to 1.

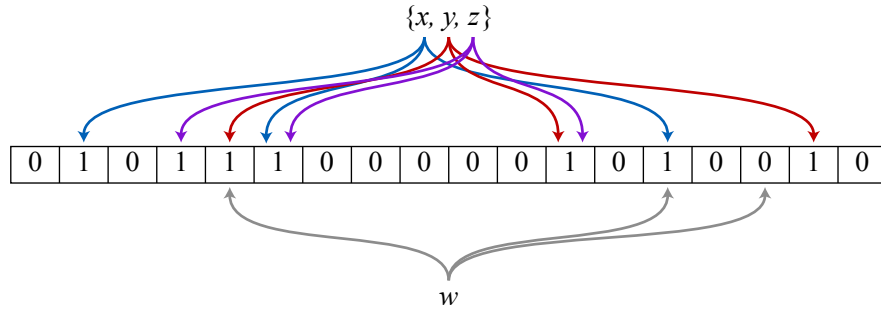


Figure 6.2: Bloom filter operation with $m = 18$ and $k = 3$ [Wikipedia, 2013b]

To test membership for an item, feed the item to each hash function. If all the corresponding indexes are 1, then the item *may* be contained within the filter - if any of them are 0, then the item is definitely not.

For a given number of entries s and a bit vector of size m , we need to use k hash functions such that:

$$k = \frac{m}{s} \ln 2 \quad (6.1)$$

The error rate is defined as:

$$\sigma = 0.5^k$$

In essence, the longer the bit vector, the more accurate the filter becomes - at the

expense of increasing space requirements. If $m = \infty$, then there are no false positives - the filter becomes accurate and deterministic (for a fixed set of k).

Swamidass and Baldi [2007] show is that the number of elements within a Bloom Filter can be estimated by:

$$E = -n \ln \frac{1 - n/n}{k}$$

It is possible to calculate the theoretical optimal length of the bit vector m from the number of expected insertions n and the *false positive rate* p . This is defined as

$$m = -n \times \log(p) \div \log(2) \times \log(2)$$

The main advantage of bloom filters is that the space complexity is constant $S = O(c)$, where c is the size of the bit vector. This will allow memory overhead to be substantially lower, whilst still retaining $T = O(1)$ lookup time complexity. As mentioned above, this is a significant (factor n) improvement over the exact approach.

However, lookup time in bloom filters is asymptotically greater than the hash-based alternative. Unlike the hashing-based approach which has, with an optimal load factor, a lookup time of $T_{\text{lookup}}^{\text{hash}} = O(1)$, a bloom filter has a lookup time of $T_{\text{lookup}}^{\text{bloom}} = O(k)$ (where k is the number of hash functions). This is unlike amongst array/vector-based approaches, which in many cases are $T = O(1)$. As the number of optimal hash functions is dependent on the number of expected entries and the bit vector size (see equation 6.1), it is predicted that execution time will increase slightly as a result, with the assumption that the hashing-based approach is operating efficiently (i.e., with an adequate load factor).

We hypothesise that, with the correct parameters, it is possible to effectively use probabilistic approaches in order to reduce both the memory and execution time overhead associated with dependency analysis.

6.3 Dependency Analysis Algorithms

Computing dependency between two operations is of critical importance to this project; an efficient algorithm is important.

6.3.1 Why are dependencies important?

Dependencies, also called *hazards* (especially in processor architecture literature) are of critical importance to automatic parallelisation as they determine whether (or not) a loop can be automatically parallelised. If there is a dependence, the statements cannot be executed in parallel without additional transformation (which may be intractable) of the program¹. However, such transformations are likely not possible for compilers to perform automatically.

6.3.2 Dependency Theory

From Allen and Kennedy [2000, p. 37], there is a data dependence from statement σ_x to statement σ_y (where $T(\sigma_x) \rightarrow T(\sigma_y)$)² if and only if:

1. both statements access the same memory location and at least one of them stores into it
2. there is a feasible run-time execution path from σ_x to σ_y

There are several kinds of inter-iteration dependency [Ibbett, 2009; Stallings, 2013, p. 526]. In this section, $R(\alpha_x)$ denotes the set of addresses/locations read by access α_x , and $W(\sigma_x)$ denotes the set of addresses read by access σ_y .

- **Write-after-write** or *output dependency*: given two statements σ_x and σ_y , there is a write-after-write dependency if a total ordering is required for σ_x and σ_y . For example:

$R3 = R1 + R2$

$R3 = R3 + R4$

If the second statement were executed before the first statement, program correctness would be violated, so speculative execution is not possible.

Formally, there is an output dependence iff $W(\sigma_x) \cap W(\sigma_y) \neq \emptyset$.

Output dependence for two statements σ_x and σ_y is denoted as $\sigma_x \delta^0 \sigma_y$.

- **Write-after-read** or *antidependency*: given two statements σ_x and σ_y , there is a write-after-read dependency if evaluation of σ_y is dependent on the evaluation of σ_x :

¹One example of where such a transformation is possible is a summation of a counter: `for int i = 0 to n: i++`. This program can be broken down into a tree of tasks and executed in a divide-and-conquer fashion, which is easily parallelisable

²For two statements x and y , $x \rightarrow y$ iff there is a viable execution path between x and y

$$R1 = R2 + R3$$

$$R2 = R4 + R5$$

In this case, if the statements are re-ordered, the wrong value of $R2$ will be used.

There is an antidependency iff $R(\sigma_x) \cap W(\sigma_y) \neq \emptyset$.

Antidependence for two statements σ_x and σ_y is denoted as $\sigma_x \delta^{-1} \sigma_y$.

- **Read-after-write** or *true dependency*: given two statements σ_x and σ_y , there is a read-after-write dependency if σ_y stores the value (or derivative value) of σ_x :

$$R1 = R2 + R3$$

$$R4 = R1 + R5$$

If the statements are reordered, either a stale value of $R1$ will be used, or $R1$ may not have been initialized.

There is a true dependency iff $W(\sigma_x) \cap R(\sigma_y) \neq \emptyset$.

True dependency for two statements σ_x and σ_y is denoted as $\sigma_x \delta \sigma_y$.

Combining these, we gain the **Bernstein Condition**. There is a dependence between σ_x and σ_y iff:

$$[R(\sigma_x) \cap W(\sigma_y)] \cup [W(\sigma_x) \cap R(\sigma_y)] \cup [W(\sigma_x) \cap W(\sigma_y)] \neq \emptyset \quad (6.2)$$

There is also a theoretical dependency that can be detected by the framework, read-after-read, but since this does not cause two iterations to be dependent, they are not considered by the framework.

For an arbitrary loop in which the loop index I runs from L to U in steps of S , the iteration number i of a specific iteration is equal to the value $\frac{I-L+S}{S}$, where I is the value of the index on that iteration.

Given a nest of n loops, the *iteration vector* i of a particular iteration of the inner-most loop is a vector of integers that contain the iteration numbers for each of the loops in order of nesting level. i is given by $i = \{i_1, i_2, \dots, i_n\}$ where i_k , $1 \leq k \leq n$, represents the iteration number for the loop at level k .

There exists a loop dependence from statement σ_1 to σ_2 in a common nest of loops iff there exists two iteration vectors i and j for the nests, such that:

1. $i < j \vee i = j$ and there is an execution path from σ_x to σ_y in the body of the loop

2. Statement σ_x accesses memory location m on iteration i (access α_x) and statement σ_y accesses location m on iteration j (access α_y)
3. Either α_x or α_y is a store

Thus, the following problem declaration can be constructed for computing dependencies:

Let i_n^l be the set of memory accesses for iteration n in loop l . For a memory access $\alpha \in i_n^l$, determine whether there exists a previous iteration such that $\alpha \in i_{n-c}^l$ for some integer c . Additionally, for an access α_i and previous access α_{i-c} , $\text{kind}(\alpha_i) \neq \text{read} \wedge \text{kind}(\alpha_{i-c}) \neq \text{read}$.

In addition, there are two kinds of loop dependence.

- **Loop-carried dependence:** σ_x can reference a common location m_c on an iteration, σ_y can reference the same location m_c
- **Loop-independent dependence:** σ_x and σ_y can both access m_c on the same iteration, but with σ_x preceding σ_y during execution

6.3.3 Offline Algorithms

An offline algorithm means that any processing is performed after the trace has been collected [Knuth, 1997, p. 525-526]. It is the simplest form of dependency analysis algorithm, but it shows poor performance. For a number of loops l with an average of i iterations each, where each iteration has an average of o operations, then $T_{\text{offline}} = O(l \cdot i \cdot o)$.

This offline algorithm has several disadvantages. Not only is its runtime particularly poor (we can achieve at least a factor l speed-up using an online algorithm), but because it requires a complete trace of accesses per iteration, it is unsuitable for detecting dependencies at run-time. Algorithm 1 outlines the offline algorithm.

6.3.4 Online Algorithms

An online algorithm for this problem is one that runs with a sequential input of values. An online algorithm is superior because it shows better performance characteristics asymptotically, $T_{\text{online}} = O(i \cdot o)$, and it runs in conjunction with the program. This allows Locomotion, in theory, to advise any JIT compilers of optimisations to perform. Figure 6.3 represents the finite state machine for the algorithm.

Algorithm 1 Offline dependency algorithm

```

1:  $d \leftarrow \emptyset$ 
2: for all loops  $l$  do
3:    $p \leftarrow \emptyset$ 
4:   for all iteration  $i \in l$  do
5:     for all access  $\alpha \in i$  do
6:       for all  $p_\alpha \in p$  do
7:         if  $\alpha \in p_\alpha$  then
8:            $d \leftarrow \alpha$ 
9:         end if
10:       $p \leftarrow \alpha$ 
11:    end for
12:  end for
13: end for
14: end for

```

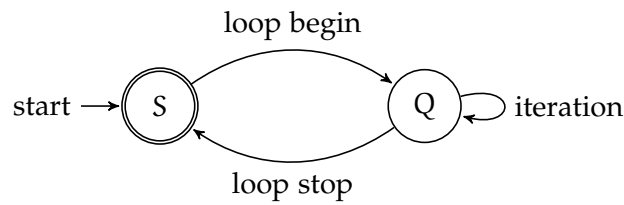


Figure 6.3: Finite state machine for the online algorithm

Algorithm 2 Online dependency algorithm

```

 $d \leftarrow \emptyset$ 
for all loops  $l$  do
   $p \leftarrow \emptyset$ 
  for all iteration  $i \in l$  do
    for all access  $\alpha \in i$  do
      for all  $p_\alpha \in p$  do
        if  $\alpha \in p_\alpha$  then
           $d \leftarrow \alpha$ 
        end if
       $p \leftarrow \alpha$ 
    end for
  end for
end for
end for

```

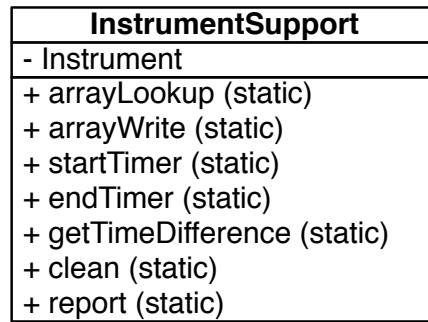


Figure 6.4: Class diagram for InstrumentSupport

6.4 Implementation Details

In this section, we consider the techniques used to implement the runtime library and online dependency algorithm from a software engineering perspective.

6.4.1 Entry Point

The main entry point to the instrumentation is the `InstrumentSupport` class. It includes an `Instrument`, which is an implementation of instrumentation. There is also a number of static methods for accessing the instrumentation framework, such as starting and ending timers (and getting time differences in nanoseconds) and flushing the instrument buffers. Figure 6.4 shows the class diagram with public methods for `InstrumentSupport`.

6.4.2 Instrument Implementation

`InstrumentSupport` is initialised at load-time with an `Instrument` (via static blocks). `Instrument` is an implementation of instrumentation, and it has a configuration. This configuration contains parameters such as whether the instrumentation should be enabled and so on.

The user-facing methods for trace collection have the following signatures:

```

1  public static <T> void
2  arrayLookup(T[] array, int index, int i, int id);
3
4  public static <T> void
5  arrayStore(T[] array, int index, T value, int i, int id);

```

Listing 6.1: Method signatures for instrumentation methods

The arguments are:

- `T[] array`: the array upon which the operation is occurring
- `int index`: the array index in question
- `int i`: the value of the iteration variable
- `T value`: the value being written to the array at the index
- `int id`: a unique loop identifier

The Java generics system was leveraged in order to reduce the amount of code required to implement the collection; it is important to recognise that the Java type system allows the trace collection mechanism to be unconcerned with the type of the array being accessed (and the type of the item being inserted if appropriate).

Although it would be possible to use the above generics-based methods for all data types, this would require the use of the object-oriented wrapper libraries (also called the reference types) for `int` \rightarrow `Integer`, `float` \rightarrow `Float` and so on. This would have the disadvantage of incurring the overhead of autoboxing [Oracle, 2010] in the library.

6.4.2.1 Java Auto(un)boxing

Despite being an object-oriented language, Java has several non-object-oriented types. These are known as *primitive values*. They are not part of the standard object-hierarchy, and no methods can be called upon them. The primitive types are `int`, `long`, `char`, `float`, `double`, `boolean` and `short`.

Each primitive type has a corresponding reference type (named as such because Java uses pass-by-reference for objects), which is an object and hence part of the class hierarchy.

Although this approach does have some advantages – primitive types require less overhead than the object-oriented value types – there are some disadvantages – namely, enable interoperability between the two.

This process is called auto boxing (in the case of primitive-to-reference type occurs) or auto unboxing (where reference-to-primitive type occurs). The advantage of this approach is that it simplifies programming and reduces ‘crufty’ code. For example, adding two `Integer`, `a` and `b` doesn’t require `Integer c = a.add(b);`, but instead the more natural `Integer c = a + b` can be used. In this case, `a` and `b` will be auto-unboxed to primitive types, the addition is performed and the result is autoboxed

back to `Integer`.

However, this automatic casting operation does occur a slight performance overhead, both in terms of memory and execution time. Although in most cases the overheads are small, for benchmarking a runtime system with somewhat significant overheads, another solution was opted for. Instead of only supporting value types, we *also* support primitive types by overloading the appropriate methods.

In order to overcome this performance problem, the instrumentation has been optimised to consider both primitive types and reference types.

6.4.3 Trace Storage and Configuration

The instrumentation includes implementations of both exact and inexact approaches to trace storage.

The `Access` class provides an abstraction for an access. An `Access` represents an array access with an associated kind, index, number (*"this iteration represents the n th access this iteration"*).

6.4.3.1 Exact - Hash Set

For the exact implementation, I chose to use the standard Java Collections Framework `HashSet` class. This is for several reasons, such as:

- **Performance:** array lookup in hash maps is $T = O(1)$, assuming an equal distribution of hash codes. For this reason, `Access` also overrides the standard `hashCode()` implementation, so that any two accesses are equivalent iff the array ID and indexes match.
- **Uniform interface:** any member of the Collections framework includes by-default various capabilities, such as being `Iterable`, meaning enhanced for-each loops can be used.
- **Standard implementation:** the built-in implementation has been tested for bugs, and is (mostly) bug-free. A custom implementation could not be subjected to this same rigorous testing.

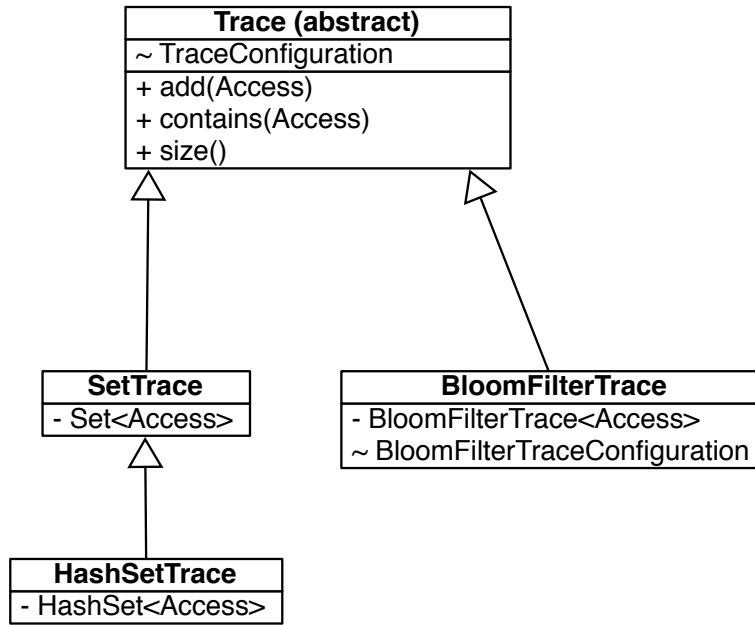


Figure 6.5: Trace class hierarchy

6.4.3.2 Inexact - Bloom Filters

Rather than creating a custom implementation of Bloom filters, I used an existing implementation in the form of Google Guava [Google, 2013].

There is a parent abstract class, `Trace` which specifies several standard methods that all interfaces must define (`add()`, `contains()` and `size()`). All implementations of traces inherit from this superclass. Figure 6.5 shows this hierarchy in detail.

A `TraceConfiguration` class is used to provide configuration details for the trace. There are no configuration variables used for hash sets, but for bloom filters this includes the initial size and the `Filter` used (an implementation detail as a result of the use of Google Guava). Dynamic configuration was preferred to constants as many experiments will need to be run, often with different variables.

When a new access α in iteration i and loop l is detected, the library checks if l is a new loop, or an existing one. If the loop is new, it instantiates a new `Loop` object, which holds the traces, as well as detecting whether an access is dependent. `Loop` stores two instantiations of the supplies `Trace` class - one for read and one for write operations. This is necessary in order to detect $\sigma_x \delta \sigma_y$, $\sigma_x \delta^0 \sigma_y$ and $\sigma_x \delta^{-1} \sigma_y$ dependencies, but not read-after-read dependencies. In order to use a single trace implementation, the semantics of the `hashCode()` and `equals()` methods, along with the `Comparable` access would need to be violated as a temporal dependence is introduced. Both bloom filters use `hashCode()` to determine an item is contained. It is not possible to use a $\sigma_x \delta \sigma_y$,

$\sigma_x \delta^0 \sigma_y$ and $\sigma_x \delta^{-1} \sigma_y$ but not read-after-read dependencies.

In order to detect dependencies, `Loop` checks all *previous* iterations, using a $T = O(1)$ check on for containment within both read and write traces. Overall, for n iterations the solution has a time complexity $T = O(n)$ and a space complexity $S = O(\alpha)$, where α is the number of accesses. Once a loop has been completed as detected by the finite state machine (see figure 6.3), the accesses are deleted in order to reduce memory location.

Dependencies are reported to the runtime via an exceptions mechanism. When a dependence is detected a special exception, `LoopDependencyException` is thrown, which is initialised with hazard metadata (which access is dependent, the kind of dependence, and the two iterations which the access occurs within). `LoopDependencyException` is a checked exception. The exception traverses up the stack until it reaches `Instrument`, which stores all exceptions of the same type in a linked list. This list is available to users of the framework, meaning full dependency information is available.

6.5 Summary

In this chapter, the theory of dependence analysis has been introduced, as well as some practical approaches to implementing them. Both kinds of dependency checking algorithm (offline and online) have been covered, as well as implementation details of online algorithms in our framework.

Lastly, a software engineering-based approach to the design of the instrumentation framework was discussed.

Chapter 7

Methodology

7.1 Introduction

In this chapter, the experimental set up and technique is outlined, as well as the specific benchmarks that were used. Rationales for each benchmark are also presented.

7.2 Experimental Setup

The hardware used for the experiments is a mid-2009 MacBook Pro with the following specifications:

- Intel Core 2 Duo Esomething
- 8GB 1067MHz DDR3
- Integrated Intel graphics
- AFS filesystem

The software configuration is as follows:

- Scientific Linux 6.3 DICE 64-bit
- Java 7 update 21

All software used was the latest version available at the time of writing.

In silico experimental design is notoriously difficult within the context of standard scientific procedure. Reproducibility of experiments is not always guaranteed, as slight and sometimes un-noticed machine intricacies and flaws can affect experimental

results. Because of this, all benchmarks were run on a separate user account (in order to reduce the number of running processes which may affect results, perhaps by consuming additional resources).

7.3 Repeats

In order to improve the results, each experiment was repeat ten times. The median of the three repeats was used for the rest of the analyses.

7.4 Measurement Methodology

7.4.1 Execution Time

Execution time was measured by taking the difference between `System.nanoTime()` before and after the experiment was run. This is superior to using other methods, such as the Unix `time` program because it computes an accurate value, instead of elapsed user-space CPU time.

7.4.2 Dependencies

When a dependency is detected, an exception is thrown to the caller, which collects dependencies in a linked list. This technique was chosen because it allows the most flexibility and separation of concern.

7.4.3 Memory Usage

The difficulties of measuring memory usage in Java programs due to the non-deterministic nature of the garbage collector are well documented in the literature [Kim and Hsu, 2000; Ogata et al., 2010]. Despite this, the Java 7 API presents several techniques [Oracle Inc, 2013] of measuring memory within the JVM:

- `freeMemory()`: the amount of free memory in the virtual machine
- `maxMemory()`: the maximum amount of memory that the virtual machine will attempt to use
- `totalMemory()`: the amount of memory currently in use by the virtual machine

In addition, there is a Java Agent for measuring memory usage of an object - the Java Agent for Memory Measurements [Ellis, 2011] (JAMM). JAMM is essentially a wrapper for the `java.lang.instrument.Instrumentation.getObjectSize()` method. There are several methods available, and the framework uses `measureDeep()` for the greatest accuracy.

`measureDeep()` crawls the object graph, calling `getObjectSize()` on each object it encounters. An `IdentityHashMap` is used to detect loops in the object graph. Unfortunately, this does affect execution time - but memory usage is recorded after execution time has been recorded. Ellis does suggest investigating the possible use of bloom filters to overcome this memory usage, but this is outside the scope of this project.

In order to take the most accurate readings possible, the JAMM library was used to measure the overhead. As an implementation detail, this led to an interesting bug in measuring the usage, which is worthy of mention.

As previously mentioned, JAMM walks the object graph, which is created by instantiation within classes - this forms an implicit directed edge between the calling context and the instantiated object. Such link is used for determining scope visibility.

In section we highlighted the use of the Google Guava library, which contains an implementation of a bloom filter. In order to decrease coupling, the Guava implementation uses so-called *Funnel*_{*T*}, which allows a class of type *T* to stream primitive values into `PrimitiveSinks`, of which the bloom filter is a member. In this way, a funnel can be thought of an interface between the filter and the type of object it contains - the `Funnel` specifies which fields should be hashed.

Initial implementations of the instrumentation used an anonymous class for the sink in a fashion similar to:

```
1 Instrument.setContainer(new BloomFilter<Access>(new Funnel<
    Access> {
2     public void funnel(Access from, PrimitiveSink to) {
3         to.putInt(from.getId())
4         .putInt(from.getIndex());
5     }
6 }));
```

This, however, results a bug in measurement of memory usage. Anonymous inner classes have an implicit edge on the object graph from their lexical context to the

context of the scope they were declared in.

```
1 public class OuterClass {  
2     public void doSomething() {  
3         // things  
4     }  
5  
6     public class InnerClass {  
7         public void test() {  
8             OuterClass.this.doSomething();  
9         }  
10    }  
11 }
```

As the funnel is stored as a field in the bloom filter, when the code to measure the size of the bloom filter was called:

```
1 public class BloomFilterTrace extends Trace {  
2     BloomFilter<Access> bloom;  
3  
4     public long getMemoryUsage() {  
5         return new MemoryMeter().measureDeep(bloom);  
6     }  
7 }
```

The memory meter would not only walk the bloom filter, but *also* the object graph of the main testing class, which contains hundreds to thousands of experiment configurations – resulting in significant memory measurements for bloom filters, despite the memory usage being very low.

7.5 Parametric Benchmarks

As we have already investigated in section 6.3, there are multiple kinds of dependency and hazards. In order to complete a thorough analysis, a new kind of benchmark has been developed in order to analyse the effect of the three kinds of dependency.

7.5.1 Motivation

In order to critically analyse the effectiveness of the bloom filter versus hash set, a controlled (i.e., non-deterministic) environment is required. In addition, a custom benchmark allows us to control the non-instrumentation workload; in order to determine an ‘upper-bound’ on the instrumentation overhead a benchmark where *only* instrumentation operations are performed.

7.5.2 Problem Statement

The problem statement is somewhat simple. Given n accesses, generate patterns of access such that there are n^δ , n^{δ^0} and $n^{\delta^{-1}}$ dependencies. These values are by definition probabilistic in nature, as the issue of array aliasing is important.

In this context, we consider an *access pattern* to be two sequential accesses, α_x and α_y , to the same index in an array. It is given that $x < y$ (i.e., $T(\alpha_x) < T(\alpha_y)$), but they may or may not be immediate. In other words, there is a viable code path between σ_x and σ_y .

7.5.3 Algorithm Design

There are five parameters to the main algorithm:

- Length of resulting arrays, l
- Percent of dependencies (expressed as a boolean in the range $0.0 \leq x \leq 1.0$), d
- Of d , what percentage of n^δ , n^{δ^0} and $n^{\delta^{-1}}$

The results of the algorithm are:

- An array, a , of random numbers within the range $0 \dots (l-1)$
- A second array, b , containing the indices of a to access at access i
- An operation array, determining which operation should be performed

To illustrate, the usage of the results are:

```

1 int[] a = getA();
2 int[] b = getB();
3 Operation[] ops = getOperations();
4
5 int temp;
```

```

6
7 for (int i = 0 to m) {
8     if op[i] == Load
9         temp = a[b[i]];
10    else if op[i] == Store
11        a[b[i]] = randomInteger(0, m);
12 }

```

This structure allows the algorithm to generate configurations of *b* and *op* to produce each kind of dependency.

The algorithm is split into three stages, one for each array. *a* is an array of random numbers.

Generating *b* can be split into two conceptual sub-arrays: the section for dependent operations, and the section for independent operations. The length of these sub-arrays depends on the parameter *d*. The independent sub-array contains simply the values $0..l_{\text{indep}}$ in incremental order, and in this dependent sub-array each value must have exactly one duplicate (in theory a value could have more than one duplicate, but this would not affect the results of the instrumentation).

To illustrate, this is an array of *b* with no dependencies:

```
0 1 2 3 4 5 6 7 8 9
```

Whereas this one has 100% dependency:

```
0 0 1 1 2 2 3 3 4 4
```

The dependent sub-array is generated using an incremental loop with a simple state machine.

After this, the arrays are merged and shuffled using the Yates-Fisher algorithm.

As a result, this stage operates in $T = O(n)$ time.

Generating *ops* Again, a state machine is used in order to determine that exactly two operations are dependent.

As the dependency kinds can be weighted with different values, a discrete probability distribution over the hazard types is used. The algorithm performs a scan over *b*, at each index performing a second scan over the array to find an equal value. If one is

found, a probability function is created, with the appropriate values being written. If a pair is not found, an operation is randomly chosen.

7.5.4 Implementation Details

The algorithm is presented as a class, with the following public interface:

```
1 public class HazardGenerator implements Generator {
2     public HazardGenerator(RandomGenerator gen);
3
4     public Generator configure(int length, int dependencies,
5 long seed, double probb_writewrite, double probb_writeread,
6 double probb_readwrite);
7
8     public Generator generate();
9
10    public int[] getA();
11
12    public int[] getB();
13
14    public AccessKind[] getAccessPattern();
15
16    public static String statistics(int[] a, int[] b,
17 AccessKind[] k);
18 }
```

The Apache Commons mathematics library is used to compute the probability distributions [Apache Commons, 2013].

Several methods return instances of `Generator` in order to fit into the rest of the instrumentation framework, as well as to allow method chaining.

7.6 Test Harness Design and Implementation

The main experimental class is `Experiments`, which contains the entry point. `Experiments` takes the parameters for the experiment, and generates an appropriate number of experiments matching the configuration. The configuration variables include the

length start, end and step size; the expected number of insertions start, end and step size for the bloom filter; the test name and so on.

There is a generic `Experiment` interface, which all experiments must implement. An experiment in this context is a specific benchmark class, not a configuration. An `Experiment` must be able to take an arbitrary number of arguments (through polymorphism to `Object[]`) and return an identifier.

An experiment class (i.e., `Class<? extends Experiment>`) is packaged with a configuration, an output format and parameters through the `Test` class, which handles instantiation of the experiment classes and so on. In addition, `Tests` can be marked to run in parallel using an `ExecutorService`, but since this may impact execution times, this feature was not used for the data collection runs.

Output is handled via the `Output` interface, for which two implementations are provided. The `Console` implementation pipes results to standard output for debugging purposes, whilst `File` uses a `PrintWriter` to pipe to a file, which can later be analysed. For performance reasons, output lines are buffered into a linked list (with $T = O(1)$ node addition time complexity) and only written to file after the experiment is over.

7.7 Result Collection

Each experiments writes it's identifier, configuration, execution time, memory usage and number of detected dependencies to a file. A separate program is then used to collate all similar results (i.e., same configuration with varying numbers of iterations and bit vector lengths if appropriate) into a single file for plotting and analysis.

7.8 Summary

In this chapter, the experimental setup for the data collection has been outlined, as well as the measurement technique for the various results. A parametric benchmark has been introduced, including a methodology, set of algorithms and some implementation details. Lastly, we covered the design of the test harness and the methodology used to collect, collate and analyse the results.

Chapter 8

Results

8.1 Introduction

In this chapter, the results of using Locomotion on the benchmarks presented in section 7.5 are presented, along with a critical analysis of the results.

8.2 Basic Testing

The first testing that was conducted was to determine the correct operation of the instrumentation and runtime library. Several tests were constructed in order to test at the extremes (for edge cases) of dependency analysis - every access is dependent, and no access is dependent.

8.2.1 No Dependencies

We first consider the case of where no dependencies are detected. We observed the following results for the exact and inexact storage mechanisms as described in section 6.2. This scenario is shown in order to prove that the instrumentation can detect lack of dependencies, and to highlight the error rate of bloom filters.

8.2.2 All Dependent

Next, the instrumentation is tested when all of the accesses are dependent. In this case, the dependencies are δ^0 dependencies, but this does not effect the results.

Figure 8.1: Overview of the memory usage when all iterations are dependent

We can see here that for both trace formats, the correct number of dependencies are detected. For the bloom filter, this is due to the ‘no false negative’ property - a false negative, in this case, is wrongly reporting that an access has not been seen before.

From figure 8.1, we can begin to see the results of using a bloom filter. When the memory usage is compared, we notice that the bloom filter is significantly lower.

Additionally, we can see from figures and ?? that the memory usage is virtually identical for both experiments, even though the number of dependencies is at a polar opposite. This is consistent with the space complexity analysis for bloom filters, presented in section ?. The memory usage does not vary with the number of accesses - again, this is to be expected as the bit vector is of a fixed length.

From these memory usage results we can also begin to understand why memory usage measurement is so difficult in the JVM – there are several anomalies in both datasets. This is likely due to the unpredictable/non-deterministic nature of the garbage collection algorithm, as well as the OS memory allocation mechanism. However, when we only take into account the non-anomalous results and see the overall trends, we see a clear pattern already – bloom filters offer a significant memory usage improvement over sets.

Next, we consider the execution time overhead of using exact versus inexact approaches. Details of the methodology used to measure the time difference can be found in section 7.4.1.

We can see from figure ?? that bloom filters do, in some cases, have a higher execution

time than hash sets. This is likely as a result of the aforementioned (section 6.4.3.2) $T_{\text{lookup}}^{\text{bloom}} = O(k)$ lookup time.

However, for a thorough analysis we must consider the implementation of the bloom filter used. The implementation used is from Google Guava, which uses an innovative technique presented by Azar and Erlebach [2006] who show that instead of $T_{\text{lookup}}^{\text{bloom}} = O(k)$, it is possible to instead use a 2-family of universal hash functions instead, without an asymptotic loss in false positive probability. The time complexity of universal hash functions for hashing n -bit strings to m -bit strings has been shown to be $T_{\text{uh}} = \Omega(mn)$ [Mansour et al., 1990]¹.

Although this is a performance improvement (from $T_{\text{lookup}}^{\text{bloom}} = O(k)$ to $T_{\text{lookup}}^{\text{bloom}} = \Omega(mn)$), there are still constants to be considered. An optimal hash map implementation (i.e., one with an adequate load factor and capacity) requires a single hash, and the hash is in-fact $T = O(1)$, whereas the bloom filter requires two $T = \Omega(mn)$ hashes. This is the likely cause of the slightly slower performance of the bloom-filter implementation.

The result for bloom filter with a bit vector length of 100 using 1000 accesses is considered to be an anomaly. Although the other results for a length 100 bit vector do show increases on time, there appears to be no correlation between number of accesses and execution time when the bit vector is 100 in length. Regardless, as we shall see in later sections, a length 100 bit vector results in a large dependency detection rate (ϵ_δ), and are therefore unsuitable for dependency analysis algorithms.

The cause of the unexpected (i.e., non-linear) execution time for the hash set is likely an sub-optimal load factor. Later, we will investigate the effects of modifying this, as well as the expected number of insertions.

In these experiments, we used a fixed ffp of 0.03. In later experiments, this value will be modified in order to determine the optimal configuration for the bloom filter.

8.3 Parametric

The parametric tests are described thoroughly in chapter ??.

The first experiments that we ran with the parametric benchmark were constant numbers of dependencies. The number of dependencies was constant at 900.

¹Note that it is possible to achieve a $T_p = \Omega(\log n)$ parallel implementation of universal hashing using CREW PRAM architecture. However, this is not likely to have a significant performance improvement.

As we can see in figure ??, the number of dependencies is a constant 900 for all accesses lengths. Naturally, this is expected.

8.4 N-Body

8.5 Overhead

8.5.1 Execution Time

8.5.2 Memory Usage

8.6 Analysis

8.7 Summary

Chapter 9

Conclusion

9.1 Concluding Remarks

9.2 Contributions

9.3 Unsolved Problems

9.4 Future Work

The work that has been presented in this dissertation is a step forwards in dynamic parallelism detection. The framework correctly identifies data-parallel loops, and we have investigated whether the use of bloom filters affects the detection rate. However, there are still significant areas of future work that are possible.

The framework presented is only currently capable of detecting parallelism at the level of loops (i.e., loop-level parallelism).

Additionally, although the framework *does* correctly detect parallelism at run-time, this information is current under-exploited by the runtime. It is possible, at least theoretically, that the framework could be combined with fellow student Ranjeet Singh's Java-to-OpenCL compiler in a JIT setting to produce a runtime system that dynamically detects parallel loops, recompiles then and executes using OpenCL (either on a CPU or GPU). However, there is an additional downside to this: the advantage (i.e., performance increase) of dynamic recompilation, moving execution to the CPU and then gathering the result must be greater than that of not doing so.

Although thread creation on CPUs is expensive in Java [Oracle Corporation, 2004], thread creation on GPUs is low [Mueller, 2009] (indeed, GPUs need 1000s of threads to operate efficiently in CUDA [Nvidia, 2011]). In order to properly perform recompilation (or not), a cost model would need to be developed.

Bibliography

Roslyn Project. URL <https://www.microsoft.com/en-us/download/details.aspx?id=27744>.

Arnold Schwaighofer. *Tail Call Optimization in the Java HotSpot VM Diplom-Ingenieur*. Masterstudium informatik thesis, Johannes Kepler University Linz, 2009. URL <http://www.ssw.uni-linz.ac.at/Research/Papers/Schwaighofer09Master/schwaighofer09master.pdf>.

Wikipedia. Hash table diagram, 2013a. URL https://en.wikipedia.org/wiki/File:Hash_table_3_1_1_0_1_0_0_SP.svg.

Wikipedia. Bloom filter diagram, 2013b. URL https://en.wikipedia.org/wiki/File:Bloom_filter.svg.

Christophe Dubach, Perry Cheng, Rodric Rabbah, David F. Bacon, and Stephen J. Fink. Compiling a high-level language for GPUs via language support for architectures and compilers. *ACM SIGPLAN Notices*, 47(6):1–12–12, August 2012. ISSN 0362-1340. doi: 10.1145/2345156.2254066. URL <http://dl.acm.org/citation.cfm?id=2345156.2254066>.

Keir Fraser. Technical Report. Technical Report 579, University of Cambridge, Cambridge, 2004. URL <http://www.cl.cam.ac.uk/techreports/UCAM-CL-TR-579.pdf>.

M Herlihy, J Eliot, and B Moss. Transactional Memory: Architectural Support For Lock-free Data Structures. In *Proceedings of the 20th Annual International Symposium on Computer Architecture*, volume 21, pages 289–300. IEEE Comput. Soc. Press, 1993. ISBN 0-8186-3810-9. doi: 10.1109/ISCA.1993.698569. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=698569>.

Murray Cole. Parallel Programming Languages and Systems, 2013. URL <http://www.inf.ed.ac.uk/teaching/courses/ppls/pplsslides.pdf>.

Tongping Liu, Charlie Curtsinger, and Emery D. Berger. Dthreads. In *Proceed-*

- ings of the Twenty-Third ACM Symposium on Operating Systems Principles - SOSP '11*, page 327, New York, New York, USA, October 2011. ACM Press. ISBN 9781450309776. doi: 10.1145/2043556.2043587. URL <http://dl.acm.org/citation.cfm?id=2043556.2043587>.
- Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. *Database Systems: The Complete Book*. Prentice Hall, first edition, 2002. ISBN 0-13-031995-3.
- Casey Marshall. Software Transactional Memory. Technical report, University of California, Santa Cruz, Santa Cruz, 2005. URL <http://www.cs.uic.edu/~ajayk/STM.pdf>.
- John Backus. The History of FORTRAN I, II and III. *IEEE Annals of the History of Computing*, 1(1):21–37, January 1979. ISSN 1058-6180. doi: 10.1109/MAHC.1979.10013. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4392880>.
- Jing Yang, Kevin Skadron, Mary Lou Soffa, and Kamin Whitehouse. Feasibility of Dynamic Binary Parallelization. *Usenix*, 2011. URL http://www.cs.virginia.edu/~skadron/Papers/yang_hotpar11.pdf.
- Cheng Wang, Youfeng Wu, Edson Borin, Shiliang Hu, Wei Liu, Dave Sager, Tin-fook Ngai, and Jesse Fang. Dynamic Parallelization of Single-Threaded Binary Programs using Speculative Slicing. 2009.
- M. Weiser. Reconstructing sequential behavior from parallel behavior projections. *Information processing letters*, 17(3):129–135. ISSN 0020-0190. URL <http://cat.inist.fr/?aModele=afficheN&cpsidt=9606727>.
- Alain Ketterlin and Philippe Clauss. Transparent Parallelization of Binary Code.
- Guoxing Dong, Kai Chen, Erzhou Zhu, Yichao Zhang, Zhengwei Qi, and Haibing Guan. A Translation Framework for Virtual Execution Environment on CPU / GPU Architecture. doi: 10.1109/PAAP.2010.53.
- Georgios Tournavitis, Zheng Wang, Björn Franke, and Michael F.P. O’Boyle. Towards a holistic approach to auto-parallelization. *ACM SIGPLAN Notices*, 44(6):177, May 2009. ISSN 03621340. doi: 10.1145/1543135.1542496. URL <http://dl.acm.org/citation.cfm?id=1543135.1542496>.
- Javabeats.com. Introduction to Java Agents, 2012. URL <http://www.javabeat.net/2012/06/introduction-to-java-agents/>.
- Eric Bruneton, Romain Lenglet, and Thierry Coupaye. ASM: a code manipulation tool to implement adaptable systems. *Adaptable and extensible component systems*, 30, 2002.

doi: 10.1.1.117.5769. URL <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.117.5769>.

Erich Gamma, Richard Helm, Ralph Johnson, and John Vlissides. *Table of Contents. Carcinogenesis*, 33(8):NP–NP, August 2012. ISSN 0143-3334. doi: 10.1093/carcin/bgs084. URL <http://www.carcin.oxfordjournals.org/cgi/doi/10.1093/carcin/bgs084>.

Jason McDonald. *Design Patterns*. Technical report, DZone, 2008. URL http://cs.franklin.edu/~whittakt/COMP311/rc008-designpatterns_online.pdf.

Apache Foundation. *Apache Bytecode Engineering Library*, 2013. URL <http://commons.apache.org/proper/commons-bcel/>.

James Strachan and The Groovy Project. *org.codehaus.groovy.classgen.asm*, 2013. URL <http://groovy.codehaus.org/gapi/org/codehaus/groovy/classgen/asm/package-summary.html>.

Björn Franke. *Compiling Techniques coursework*, 2013. URL http://www.inf.ed.ac.uk/teaching/courses/ct/Coursework/Coursework_2013.pdf.

Shigeru Chiba. *Javassist - A Reflection-Based Programming Wizard for Java*. In *Proceedings of the OOPSLA workshop on Reflective Programming in C and*, 1998. URL <http://www.csg.ci.i.u-tokyo.ac.jp/~chiba/oopsla98/proc/chiba.pdf>.

Gregor Kiczales. *Aspect-oriented programming*. *ACM Computing Surveys*, 28(4es): 154–es, December 1996. ISSN 03600300. doi: 10.1145/242224.242420. URL <http://portal.acm.org/citation.cfm?doid=242224.242420>.

C Constantinides, T Skotiniotis, and M Stoerzer. *AOP considered harmful*. Technical report, Concordia University, 2004. URL <http://pp.info.uni-karlsruhe.de/uploads/publikationen/constantinides04eiwas.pdf>.

Edsger W Dijkstra. *Letters to the editor: go to statement considered harmful*. *Communications of the ACM*, 11(3):147–148, March 1968. ISSN 00010782. doi: 10.1145/362929.362947. URL <http://portal.acm.org/citation.cfm?doid=362929.362947>.

Gregor Kiczales, Erik Hilsdale, Jim Hugunin, Mik Kersten, Jeffrey Palm, and William G Griswold. *An Overview of AspectJ*. *Main*, 2072:327–353, 2001. ISSN 03029743. doi: 10.1007/3-540-45337-7_18. URL <http://www.cs.ubc.ca/~kdvolder/binaries/aspectj-overview.pdf>.

Chris Allan, Pavel Avgustinov, Aske Simon Christensen, Laurie J Hendren, Sascha Kuzins, Jennifer Lhoták, Ondrej Lhoták, Oege De Moor, Damien Sereni, Ganesh

- Sittampalam, and Julian Tibble. abc: The AspectBench Compiler for AspectJ. In *Generative Programming and Component Engineering 4th International Conference GPCE*, volume 3676, pages 10–16, 2005. ISBN 3540291385. doi: 10.1007/11561347\2.
- Bruno Harbulot and John R Gurd. A join point for loops in AspectJ. In *Computer*, pages 63–74, 2005. ISBN 159593300X. doi: 10.1145/1119655.1119666.
- Kung Chen and Chin-hung Chien. Extending the Field Access Pointcuts of AspectJ to Arrays. *Journal of Software Engineering Studies*, 2(2):2–11, 2007. URL http://www.geocities.ws/m8809301/pub/JSESv2n2_KungChen_970214.pdf.
- Lukáš Marek, Alex Villazón, Yudi Zheng, Danilo Ansaloni, Walter Binder, and Zhengwei Qi. DiSL. In *Proceedings of the 11th annual international conference on Aspect-oriented Software Development - AOSD '12*, page 239, New York, New York, USA, March 2012. ACM Press. ISBN 9781450310925. doi: 10.1145/2162049.2162077. URL <http://dl.acm.org/citation.cfm?id=2162049.2162077>.
- Carlo A. Furia and Sebastian Nanz, editors. *Objects, Models, Components, Patterns*, volume 7304 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012. ISBN 978-3-642-30560-3. doi: 10.1007/978-3-642-30561-0. URL <http://www.springerlink.com/index/10.1007/978-3-642-30561-0>.
- S M Blackburn, R Garner, C Hoffman, A M Khan, K S McKinley, R Bentzur, A Diwan, D Feinberg, D Frampton, S Z Guyer, M Hirzel, A Hosking, M Jump, H Lee, J E B Moss, A Phansalkar, D Stefanović, T VanDrunen, D Von Dincklage, and B Wiedermann. The DaCapo benchmarks: Java benchmarking development and analysis. *ACM Sigplan Notices*, 41:169–190, 2006. ISSN 03621340. doi: 10.1145/1167515.1167488.
- Oracle and OpenJDK. Graal Project, 2012. URL <http://openjdk.java.net/projects/graal/>.
- Tim Lindholm, Alex Buckley, Gilad Bracha, and Frank Yellin. The Java Virtual Machine Specification Java SE 7 Edition. Technical report, Oracle, Inc, Redwood City, 2013. URL <http://docs.oracle.com/javase/specs/jvms/se7/jvms7.pdf>.
- Alfred V. Aho, Monica S. Lam, R Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques and Tools*. Pearson, second edition, 2007. ISBN 0-321-49169-6.
- Steven S. Muchnick. *Advanced Compiler Design and Implementation*. Morgan Kaufmann, San Francisco, first edition, 1997. ISBN 1-55860-320-4.
- Neal Glew and Jen Palsberg. Type-Safe Method Inlining. In Boris Magnusson, editor,

- ECOOP 2002 *Object-Oriented Programming*, volume 2374 of *Lecture Notes in Computer Science*, pages 525–544. Springer Berlin Heidelberg, Berlin, Heidelberg, May 2002. ISBN 978-3-540-43759-8. doi: 10.1007/3-540-47993-7. URL <http://link.springer.com/10.1007/3-540-47993-7>.
- Chang Peng. Loop Prediction, 2010. URL <https://wikis.oracle.com/display/HotSpotInternals/LoopPredication>.
- Urs Hölzle, Craig Chambers, and David Ungar. Debugging optimized code with dynamic deoptimization. *ACM Sigplan Notices*, 27:32–43, 1992. ISSN 03621340. doi: 10.1145/143103.143114.
- Javassist Project. Javassist documentation, 2013.
- Sahni Sartaj. *Data Structures, Algorithms and Applications in C++*. McGraw-Hill, first edition, 1998. ISBN 0-07-109219-6.
- Burton H. Bloom. Space/time trade-offs in hash coding with allowable errors. *Communications of the ACM*, 13(7):422–426, July 1970. ISSN 00010782. doi: 10.1145/362686.362692. URL <http://dl.acm.org/citation.cfm?id=362686.362692>.
- J. Lawrence Carter and Mark N Wegman. Universal classes of hash functions. *Journal of Computer and System Sciences*, 18(2):143–154, April 1979. ISSN 00220000. doi: 10.1016/0022-0000(79)90044-8. URL <http://linkinghub.elsevier.com/retrieve/pii/0022000079900448>.
- S Joshua Swamidass and Pierre Baldi. Mathematical correction for fingerprint similarity measures to improve chemical retrieval. *Journal of chemical information and modeling*, 47(3):952–64, 2007. ISSN 1549-9596. doi: 10.1021/ci600526a. URL <http://www.ncbi.nlm.nih.gov/pubmed/17444629>.
- Randy Allen and Ken Kennedy. *Optimising Compilers for Modern Architectures: A Dependence-based Approach*. Morgan Kaufmann, San Francisco, first edition, 2000. ISBN 1-55860-286-0.
- Ronald Ibbett. High Performance Computer Architectures, 2009. URL <http://homepages.inf.ed.ac.uk/cgi/rni/comp-arch.pl?Paru/depend.html,Paru/depend-f.html,Paru/menu.html>.
- William Stallings. *Computer Organisation and Architecture: Design for Performance*. Pearson, Harlow, ninth edition, 2013. ISBN 0-273-76919-7.
- Donald E. Knuth. *The Art of Computer Programming Volume 2: Seminumerical Algorithms*. Addison Wesley, third edition, 1997. ISBN 0-201-89683-2.

- Oracle. Autoboxing, 2010. URL <http://docs.oracle.com/javase/1.5.0/docs/guide/language/autoboxing.html>.
- Inc Google. Class BloomFilter;T¿, 2013. URL <http://docs.guava-libraries.googlecode.com/git-history/release/javadoc/com/google/common/hash/BloomFilter.html>.
- Jin-Soo Kim and Yarsun Hsu. Memory system behavior of Java programs. *ACM SIGMETRICS Performance Evaluation Review*, 28(1):264–274, June 2000. ISSN 01635999. doi: 10.1145/345063.339422. URL <http://dl.acm.org/citation.cfm?id=345063.339422>.
- Kazunori Ogata, Dai Mikurube, Kiyokuni Kawachiya, Scott Trent, and Tamiya Onodera. A study of Java’s non-Java memory. *ACM SIGPLAN Notices*, 45(10):191, October 2010. ISSN 03621340. doi: 10.1145/1932682.1869477. URL <http://dl.acm.org/citation.cfm?id=1932682.1869477>.
- Oracle Inc. Runtime (Java Platform SE 7), 2013. URL <http://docs.oracle.com/javase/7/docs/api/java/lang/Runtime.html>.
- Jonathan Ellis. Java Agent for Memory Measurements, 2011. URL <https://github.com/jbellis/jamm/>.
- Apache Commons. EnumeratedDistribution;T¿ documentation, 2013. URL <http://commons.apache.org/proper/commons-math/javadocs/api-3.2/org/apache/commons/math3/distribution/EnumeratedDistribution.html>.
- Yossi Azar and Thomas Erlebach, editors. *Algorithms ESA 2006*, volume 4168 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2006. ISBN 978-3-540-38875-3. doi: 10.1007/11841036. URL <http://www.springerlink.com/index/10.1007/11841036>.
- Y. Mansour, N. Nisan, and P. Tiwari. The computational complexity of universal hashing. In *Proceedings of the twenty-second annual ACM symposium on Theory of computing - STOC '90*, pages 235–243, New York, New York, USA, April 1990. ACM Press. ISBN 0897913612. doi: 10.1145/100216.100246. URL <http://dl.acm.org/citation.cfm?id=100216.100246>.
- Oracle Corporation. JSR-133: Java Memory Model and Thread Specification. 2004. URL <http://www.cs.umd.edu/~pugh/java/memoryModel/jsr133.pdf>.
- Klaus Mueller. GPU Programming: CUDA Threads. Technical report, Stony Brook

University, New York, New York, USA, 2009. URL http://www.cs.sunysb.edu/~mueller/teaching/cse591_GPU/threads.pdf.

C Nvidia. NVIDIA CUDA C Programming Guide. *Changes*, page 173, 2011.