
When and How Java Developers Give Up Static Type Safety

Doctoral Dissertation submitted to the
Faculty of Informatics of the Università della Svizzera Italiana
in partial fulfillment of the requirements for the degree of
Doctor of Philosophy

presented by
Luis Mastrangelo

under the supervision of
Prof. Matthias Hauswirth and Prof. Nathaniel Nystrom

June 2019

Dissertation Committee

Prof. Antonio Carzaniga	Università della Svizzera Italiana, Switzerland
Prof. Gabriele Bavota	Università della Svizzera Italiana, Switzerland
Prof. Jan Vitek	Northeastern University & Czech Technical University
Prof. Hridayesh Rajan	Iowa State University

Dissertation accepted on 13 June 2019

Research Advisor

Prof. Matthias Hauswirth

Co-Advisor

Prof. Nathaniel Nystrom

Ph.D. Program Director

Prof. Walter Binder

Ph.D. Program Director

Prof. Olaf Schenk

I certify that except where due acknowledgement has been given, the work presented in this thesis is that of the author alone; the work has not been submitted previously, in whole or in part, to qualify for any other academic award; and the content of the thesis is the result of work which has been carried out since the official commencement date of the approved research program.

Luis Mastrangelo
Lugano, 13 June 2019



*To Ankur Singhal,
In Memoriam*

“It is not only the violin that shapes the violinist, we are all shaped by the tools we train ourselves to use, and in this respect programming languages have a devious influence: they shape our thinking habits.”

– Edsger W. Dijkstra, 2001, To the members of the Budget Council^a

^a<http://www.cs.utexas.edu/users/EWD/transcriptions/OtherDocs/Haskell.html>

Abstract

The main goal of a static type system is to prevent certain kind of errors from happening at run-time. A type system is formulated as a set of constraints that gives any expression or term in a program a well-defined type. Yet mainstream programming languages are endowed with type systems that provide the means to loosen their static typing constraints through the *unsafe intrinsics* and *reflective capabilities* mechanisms.

We want to understand how and when developers give up these static constraints. This knowledge can be: a) a recommendation for current and future language designers to make informed decisions b) a reference for tool builders, *e.g.*, by providing more precise or new refactoring analyses, c) a guide for researchers to test new language features, or to carry out controlled programming experiments, and d) a guide for developers for better practices.

In this thesis we conducted two empirical studies to understand how these mechanisms—unsafe intrinsics and reflective capabilities—are used by Java developers when the static type system becomes too strict. We have devised usage patterns for both a subset of unsafe intrinsics and reflective capabilities. Usage patterns are recurrent programming idioms to solve a specific issue. We believe that having usage patterns can help us to better categorize use cases and thus understand how those features are used.

Acknowledgements

First of all, I want to thank both my advisors, Prof. Matthias Hauswirth and Prof. Nate Nystrom. Their expertise helped me to become a better researcher. They taught me how to ask the proper questions when facing a problem, that not only will help me in my professional career, but in everyday life.

In our *Unsafe* study, Luca Ponzanelli, Andrea Mocci and Prof. Michele Lanza collaborated with us to correlate *Unsafe* comments in the Questions and Answers database, Stack Overflow, mentioned in Chapter 3. Thanks to them.

Special thanks to Max Schäfer from Semmle. He executed our QL queries in the entire *lgtm* database used in Chapter 4. Without his help, this thesis would not have been possible.

Prof. Antonio Carzaniga and Prof. Gabriele Bavota made valuable comments during the Ph.D. proposal. They contributed to find an appropriate sample size in the cast study.

My group colleagues, Mohammad, Dmitri, Andrea, and Igor. With whom we had uncountable talks during my Ph.D. Special mention for Amanj and Nosheen, I learnt a lot in our coffee breaks. Thank you guys.

On the personal side, Vanesa, Cinzia, Mariano, Esteban, Dino, Olaf, Dimos, Jeff, Marisa, Rali, Gabi, Dilly, Sebastian, Eva, Stefania, and Grace who make me feel I found a family in Lugano. My former and current flat mates, Nora, Marie-line, Lea, and Celine; and Tanya and Matej, living with you guys was a pleasure.

My flat mate and dearest friend, Ankur Singhal, I wish you were here to see this moment. The Singhal family, Anil, Anita, and Ankita. I appreciate your support in the toughest moments.

And last but not least, I want to thank my family, for their unconditional support, and for always believing in me.

x

Contents

Contents	xi
List of Figures	xv
List of Tables	xvii
List of Listings	xvii
1 Introduction	1
1.1 Beyond Static Type Checking	2
1.2 Research Question	4
1.3 Thesis Outline	6
2 Literature Review	7
2.1 Benchmarks and Corpora	8
2.2 Tools for Mining Software Repositories	9
2.3 Large-scale Codebase Empirical Studies	11
2.3.1 Unsafe Intrinsic in Java	14
2.3.2 Reflective Capabilities	14
2.4 Conclusions	16
3 Empirical Study on the <i>Unsafe</i> API	17
3.1 The Risks of Compromising Safety	19
3.2 Is Unsafe Used?	22
3.3 Finding sun.misc.Unsafe Usage Patterns	28
3.4 Usage Patterns of sun.misc.Unsafe	30
3.4.1 Allocate an Object without Invoking a Constructor .	33
3.4.2 Process Byte Arrays in Block	33
3.4.3 Atomic Operations	34
3.4.4 Strongly Consistent Shared Variables	34

3.4.5	Park/Unpark Threads	35
3.4.6	Update Final Fields	35
3.4.7	Non-Lexically-Scoped Monitors	36
3.4.8	Serialization/Deserialization	36
3.4.9	Foreign Data Access and Object Marshaling	37
3.4.10	Throw Checked Exceptions without Being Declared	37
3.4.11	Get the Size of an Object or an Array	38
3.4.12	Large Arrays and Off-Heap Data Structures	38
3.4.13	Get Memory Page Size	38
3.4.14	Load Class without Security Checks	39
3.5	What is the Unsafe API Used for?	39
3.6	Conclusions	43
4	Empirical Study on the Cast Operator	45
4.1	Casts in Java	46
4.2	Issues Developers have using Cast Operations	50
4.3	Finding Cast Usage Patterns	52
4.4	Overview of the Sampled Casts	55
4.5	Cast Usage Patterns	56
4.5.1	TYPECASE	59
4.5.2	STASH	69
4.5.3	FACTORY	74
4.5.4	FAMILY	77
4.5.5	USERAWTYPE	79
4.5.6	EQUALS	82
4.5.7	REDUNDANT	86
4.5.8	COVARIANTRETURNType	88
4.5.9	SELECTOVERLOAD	90
4.5.10	KNOWNRETURNType	92
4.5.11	DESERIALIZATION	94
4.5.12	VARIABLESUPERTYPE	95
4.5.13	SOLESUBCLASSIMPLEMENTATION	98
4.5.14	NEWDYNAMICINSTANCE	99
4.5.15	OBJECTASARRAY	102
4.5.16	IMPLICITINTERSECTIONType	103
4.5.17	REMOVESWILDCARD	105
4.5.18	OPERANDSTACK	106
4.5.19	REFLECTIVEACCESSIBILITY	107
4.5.20	FLUENTAPI	108

4.5.21	COVARIANTGENERIC	110
4.5.22	COMPOSITE	112
4.5.23	GENERICARRAY	113
4.5.24	ACCESSSUPERCLASSFIELD	115
4.5.25	UNOCCUPIEDTYPEPARAMETER	117
4.6	Discussion	118
4.7	Conclusions	122
5	Conclusions	123
5.1	Java's Evolution	124
5.2	Limitations and Future Work	124
5.3	Lessons Learnt	126
5.4	Artifacts	126
A	Automatic Detection of Patterns using QL	129
A.1	Introduction to QL	129
A.2	Additional QL Classes and Predicates	131
B	JNIF: Java Native Instrumentation	139
B.1	Introduction	139
B.2	Related Work	141
B.3	Using JNIF	144
B.4	JNIF Design and Implementation	147
B.5	Validation	150
B.6	Performance Evaluation	152
B.7	Limitations	154
B.8	Conclusions	155
	Bibliography	159

Figures

3.1	sun.misc.Unsafe method usage on <i>Maven Central</i>	24
3.2	sun.misc.Unsafe field usage on <i>Maven Central</i>	25
3.3	<i>com.lmax:disruptor</i> call sites	29
3.4	<i>org.scala-lang:scala-library</i> call sites	30
3.5	Classes using off-heap large arrays	31
4.1	Cast Usage Pattern Occurrences	58
4.2	TYPECASE Variant Occurrences	59
4.3	STASH Variant Occurrences	70
4.4	EQUALS Variant Occurrences	82
B.1	Instrumentation time on DaCapo and Scala benchmarks . . .	157

Tables

1.1	Safe/Unsafe and Statically/Dynamically checked languages	2
3.1	Patterns and their occurrences in the Maven Central repository	32
3.2	Patterns and their alternatives	40
4.1	Statistics on Sampled Casts	55
4.2	Cast Usage Patterns	57
4.3	Categorization of Cast Usage Patterns	120

Listings

3.1	Instantiating an Unsafe object	19
3.2	sun.misc.Unsafe can violate type safety	20
3.3	sun.misc.Unsafe can crash the VM	20
3.4	sun.misc.Unsafe can violate a method contract	20
3.5	sun.misc.Unsafe can lead to uninitialized objects	21
3.6	sun.misc.Unsafe can lead to monitor deadlocks	21
4.1	Variable o (defined as Object) cast to String.	47
4.2	Compiled bytecode to the checkcast instruction.	47
4.3	Catch ClassCastException when a cast fails.	48
4.4	Runtime type test using instanceof before applying a cast.	48
4.5	Runtime type test using getClass before applying a cast.	48
4.6	Cast throws ClassCastException because of a forgotten guard.	50
4.7	Cast throws ClassCastException because of wrong cast target.	51
4.8	Cast throws ClassCastException because of generic inference.	51
4.9	Query for the <i>GuardByInstanceOf</i> variant.	66
4.10	Detection of a cast inside a switch statement	66
4.11	Detection of a cast guarded by the Class.isArray method.	67
4.12	Query for the <i>GuardByClassLiteral</i> variant.	67
4.13	Detection of the <i>LookupById</i> variant	73
4.14	Detection of the USERAWTYPE pattern.	81
4.15	Detection of the EQUALS pattern.	85
4.16	Detection query for the REDUNDANT pattern	88
4.17	COVARIANTRETURNTYPE detection query.	90
4.18	Query to detect the SELECTOVERLOAD pattern.	92
4.19	Detection of the DESERIALIZATION pattern.	95
4.20	Detection of the VARIABLESUPERTYPE pattern.	98
4.21	Detection of the SOLESUBCLASSIMPLEMENTATION pattern.	99
4.22	Detection of the NEWDYNAMICINSTANCE pattern.	101

4.23	Detection of the OBJECTAsARRAY pattern.	103
4.24	Detection of the IMPLICITINTERSECTIONTYPE pattern.	104
4.25	Detection of the REMOVEWILDCARD pattern.	105
4.26	Detection of the REFLECTIVEACCESSIBILITY pattern.	108
4.27	Detection of the FLUENTAPI pattern.	110
4.28	Query to detect the COVARIANTGENERIC pattern.	112
4.29	Detection of the ACCESSSUPERCLASSFIELD pattern.	116
A.1	Query to fetch all cast expressions in a project.	130
A.2	Query to fetch unused parameters.	130
A.3	Query to count methods with implementation.	130
A.4	Cast class definition.	131
A.5	Upcast class definition	132
A.6	OverloadedArgument class definition.	132
A.7	VarCast class definition.	133
A.8	controlByEqualityTest predicate definition.	133
A.9	controlByEqualsMethod predicate definition.	134
A.10	isSubtype predicate definition.	134
A.11	GetClassGuardsVarCast class definition.	134
A.12	AutoValueGenerated class definition.	135
A.13	NewDynamicInstanceAccess class definition.	135
A.14	ReflectiveMethodAccess class definition.	136
A.15	SetAccessibleTrueMethodAccess class definition.	137
A.16	notGenericRelated predicate definition.	137
B.1	Decoding a class	145
B.2	Encoding a class	145
B.3	Traversing all methods in a class	146
B.4	Instrumenting constructor entries	146
B.5	Instrumenting <init> methods	147
B.6	Running testapp	154
B.7	Running dacapo	154
B.8	Running dacapo	154
B.9	Running full eval five times	154
B.10	Plots	154

Chapter 1

Introduction

In programming language design, the main goal of a *static* type system is to prevent certain kind of errors from happening at run-time. A type system is formulated as a set of constraints that gives any expression or term in a program a well-defined type. As Pierce [2002] states: “A type system can be regarded as calculating a kind of *static* approximation to the run-time behaviors of the terms in a program.” These constraints are enforced by the *type checker* either when compiling or linking the program. Thus, any program not satisfying the constraints stated within a type system is simply rejected by the type checker.

Besides detecting early errors, a type checker can also be an invaluable *maintenance* tool. For instance, it can assist an IDE to perform refactor analysis, such as renaming a method or field. A static type system can be helpful to enforce disciplined programming. When composing large-scale software, *modular languages* are built-up of types, shown in the interfaces of modules. Along these lines, type systems can be useful for *documenting* purposes. Type annotations, *e.g.*, in method and fields declarations, can provide useful hints to the developer. Since type annotations are meant to be checked every time the program is compiled, this information cannot be outdated, unlike comments in the source text.

Static type systems can aid to generate more *efficient* machine code, *e.g.*, choosing a different representation for integer or real values at run-time. Furthermore, in statically checked languages, *e.g.*, Java or Rust, many checks are performed at compile-time, instead of being performed otherwise at run-time. Compare this to dynamically checked languages, where all checks need to be performed at runtime, degrading performance. Table 1.1 shows where mainstream languages fit in the safe/unsafe and stati-

cally/dynamically checked spectrum.

Table 1.1. Safe/Unsafe and Statically/Dynamically checked languages

	Statically checked	Dynamically checked
Safe	Haskell, SML, Java, C#, Rust, <i>etc.</i>	Python, Lisp, Racket, <i>etc.</i>
Unsafe	C, C++, <i>etc.</i>	

1.1 Beyond Static Type Checking

Nevertheless, often the static approximation provided by a type system is not precise enough. Being static, the analysis done by the type checker needs to be conservative: It is better to reject programs that are valid, but whose validity cannot be ensured by the type checker, rather than accept some invalid programs. However, there are situations when the developer has more information about the program that is too complex to explain in terms of typing constraints. To that end, programming languages often provide *mechanisms* that make the typing constraints less strict to permit more programs to be valid, at the expense of causing more errors at run-time. These mechanisms are essentially two: *Unsafe Intrinsic*s and *Reflective Capabilities*.

Unsafe Intrinsic

Some programming languages provide *unsafe intrinsic*s, the ability to perform certain operations *without* being checked by the compiler. They are *unsafe* because any misuse made by the programmer can compromise the entire system, *e.g.*, corrupting data structures without notice, or crashing the run-time system. In other words, all guarantees provided by a static type system are undermined by the inclusion of unsafe intrinsic

Unsafe intrinsic can be seen in safe languages, *e.g.*, Java, C#, Rust, or Haskell. Foreign Function Interface (FFI), *i.e.*, calling native code from within a safe environment is unsafe. It is so because the run-time system cannot guarantee anything about the native code. In addition to FFI, some safe languages offer so-called *unsafe blocks*, *i.e.*, making unsafe operations

within the language itself, *e.g.*, C#¹ and Rust.² For instance, when using unsafe blocks in Rust, the developer can dereference a raw pointer, making the application crash.

Other languages instead provide an API to perform unsafe operations, *e.g.*, Haskell³ and Java. But in the case of Java, the API to make unsafe operations, `sun.misc.Unsafe`, is unsupported⁴ and undocumented. For instance, by invoking the `allocateInstance` on an instance of `sun.misc.Unsafe` the developer can allocate an object without calling any constructor, thus, violating Java's type system guarantees. It was originally intended for internal use within the JDK, but as we shall see later on, it is used outside the JDK as well.

Reflective Capabilities

Many programming languages provide some sort of *reflective capabilities*, *i.e.*, it enables an executing program to examine or “introspect” upon itself. Much of the Java Reflection API resides in the `java.lang.reflect` package, allowing the running program to obtain information about classes and objects. By using reflection, it is possible to dynamically create instances of a class at run-time as well, *e.g.*, through the `Class` class. C# provides analogous classes, *e.g.*, the `Type` class, to achieve the same functionality. When reflection is used, many checks that were done by the type checker statically (at compile-time) now need to be performed dynamically (at run-time).

Programming languages with subtyping such as Java or C++ provide a mechanism to *view* an expression as a different type as it was defined, a form of *lightweight* reflection. This mechanism is often called *casting* and usually takes the form $(T) \ t$. Casting can be in two directions: *upcast* and *downcast*. An upcast conversion happens when converting from a reference type S to a reference type T , provided that T is a *supertype* of S . An upcast does not require any explicit casting operation nor compiler check. However, as we shall see later on, there are situations where an upcast requires an explicit casting operation. On the other hand, a downcast happens when converting from a reference type S to a reference type T , provided that T is a *subtype* of S . Unlike upcasts, downcasts require a run-time check to ver-

¹<https://docs.microsoft.com/en-us/dotnet/csharp/language-reference/language-specification/unsafe-code>

²<https://doc.rust-lang.org/book/second-edition/ch19-01-unsafe-rust.html>

³<http://hackage.haskell.org/package/base-4.11.1.0/docs/System-IO-Unsafe.html>

⁴<http://www.oracle.com/technetwork/java/faq-sun-packages-142232.html>

ify that the conversion is indeed valid. For instance, Java provides the cast operator; while C++ with run-time type information (RTTI) provides the `dynamic_cast<>` operator. This implies that downcasts provide the means to bypass the static type system. By avoiding the type system, downcasts can pose potential threats, because it is like the developer saying to the compiler: “*Trust me here, I know what I’m doing*”. Being an escape-hatch to the type system, a cast is often seen as a design flaw or code smell [Tufano et al., 2015] in an object-oriented system.

Featherweight Java [Igarashi et al., 2001] is a minimal core calculus for Java. The language features it models are “mutually recursive class definitions, object creation, field access, method invocation, method override, method recursion through this, subtyping, and casting.” In their model, there are three ways a computation may get stuck. When accessing a field not declared for a class, by invoking a nonexistent method for a class, or when casting to a class that is *not* a supertype of the object’s *runtime* class. They proved that on any well-typed program, the first two never happen, and the third one never happens on well-typed programs *without downcasts*. From the theoretical point of view, casts make typing rules unsound, *i.e., a well-typed program with casts may get stuck*.

1.2 Research Question

If static type systems aim to prevent certain kind of errors from happening at run-time, yet they provide the means to loosen their typing constraints, why exactly does one need to do so? Are these mechanisms actually used in real-world code? If yes, then how so? This triggers our **main research question**:

MRQ

For what purpose do developers give up static type checking?

We have confidence that this knowledge can be: a) a reference for current and future language designers to make informed decisions about programming languages, *e.g.*, the adoption of *Variable Handles* in Java 9 [Lea, 2014], or the addition of *Smart Casts* in Kotlin,⁵ b) a reference for tool builders, *e.g.*, by providing more precise or new refactoring analyses, c) a guide for researchers to test new language features, *e.g.*, Winther [2011] or

⁵<https://kotlinlang.org/docs/reference/typecasts.html#smart-casts>

to carry out controlled experiments about programming, *e.g.*, Stuchlik and Hanenberg [2011] and d) a guide for developers for best or better practices.

To answer our question above, we empirically studied how the two aforementioned mechanisms—unsafe intrinsics and reflective capabilities—are used by developers. Since we seek to *understand* how these mechanisms are used, our methodology is based on qualitative analysis. Our qualitative data to analyse is source code text (to study unsafe intrinsics we performed a preliminary analysis on intermediate code). In particular, in both studies we performed manual qualitative (static) analysis. It is static because we have analysed *only* the *source text* (*cf.* dynamic analysis). We performed repository mining to gather the source code text to analyse.

Since any kind of language study must be language-specific, we focus on Java given its wide usage and relevance for both research and industry.⁶ Moreover, we focus on the Java Unsafe API to study unsafe intrinsics, given that the Java Native Interface already has been studied in Tan et al. [2006]; Tan and Croft [2008]; Kondoh and Onodera [2008]; Sun and Tan [2014]; Li and Tan [2009]. Similarly, although casting uses run-time type information like the Java’s reflection API, the reflection API has been studied in Livshits [2006]; Livshits et al. [2005]; Landman et al. [2017].

To better drive our *main research question*, we propose to answer the following set of sub-questions. To answer these research sub-questions, we have devised *usage patterns* for both the Unsafe API and casting. Usage patterns are *recurrent programming idioms* used by developers to solve a specific issue. We believe that having usage patterns can help us to better categorize use cases and thus understand how these mechanisms are used. These patterns can provide an insight into how the language is being used by developers in real-world applications. Overall these sub-questions will help us to answer our MRQ:

Unsafe API

RQ/U1 : To what extent does the Unsafe API impact common application code? We want to understand to what extent code actually uses Unsafe or depends on it.

RQ/U2 : How and when are Unsafe features used? We want to investigate what functionality third-party libraries require from Unsafe. This could point out ways in which the Java language

⁶<https://www.tiobe.com/tiobe-index/>

and/or the JVM can be evolved to provide the same functionality, but in a safer way.

These questions have been already answered in our previous published study on the Unsafe API in Java [Mastrangelo et al., 2015].

Casting

RQ/C1 : How frequently is casting used in common application code?

To what extent does application code actually use casting operations?

RQ/C2 : How and when casts are used? If casts are used in application code, how and when do developers use them?

RQ/C3 : How recurrent are the patterns for which casts are used? In addition to understand how and when casts are used, we want to measure how often developers need to resort to certain idioms to solve a particular problem.

The results of this study have been submitted for publication to OOP-SLA'19.

1.3 Thesis Outline

The rest of this thesis is as follows. In Chapter 2 we give a review of the literature in empirical studies of programming languages features. In particular, Sections 2.3.1 and 2.3.2 review the *state-of-the-art* of the different aspects related to the two proposed studies. Chapter 3 presents a summary of our *Unsafe* study; while in Chapter 4 we present our *casting* study. Finally, Chapter 5 presents the conclusions for the thesis.

The Appendix A contains an introduction to QL—the language we used to approximate automatic detection of patterns—and reference material used in our casting study. Appendix B—although not directly related—describes our bytecode analysis library used in some experiments in both Chapters 3 and 4.

Chapter 2

Literature Review

Understanding how developers use language features and APIs is a broad topic. There is plenty of research in the computer science literature about empirical studies of programs which involves multiple *dimensions* directly related to our plan. Over the last decades, researchers always have been interested in understanding what kind of programs developers write. The motivation behind these studies is quite broad, and has been shifted to the needs of researchers, together with the evolution of computer science itself.

For instance, to measure the advantages between compilation and interpretation in Basic, Hammond [1977] studied a representative dataset of programs. Knuth [1971] started to study Fortran programs. By knowing what kind of programs arise in practice, a compiler optimizer can focus in those cases, and therefore can be more effective. Adding to Knuth’s work, Shen et al. [1990] conducted an empirical study for parallelizing compilers. Similar works have been done for Cobol Salvadori et al. [1975]; Chevance and Heidet [1978], Pascal Cook and Lee [1982], and APL Saal and Weiss [1975, 1977] programs. Miller et al. [1990, 1995]; Forrester and Miller [2000] studied the reliability of programs using *fuzz* testing. Dieckmann and Hölzle [1999] studied the memory allocating behavior in the SPECjvm98 benchmarks.¹

But there is more than empirical studies at the source code level. A machine instruction set is effectively another kind of language. Therefore, its design can be affected by how compilers generate machine code. Several studies targeted the JVM Collberg et al. [2007]; O’Donoghue et al. [2002]; Antonioli and Pilz [1998]; while Cook [1989] did a similar study for Lilith in the past.

¹<https://www.spec.org/jvm98/>

The importance of conducting empirical studies of programs gave rise to the International Conference on Mining Software Repositories² in 2004.

Outline

When conducting empirical studies about programs, multiple dimensions are involved. The first one is *What to analyse?* Benchmarks and corpora are used as a source of programs to analyse. Another aspect is how to select good candidates projects from a large-base software repository. This is presented in Section 2.1. After the selection of programs to analyse is set, comes the question *how to analyse them?* An overview of what tools are available to extract information from software repositories is given in Section 2.2. With this infrastructure, *what questions do researchers ask?* In Section 2.3, we give an overview of large-scale empirical studies that show what kind of questions researchers ask. In particular, this section ends by presenting the related work more specific to the Unsafe API and Casting in Sections 2.3.1 and 2.3.2 respectively. Finally, Section 2.4 concludes this chapter.

2.1 Benchmarks and Corpora

Benchmarks are crucial to properly evaluate and measure product development. This is key for both research and industry. One popular benchmark suite for Java is the DaCapo Benchmark [Blackburn et al., 2006]. This suite has been already cited in more than thousand publications, showing how important is to have reliable benchmark suites. The SPECjvm2008³ (Java Virtual Machine Benchmark) and SPECjbb2000⁴ (Java Business Benchmark) are another popular Java benchmark suite.

Another suite has been developed by Tempero et al. [2010]. They provide a corpus of curated open source systems to facilitate empirical studies on source code. On top of Qualitas Corpus, Dietrich et al. [2017b] provide an executable corpus of Java programs. This allows any researcher to experiment with both static and dynamic analysis.

²<http://www.msrconf.org/>

³<https://www.spec.org/jvm2008/>

⁴<https://www.spec.org/jbb2000/>

For any benchmark or corpus to be useful and reliable, it must faithfully represent real world code. For instance, DaCapo applications were selected to be diverse real applications and ease of use, but they “excluded GUI applications since they are difficult to benchmark systematically.” Along these lines, Allamanis and Sutton [2013] go one step further and provide a large-scale (14,807) curated corpus of open source Java projects.

With the advent of cloud computing, several source code management (SCM) hosting services have emerged, *e.g.*, *GitHub*, *GitLab*, *Bitbucket*, and *SourceForge*. These services allow the developer to work with different SCMs, *e.g.*, *Git*, *Mercurial*, *Subversion* to host their open source projects. These projects are usually taken as a representation of real-world applications. Thus, while not curated corpora, these hosting services are commonly used to conduct empirical studies.

Another dimension to consider when analysing large codebases, is how relevant the repositories are. Lopes et al. [2017] conducted a study to measure code duplication in *GitHub*. They found out that much of the code there is actually duplicated. This raises a flag when considering which projects to analyse when mining software repositories.

Baxter et al. [1998] propose a clone detection algorithm using Abstract Syntax Trees, while Rieger and Ducasse propose a visual detection for clones. Yuan and Guo [2011]; Chen et al. instead propose Count Matrix-based approach to detect code clones.

Nagappan et al. [2013] have developed the Software Projects Sampling (SPS) tool. SPS tries to find a maximal set of projects based on representativeness and diversity. Diversity dimensions considered include total lines of code, project age, activity, number of contributors, total code churn, and number of commits.

2.2 Tools for Mining Software Repositories

When talking about mining software repositories, we refer to extracting any kind of information from large-scale codebase repositories. Usually doing so requires several engineering but challenging tasks. The most common being downloading, storing, parsing, analysing and properly extracting information from different kinds of artifacts. In this scenario, there are several tools that allows a researcher or developer to query information about software repositories.

Urma and Mycroft [2012] evaluated seven source code query languages: *Java Tools Language* [Cohen and Maman], *SOUL* [De Roover et al., 2011], *Browse-By-Query*⁵, *JQuery* [Volder, 2006], *.QL* [de Moor et al., 2007], *Jackpot*⁶, and *PMD*⁷. They have implemented, whenever possible, four use cases using the tools mentioned above. They concluded that only *SOUL* and *.QL* have the minimal features to implement all their use cases.

Dyer et al. [2013a,c] built *Boa*, both a domain-specific language and an online platform⁸. It is used to query software repositories on two popular hosting services, *GitHub* and *SourceForge*. The same authors of *Boa* conducted a study on how new Java features, e.g., *Assertions*, *Enhanced-For Loop*, *Extends Wildcard*, were adopted by developers over time [Dyer et al., 2014]. This study is based *SourceForge* data. The current problem with *SourceForge* is that is outdated.

To this end, Gousios [2013] provides an offline mirror of *GitHub* that allows researchers to query any kind of that data. Later on, Gousios et al. [2014] published the dataset construction process of *GitHub*.

Similar to *Boa*, *lgtm*⁹ is a platform to query software projects properties. It works by querying repositories from *GitHub*. But it does not work at a large-scale, i.e., *lgtm* allows the user to query just a few projects. Unlike *Boa*, *lgtm* is based on QL—before named *.QL*—, an object-oriented domain-specific language to query recursive data structures based on Datalog [Augustinov et al., 2016]. Another static analysis framework based on Datalog is *Doop* [Bravenboer and Smaragdakis, b]. Since QL and *Doop* are based on Datalog, both are well-suited to perform points-to analysis and data-flow analysis. However, scaling such analysis to a large-scale study remains an open problem.

On top of *Boa*, Tiwari et al. [2017] built *Candoia*¹⁰. Although it is not a mining software repository *per se*, it eases the creation of mining applications.

Another tool to analyse large software repositories is presented in Brandauer and Wrigstad [2017]. In this case, the analysis is dynamic, based on program traces. At the time of this writing, the service¹¹ was unavailable

⁵<http://browsebyquery.sourceforge.net/>

⁶<http://wiki.netbeans.org/Jackpot>

⁷<https://pmd.github.io/>

⁸<http://boa.cs.iastate.edu/>

⁹<https://lgtm.com/>

¹⁰<http://candoia.github.io/>

¹¹<http://www.spencer-t.racing/datasets>

for testing.

Bajracharya et al. [2009] provide a tool to query large code bases by extracting the source code into a relational model. Sourcegraph¹² is a tool that allows regular expression and diff searches. It integrates with source repositories to ease navigate software projects.

Posnett et al. [2010] have extended ASM [Bruneton et al., 2002b] to detect meta-patterns, *i.e.*, structural patterns of object-oriented interaction. Hu and Sartipi [2008] used both dynamic and static analysis to discover design patterns, while Arcelli et al. [2008] used only dynamic analysis.

Trying to unify analysis and transformation tools, Vinju and Cordy [2006] and Klint et al. [2009] built *Rascal*, a DSL that aims to bring them together by querying the AST of a program. Spoon is a Java library “to analyse, rewrite, transform, transpile Java source code” [Pawlak et al., 2015]. It supports symbol resolution natively and match code elements using Spoon patterns. Probably one of the most mature libraries to parse and manipulate Java source code is Eclipse Java Development Tools (JDT).¹³ Besides parsing Java source code, Eclipse JDT has the ability to compile, run, and debug Java source code. ExtendJ [Ekman and Hedin, 2007] is an extensible Java compiler. It supports semantic analysis and bytecode generation. With Java 8, it is possible to write plug-ins for the javac compiler. By writing a compiler plug-in, it is possible to add extra-compile checks, perform code transformations and custom analysis. JavaParser,¹⁴ as its name suggests, is a parser for Java. The main issue with JavaParser is that it lacks the ability to perform symbol resolution integrated with project dependencies. javalang¹⁵ is a library written in Python to parse Java source code.

In early prototypes of our cast study we have used javalang and a custom javac plug-in,¹⁶ and JavaParser¹⁷ to parse and analyse Java source code.

2.3 Large-scale Codebase Empirical Studies

In the same direction as our plan, Callaú et al. [2013] performed an empirical study to assess how much the dynamic and reflective features of

¹²<https://sourcegraph.com>

¹³<https://www.eclipse.org/jdt/>

¹⁴<http://javaparser.org/>

¹⁵<https://github.com/c2nes/javalang>

¹⁶<https://gitlab.com/acuarica/java-cast-inspection>

¹⁷<https://gitlab.com/acuarica/java-cast-study>

Smalltalk are actually used in practice. Analogously, Richards et al. [2010, 2011]; Wei et al. [2016] conducted a similar study, but in this case targeting JavaScript’s dynamic behavior and in particular the `eval` function. Also, for JavaScript, Madsen and Andreasen [2014] analysed how fields are accessed via strings, while Jang et al. [2010] analysed privacy violations. Similar empirical studies were done for PHP [Hills et al., 2013; Dahse and Holz, 2015; Doyle and Walden, 2011] and Swift [Rebouças et al., 2016]. Pinto et al. [2015] conducted a large-scale study on how concurrency is used in Java

Going one step forward, Ray et al. [2017] studied the correlation between programming languages and defects. One important note is that they choose relevant projects by popularity, measured by how many times was *starred* in *GitHub*. We argue that it is more important to analyse projects that are *representative*, not *popular*.

Gorla et al. [2014] mined a large set of Android applications, clustering applications by their description topics and identifying outliers in each cluster with respect to their API usage. Grechanik et al. [2010] also mined large scale software repositories to obtain several statistics on how source code is actually written.

For Java, Dietrich et al. [2017a] conducted a study about how programmers use contracts in *Maven Central*¹⁸. Dietrich et al. [2014] have studied how API changes impact Java programs. They have used the Qualitas Corpus [Tempero et al., 2010] mentioned above for their study.

Tufano et al. [2015, 2017] studied when code smells are introduced in source code. Palomba et al. [2015] contribute a dataset of five types of code smells together with a systematic procedure for validating code smell datasets. Palomba et al. [2013] propose to detect code smells using change history information.

Nagappan et al. [2015] conducted a study on how the `goto` statement is used in C. They used *GitHub* as a data source for C programs. They concluded that `goto` statements are most used for *handling errors* and *cleaning up resources*.

Static vs. Dynamic Analysis. Given the dynamic nature of JavaScript, most of the studies mentioned above for JavaScript perform dynamic analysis. However, Callaú et al. [2013] uses static analysis to study a dynamically checked language. For Java, most empirical studies use static analysis. This is due the fact of the availability of input data. Finding valid input data for test cases is not a trivial task, even less to make it scale. For JavaScript,

¹⁸<http://central.sonatype.org/>

having a big corpus of web-sites generating valid input data makes more feasible to implement dynamic analysis.

Exceptions

Kery et al. [2016]; Asaduzzaman et al. [2016] focus on exceptions. They conducted empirical studies on how programmers handle exceptions in Java code. The work done by Nakshatri et al. [2016] categorized them into patterns. Coelho et al. [2015] used a more dynamic approach by analysing stack traces and code issues in *GitHub*. Kechagia and Spinellis [2014] analysed how undocumented and unchecked exceptions cause most of the exceptions in Android applications.

Programming Language Features

Programming language design has been always a hot topic in computer science literature. It has been extensively studied in the past decades. There is a trend in incorporating programming features into mainstream object-oriented languages, *e.g.*, lambdas in Java 8¹⁹, C++11²⁰ and C# 3.0²¹; or parametric polymorphism, *i.e.*, generics, in Java 5.^{22,23} For instance, Java generics were designed to extend Java's type system to allow "a type or method to operate on objects of various types while providing compile-time type safety" [Gosling et al., 2013]. However, it was later shown [Amin and Tate, 2016] that compile-time type safety was not fully achieved.

Mazinanian et al. [2017] and Uesbeck et al. [2016] studied how developers use lambdas in Java and C++ respectively. The inclusion of generics in Java is closely related to collections. Parnin et al. [2011, 2013] studied how generics were adopted by Java developers. They found that the use of generics does not significantly reduce the number of type casts.

Costa et al. [2017] have mined *GitHub* corpus to study the use and performance of collections, and how these usages can be improved. They found that in most cases there is an alternative usage that improves performance.

¹⁹<https://docs.oracle.com/javase/specs/jls/se8/html/jls-15.html#jls-15.27>

²⁰<http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2006/n1968.pdf>

²¹https://msdn.microsoft.com/en-us/library/bb308966.aspx#csharp3.0overview_topic7

²²<https://docs.oracle.com/javase/1.5.0/docs/guide/language/generics.html>

²³<http://www.oracle.com/technetwork/java/javase/generics-tutorial-159168.pdf>

Another study about how a programming language feature is used is done in Tempero et al. [2008]. They conducted a study on how inheritance is used in Java programs.

This kind of studies give an insight of the adoption of lambdas and generics; which can drive future direction for language designers and tool builders, while providing developers with best practices.

2.3.1 Unsafe Intrinsic in Java

Oracle provides the `sun.misc.Unsafe` class for low-level programming, *e.g.*, synchronization primitives, direct memory access methods, array manipulation and memory usage. Although `sun.misc.Unsafe` is not officially documented, it is being used in both industrial applications and research projects [Korland et al., 2010; Pukall et al.; Gligoric et al., 2011] outside the JDK, compromising the safety of the Java ecosystem.

Oracle’s software engineer Paul Sandoz performed an informal analysis of *Maven Central* artifacts and usages in Greptime [Sandoz, 2015] and conducted a unscientific user survey to study how the *Unsafe* API is used [Sandoz, 2014]. The survey consists of 7 questions²⁴ that help to understand what pieces of `sun.misc.Unsafe` should be mainstreamed.

Tan et al. [2006] propose a combination of static and dynamic checks to provide a safe variant of the Java Native Interface (JNI). They have identified several loopholes that may cause unsafe interoperability between Java and native code. The language extension provided by Bubak and Kurzyniec [2000] allows the developer to interleave Java and native code in the same compilation unit. However, the native code is not—statically nor dynamically—checked, causing a possible JVM crash. Tan and Croft [2008] and Kondoh and Onodera [2008] conducted an empirical security study to describe a taxonomy to classify bugs when using JNI. Sun and Tan [2014] develop a method to isolate native components in Android applications. Li and Tan [2009] analyse the discrepancy between how exceptions are handled in native code and Java.

2.3.2 Reflective Capabilities

Livshits [2006]; Livshits et al. [2005] “describes an approach to call graph construction for Java programs in the presence of reflection.” He has de-

²⁴<http://www.infoq.com/news/2014/02/Unsafe-Survey>

vised some common usage patterns for reflection. Most of the patterns use casts. We plan to categorize all cast usages, not only where reflection is used.

Landman et al. [2017] have analysed the relevance of static analysis tools with respect to reflection. They conducted an empirical study to check how often the reflection API is used in real-world code. They have devised reflection AST patterns, which often involve the use of casts. Finally, they argue that controlled programming experiments on subjects need to be correlated with real-world use cases, *e.g.*, *GitHub* or *Maven Central*.

Casting operations in Java²⁵ allows the developer to view a reference at a different type as it was declared. The related `instanceof` operator²⁶—written `e instanceof T`—tests whether a reference `e` could be cast to a different type `T` without throwing `ClassCastException` at run-time.

Winther [2011] has implemented a path sensitive analysis that allows the developer to avoid casting once a guarded `instanceof` is provided. He proposes four cast categorizations according to their run-time type safety: *Guarded Casts*, *Semi-Guarded Casts*, *Unguarded Casts*, and *Safe Casts*.

Tsantalis et al. [2008] present an Eclipse plug-in that identifies type-checking bad smells, a “variation of an algorithm that should be executed, depending on the value of an attribute”. They provide refactoring analysis to remove the detected smells by introducing inheritance and polymorphism. This refactoring will introduce casts to select the right type of the object.

Controlled Experiments on Subjects. There is an extensive literature *per se* in controlled experiments on subjects to understand several aspects in programming, and programming languages. For instance, Soloway and Ehrlich [1984] tried to understand how expert programmers face problem solving. Budd et al. [1980] conducted an empirical study on how effective mutation testing is. Prechelt [2000] compared how a given—fixed—task was implemented in several programming languages. LaToza and Myers [2010] realize that, in essence, programmers need to answer reachability questions to understand large codebases. Several authors Stuchlik and Hanenberg [2011]; Mayer et al. [2012]; Harlin et al. [2017] measure whether using a static-type system improves programmers productivity. They compare how a static and a dynamic type system impact on productivity. The common setting for these studies is to have a set of programming problems.

²⁵<https://docs.oracle.com/javase/specs/jls/se8/html/jls-15.html#jls-15.16>

²⁶<https://docs.oracle.com/javase/specs/jls/se8/html/jls-15.html#jls-15.20.2>

Then, let a group of developers solve them in both a static and dynamic languages. For this kind of studies to reflect reality, the problems to be solved need to be representative of the real-world code. Having artificial problems may lead to invalid conclusions. The work by Wu and Chen [2017]; Wu et al. [2017] goes towards this direction. They have examined programs written by students to understand real debugging conditions. Their focus is on ill-typed programs written in Haskell.

2.4 Conclusions

The Java Native Interface and Java’s reflection API are well-studied topics. Several studies have been conducted to understand why developers use these features, and several analyses have been devised to check whether their usage is correct.

But Java’s unsafe intrinsics and reflection capabilities comprise more than JNI and reflection API. Unsafe operations can be performed by using the undocumented `sun.misc.Unsafe` class. The cast operator provides a lightweight form of reflection. However—to our knowledge—these features have never been studied before in the literature. This thesis provides the first empirical studies on the *Unsafe* API and cast operator in Java. In our work [Mastrangelo et al., 2015] we extend Sandoz’ work by performing a comprehensive study of the *Maven Central* software repository to analyse how and when `sun.misc.Unsafe` is being used. This study is summarized in Chapter 3. We refined the categorization performed by Winther [2011] to answer our *RQ/C2 (How and when casts are used?)*. This is described in Chapter 4. We believe that understanding how and when developers use these features can provide informed decisions for the future of Java while providing a guide for developers with better or best practices.

Chapter 3

Empirical Study on the *Unsafe* API

The Java Virtual Machine (JVM) executes Java bytecode and provides other services for programs written in many programming languages, including Java, Scala, and Clojure. The JVM was designed to provide strong safety guarantees. However, many widely used JVM implementations expose an API that allows the developer to access low-level, unsafe features of the JVM and underlying hardware, features that are unavailable in safe Java bytecode. This API is provided through an undocumented¹ class, `sun.misc.Unsafe`, in the Java reference implementation produced by Oracle.

Other virtual machines provide similar functionality. For example, the C# language provides an unsafe construct on the .NET platform,² and Racket provides unsafe operations.³

The operations `sun.misc.Unsafe` provides can be dangerous, as they allow developers to circumvent the safety guarantees provided by the Java language and the JVM. If misused, the consequences can be resource leaks, deadlocks, data corruption, and even JVM crashes.^{4 5 6 7 8}

We believe that `sun.misc.Unsafe` was introduced to provide better performance and more capabilities to the writers of the Java runtime library. However, `sun.misc.Unsafe` is increasingly being used in third-party frameworks and libraries. Application developers who rely on Java's safety guar-

¹<http://www.oracle.com/technetwork/java/faq-sun-packages-142232.html>

²[https://msdn.microsoft.com/en-us/en-en/library/chfa2zb8\(v=vs.90\).aspx](https://msdn.microsoft.com/en-us/en-en/library/chfa2zb8(v=vs.90).aspx)

³<http://docs.racket-lang.org/reference/unsafe.html>

⁴<https://groups.google.com/d/msg/elasticsearch/Nh-kXI5J6Ek/WXIZKhhGVHkJ>

⁵<https://github.com/EsotericSoftware/kryo/issues/219>

⁶<https://github.com/dain/snappy/issues/24>

⁷https://netbeans.org/bugzilla/show_bug.cgi?id=229655

⁸https://netbeans.org/bugzilla/show_bug.cgi?id=244914

antees have to trust the implementers of the language runtime environment (including the core runtime libraries). Thus the use of `sun.misc.Unsafe` in the runtime libraries is no more risky than the use of an unsafe language to implement the JVM. However, the fact that more and more “normal” libraries are using `sun.misc.Unsafe` means that application developers have to trust a growing community of third-party Java library developers to not inadvertently tamper with the fragile internal state of the JVM.

Given that the benefits of safe languages are well known, and the risks of unsafe languages are obvious, why exactly does one need unsafe features in third-party libraries? Are those features used in real-world code? If yes, how are they used, and what are they used for? More precisely, we want to answer the following research questions:

RQ/U1 : To what extent does the Unsafe API impact common application code? We want to understand to what extent code actually uses `Unsafe` or depends on it.

RQ/U2 : How and when are Unsafe features used? We want to investigate what functionality third-party libraries require from `Unsafe`. This could point out ways in which the Java language and/or the JVM can be evolved to provide the same functionality, but in a safer way.

If *Unsafe* is not just dangerous, but also confusing or difficult to use, then its use by third-party developers is particularly problematic. If there are specific *Unsafe* features or usage patterns that developers worry about, it would make sense to evolve Java or the JVM to provide safer alternatives in that direction.

We studied a large repository of Java code, *Maven Central*, to answer these questions. We have analysed 74 GB of compiled Java code, spread over 86,479 Java archives, to determine how Java’s unsafe capabilities are used in real-world libraries and applications. We found that 25% of Java bytecode archives depend on unsafe third-party Java code, and thus Java’s safety guarantees cannot be trusted. We identify 14 different usage patterns of Java’s unsafe capabilities, and we provide supporting evidence for why real-world code needs these capabilities. Our long-term goal is to provide a strong foundation to make informed decisions in the future evolution of the Java language and virtual machine, and for the design of new language features to regain safety in Java.

Outline

We have already published our work on how developers use the Unsafe API in Java [Mastrangelo et al., 2015]. In this thesis we outline the risks of using the *Unsafe* API in Section 3.1. Then we answer *RQ/U1* (To what extent does the Unsafe API impact common application code?) in Section 3.2. To answer *RQ/U2* (How and when are Unsafe features used?), first we introduce our methodology and the patterns we found in Sections 3.3 and 3.4 respectively, to then present how the patterns we found could be implemented in a safer way in Section 3.5. Finally, Section 3.6 concludes this chapter.

3.1 The Risks of Compromising Safety

We outline the risks of *Unsafe* by illustrating how the improper use of *Unsafe* violates Java’s safety guarantees.

In Java, the unsafe capabilities are provided as instance methods of class `sun.misc.Unsafe`. Access to the class has been made less than straightforward. Class `sun.misc.Unsafe` is final, and its constructor is not public. Thus, creating an instance requires some tricks. For example, one can invoke the private constructor via reflection. This is not the only way to get hold of an unsafe object, but it is the most portable.

```
1 Constructor<Unsafe> c = Unsafe.class.getDeclaredConstructor();  
2 c.setAccessible(true);  
3 Unsafe unsafe = c.newInstance();
```

Listing 3.1. Instantiating an Unsafe object

Given the unsafe object, one can now simply invoke any of its methods to directly perform unsafe operations.

Violating Type Safety

In Java, variables are strongly typed. For example, it is impossible to store an int value inside a variable of a reference type. *Unsafe* can violate that guarantee: it can be used to store a value of any type in a field or array element.

```
1 class C {  
2     private Object f = new Object();  
3 }  
4 long fieldOffset = unsafe.objectFieldOffset(  
5     C.class.getDeclaredField("f") );  
6 C o = new C();  
7 unsafe.putInt(o, fieldOffset, 1234567890);    // f now points to nirvana
```

Listing 3.2. sun.misc.Unsafe can violate type safety

Crashing the Virtual Machine

A quick way to crash the VM is to free memory that is in a protected address range, for example by calling `freeMemory` as follows.

```
unsafe.freeMemory(1);
```

Listing 3.3. sun.misc.Unsafe can crash the VM

In Java, the normal behavior of a method to deal with such situations is to throw an exception. Being unsafe, instead of throwing an exception, this invocation of `freeMemory` crashes the VM.

Violating Method Contracts

In Java, a method that does not declare an exception cannot throw any checked exceptions. *Unsafe* can violate that contract: it can be used to throw a checked exception that the surrounding method does not declare or catch.

```
1 void m() {  
2     unsafe.throwException(new Exception());  
3 }
```

Listing 3.4. sun.misc.Unsafe can violate a method contract

Uninitialized Objects

Java guarantees that an object allocation also initializes the object by running its constructor. *Unsafe* can violate that guarantee: it can be used to

allocate an object without ever running its constructor. This can lead to objects in states that the objects' classes would not seem to admit.

```
1 class C {  
2     private int f;  
3     public C() { f = 5; }  
4     public int getF() { return f; }  
5 }  
6  
7 C c = (C)unsafe.allocateInstance(C.class);  
8 assert c.getF()==5; // violated
```

Listing 3.5. `sun.misc.Unsafe` can lead to uninitialized objects

Monitor Deadlock

Java provides synchronized methods and synchronized blocks. These constructs guarantee that monitors entered at the beginning of a section of code are exited at the end. *Unsafe* can violate that contract: it can be used to asymmetrically enter or exit a monitor, and that asymmetry might be not immediately obvious.

```
1 void m() {  
2     unsafe.monitorEnter(o);  
3     if (c) return;  
4     unsafe.monitorExit(o);  
5 }
```

Listing 3.6. `sun.misc.Unsafe` can lead to monitor deadlocks

The examples shown above are just the most straightforward violations of Java's safety guarantees. The `sun.misc.Unsafe` class provides a multitude of methods that can be used to violate most guarantees Java provides.

To sum it up: *Unsafe* is dangerous. But should anybody care? In the next sections we present a study to determine whether and how *Unsafe* is used in real-world third-party Java libraries, and to what degree real-world applications directly and indirectly depend on it.

3.2 Is Unsafe Used?

To answer *RQ/U1* (*To what extent does the Unsafe API impact common application code?*) we need to determine whether and how Unsafe is actually used in real-world third-party Java libraries, and to what degree real-world applications directly and indirectly depend on such unsafe libraries. To achieve our goal, several elements are needed.

Code Repository. As a code base representative of the “real world”, we have chosen the *Maven Central* software repository. The rationale behind this decision is that a large number of well-known Java projects deploy to *Maven Central* using Apache Maven. Besides code written in Java, projects written in Scala are also deployed to *Maven Central* using the Scala Build Tool (sbt). Moreover, *Maven Central* is the largest Java repository,⁹ and it contains projects from the most popular source code management repositories, like *GitHub* and *SourceForge*.

Artifacts. In Maven, an artifact is the output of the build procedure of a project. An artifact can be any type of file, ranging from a *.pdf* to a *.zip* file. However, Artifacts are usually *.jar* files, which archive compiled Java bytecode stored in *.class* files.

Bytecode Analysis. We examine these kinds of artifacts to analyse how they use `sun.misc.Unsafe`. We use a bytecode analysis library to search for method call sites and field accesses of the `sun.misc.Unsafe` class [Bruneton et al., 2002a,b; Kuleshov, 2007].

However, our first attempt was to use JNIF, our own bytecode analysis library [Mastrangelo and Hauswirth, 2014]. JNIF is described in Appendix B. Due to its own limitations, we decided to use the aforementioned analysis library.

Dependency Analysis. We define the impact of an artifact as how many artifacts depend on it, either directly or indirectly. This helps us to define the impact of artifacts that use `sun.misc.Unsafe`, and thus the impact `sun.misc.Unsafe` has on real-world code overall.

Usage Pattern Detection. After all call sites and field accesses are found, we analyse this information to discover usage patterns. It is common that an artifact exhibits more than one pattern. Our list of patterns is not exhaustive. We have manually investigated the source code of the 100 highest-impact artifacts using `sun.misc.Unsafe` to understand why and how they are using it.

⁹<http://www.modulecounts.com/>

Our analysis found 48,490 uses of `sun.misc.Unsafe`—48,139 call sites and 351 field accesses—distributed over 817 different artifacts. This initial result shows that `Unsafe` is indeed used in third-party code.

We use the dependency information to determine the impact of the artifacts that use `sun.misc.Unsafe`. We rank all artifacts according to their impact (the number of artifacts that directly or indirectly depend on them). High-impact artifacts are important; a safety violation in them can affect any artifact that directly or indirectly depends on them. We find that while overall about 1% of artifacts directly use `Unsafe`, for the top-ranked 1000 artifacts, 3% directly use `Unsafe`. Thus, `Unsafe` usage is particularly prevalent in high-impact artifacts, artifacts that can affect many other artifacts.

Moreover, we found that 21,297 artifacts (47% of the 47,127 artifacts with dependency information, or 25% of the 86,479 artifacts we downloaded) directly or indirectly depend on `sun.misc.Unsafe`. Excluding language artifacts, numbers do not change much: Instead of 21,297 artifacts, we found 19,173 artifacts, 41% of the artifacts with dependency information, or 22% of artifacts downloaded. Thus, `sun.misc.Unsafe` usage in third-party code indeed impacts a large fraction of projects.

The complete scripts and results used for this study are available online.¹⁰ Moreover, this study contains a companion artifact that aims to reproduce the results shown in this chapter. The companion artifact can be found in the *Source Materials* tab in the ACM Digital Library.¹¹

Which Features of *Unsafe* Are Actually Used?

Figures 3.1 and 3.2 show all instance methods and static fields of the `sun.misc.Unsafe` class. For each member we show how many call sites or field accesses we found across the artifacts. The class provides 120 public instance methods and 20 public fields (version 1.8 update 40). The figure only shows 93 methods because the 18 methods in the *Heap Get* and *Heap Put* groups, and *staticFieldBase* are overloaded, and we combine overloaded methods into one bar.

We show two columns, *Application* and *Language*. The *Language* column corresponds to language implementation artifacts while the *Application* column corresponds to the rest of the artifacts.

We categorized the members into groups, based on the functionality

¹⁰<https://gitlab.com/acuarica/java-unsafe-analysis>

¹¹<http://dx.doi.org/10.1145/2814270.2814313>

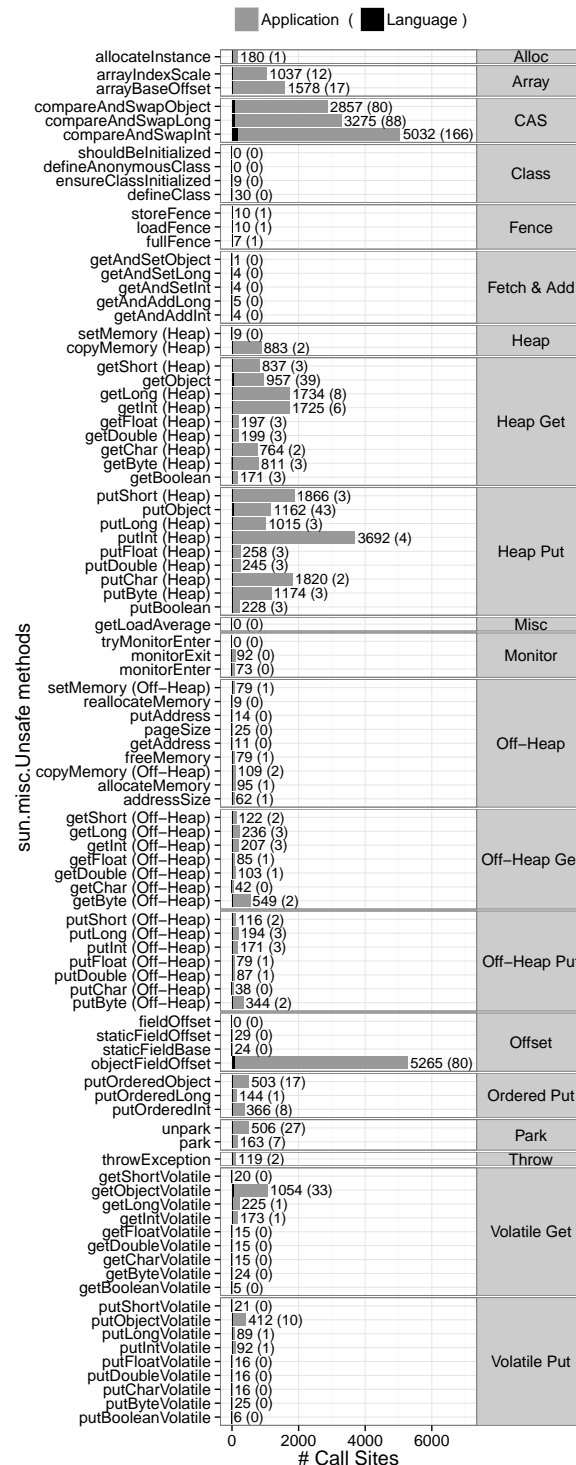


Figure 3.1. sun.misc.Unsafe method usage on Maven Central

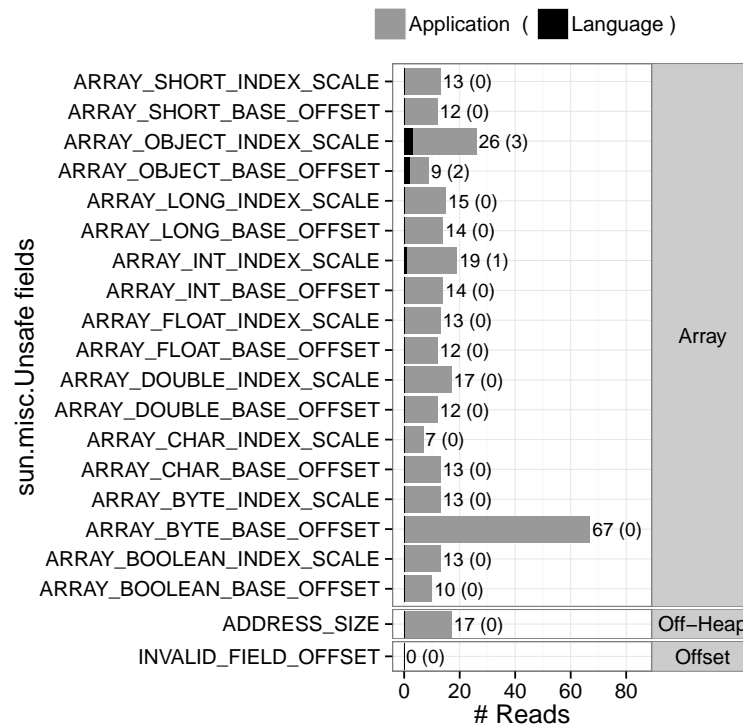


Figure 3.2. `sun.misc.Unsafe` field usage on *Maven Central*

they provide:

- The *Alloc* group contains only the *allocateInstance* method, which allows the developer to allocate a Java object without executing a constructor. This method is used 181 times: 180 in *Application* and 1 in *Language*.
- The *Array* group contains methods and fields for computing relative addresses of array elements. The fields were added as a simpler and potentially faster alternative in a more recent version of *Unsafe*. The value of all fields in this group are constants initialized with the result of a call to either *arrayBaseOffset* or *arrayIndexScale* in the *Array* group. The figures show that the majority of sites still invoke the methods instead of accessing the corresponding constant fields.
- The *CAS* group contains methods to atomically compare-and-swap a Java variable. These methods are implemented using processor-specific atomic instructions. For instance, on *x86* architectures, *com-*

pareAndSwapInt is implemented using the *CMPXCHG* machine instruction. Figure 3.1 shows that these methods represent the most heavily used feature of *Unsafe*.

- Methods of the *Class* group are used to dynamically load and check Java classes. They are rarely used, with *defineClass* being used the most.
- The methods of the *Fence* group provide memory fences to ensure loads and stores are visible to other threads. These methods are implemented using processor-specific instructions. These methods were introduced only recently in Java 8, which explains their limited use in our data set. We expect that their use will increase over time and that other operations, such as those in the *Ordered Put*, or *Volatile Put* groups will decrease as programmers use the lower-level fence operations.
- The *Fetch & Add* group, like the *CAS* group, allows the programmer to atomically update a Java variable. This group of methods was also added recently in Java 8. We expect their use to increase as programmers replace some calls to methods in the *CAS* group with the new functionality.
- The *Heap* group methods are used to directly access memory in the Java heap. The *Heap Get* and *Heap Put* groups allow the developer to load and store a Java variable. These groups are among the most frequently used ones in *Unsafe*.
- The *Misc* group contains the method *getLoadAverage*, to get the load average in the operating system run queue assigned to the available processors. It is not used.
- The *Monitor* group contains methods to explicitly manage Java monitors. The *tryMonitorEnter* method is never used.
- The *Off-Heap* group members provide access to unmanaged memory, enabling explicit memory management. Similarly to the *Heap Get* and *Heap Put* groups, the *Off-Heap Get* and *Off-Heap Put* groups allow the developer to load and store values in Off-Heap memory. The usage of these methods is non-negligible, with *getByte* and *putByte* dominating the rest. The value of the *ADDRESS_SIZE* field is the result of the method *addressSize()*.

- Methods of the *Offset* group are used to compute the location of fields within Java objects. The offsets are used in calls to many other `sun.misc.Unsafe` methods, for instance those in the *Heap Get*, *Heap Put*, and the *CAS* groups. The method *objectFieldOffset* is the most called method in `sun.misc.Unsafe` due to its result being used by many other `sun.misc.Unsafe` methods. The *fieldOffset* method is deprecated, and indeed, we found no uses. The `INVALID_FIELD_OFFSET` field indicates an invalid field offset; it is never used because code using *objectFieldOffset* is not written in a defensive style (given that *Unsafe* is used when performance matters, and extra checks might negatively affect performance).
- The *Ordered Put* group has methods to store to a Java variable without emitting any memory barrier but guaranteeing no reordering across the store.
- The *park* and *unpark* methods are contained in the *Park* group. With them, it is possible to block and unblock a thread's execution.
- The *throwException* method is contained in the *Throw* group, and allows one to throw checked exceptions without declaring them in the `throws` clause.
- Finally, the *Volatile Get* and *Volatile Put* groups allow the developer to store a value in a Java variable with volatile semantics.

It is interesting to note that despite our large corpus of code, there are several *Unsafe* methods that are never actually called. If *Unsafe* is to be used in third-party code, then it might make sense to extract those methods into a separate class to be only used from within the runtime library.

Beyond Maven Central

While *Maven Central* is a large repository, we wanted to check whether our results generalize to other common repositories. Thus, we performed a similar analysis of method usage using the *Boa* Dyer et al. [2013a,b] infrastructure. *Boa* allows the developer to mine ASTs of Java projects in *SourceForge*.

The usage profile of *Unsafe* methods we obtained from *Boa* was similar in shape, but at a different scale, compared to the one obtained from

Maven Central. Using *Boa*'s *SourceForge* dataset, for instance, the most called method, `objectFieldOffset`, is called 200 times in 50 projects. This is two orders of magnitude lower than the count we found on *Maven Central*. Although *Boa* enables the mining of source code in a convenient way, the data it analyses probably is not current enough to include the more recent Java code that uses `sun.misc.Unsafe` more heavily.

In recent versions [Dyer et al., 2015], *Boa* added support to conduct empirical studies on open source projects from *GitHub* and *Qualitas Corpus* [Tempero et al., 2010]. However, at the time we conducted our study on *Unsafe*, this support was not included yet.

3.3 Finding `sun.misc.Unsafe` Usage Patterns

We examined the artifacts in the *Maven Central* software repository to identify usage patterns for *Unsafe*. This section describes our methodology for identifying these patterns.

Our first step is to visualize how an artifact uses *Unsafe*. To this end, we count the *Unsafe* call sites and field usages per class in each artifact. Figures 3.3 and 3.4 show two examples of call sites usages for *com.lmax:disruptor* and *org.scala-lang:scala-library* respectively. Each row shows a fully qualified class name and their usage of `sun.misc.Unsafe`.

After determining the call sites and field usage per artifact, we tried to find a way to group artifacts by how they use `sun.misc.Unsafe`. The first issue is to determine which method calls work together to achieve a goal. These calls might all be located within a single class, be spread across different classes within a package, or be spread across different packages within the whole artifact. After trying different combinations, we decided to group together calls occurring within a single class and its inner classes.

We cluster classes and their inner classes by *Unsafe* method usage using a dendrogram. Because a dendrogram can result in different clusters depending on at which height the dendrogram is cut, we experimented with various clusterings until settling on 31 clusters. An example of a cluster and its dendrogram is shown in Figure 3.5. In the figure we can see classes using methods of the *Off-Heap*, *Off-Heap Get*, and *Off-Heap Put* groups to implement large arrays.

Once we had a clustering of the artifacts by method usage, we manually inspected a sample of artifacts in each cluster to identify patterns. Some artifacts contained more than one pattern. For instance the cluster

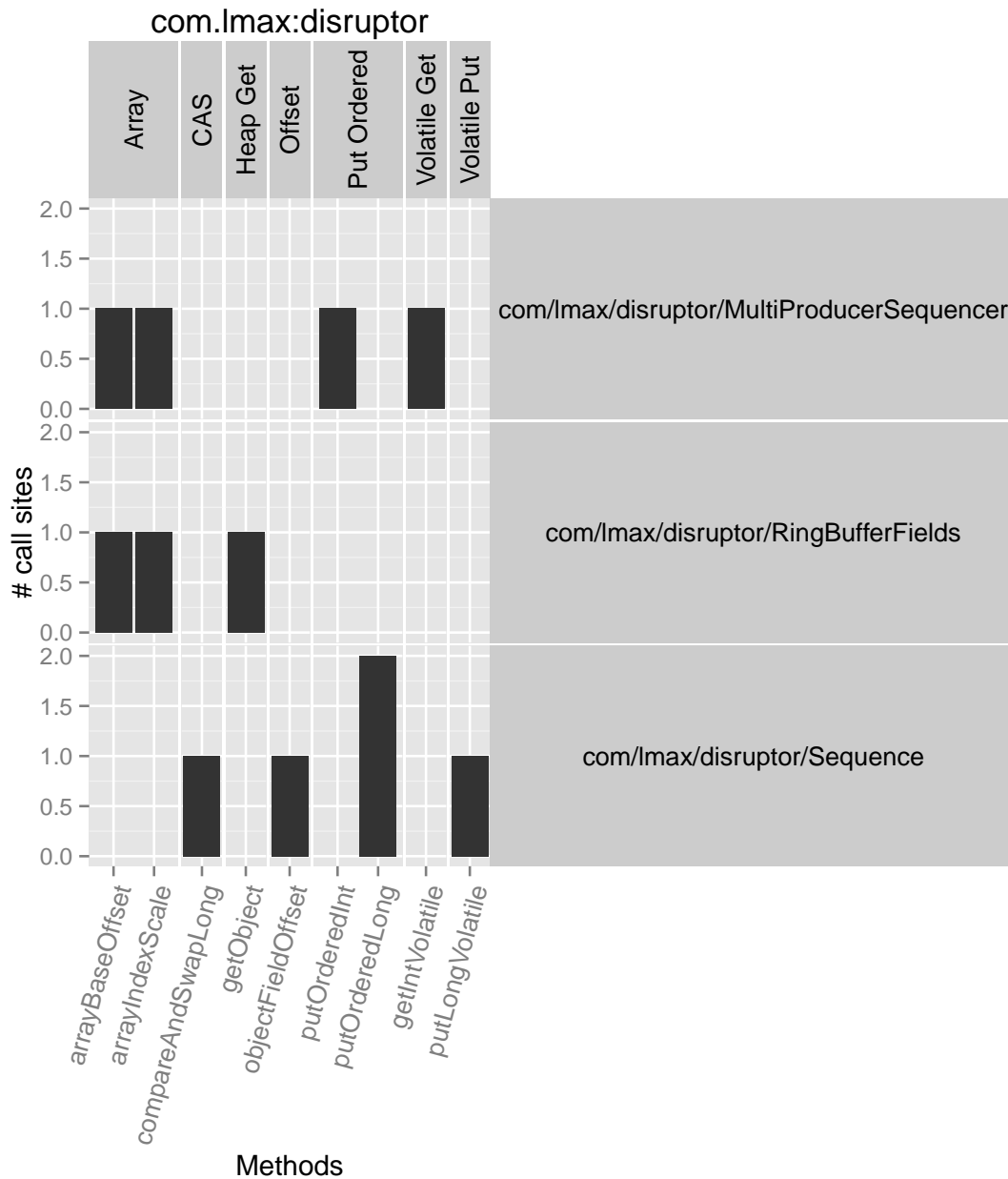


Figure 3.3. *com.lmax:disruptor* call sites

in Figure 3.5 contains classes that use *Unsafe* to implement large off-heap arrays, but also contains calls to methods of the *Put Volatile* group used to implement strongly shared consistent variables. We tagged each artifact manually inspected with the set of patterns that it exhibits.

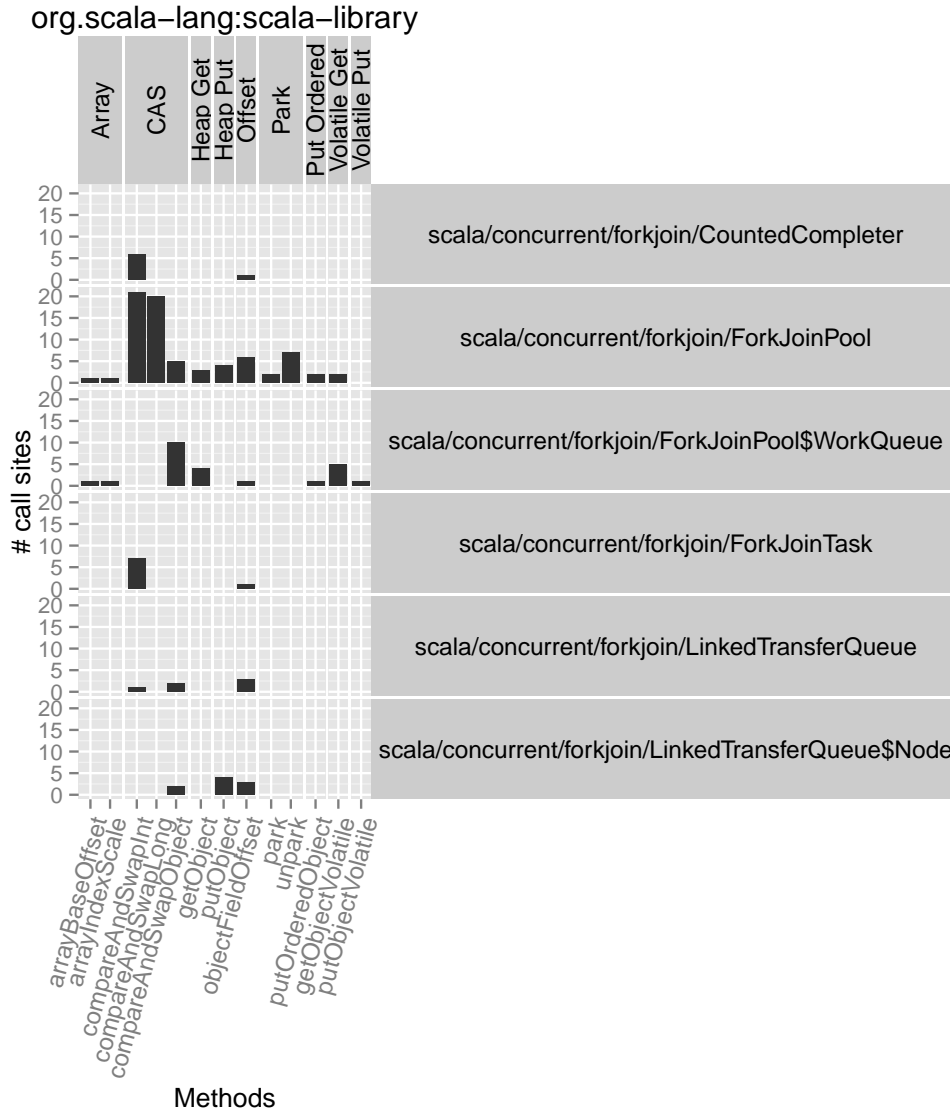


Figure 3.4. `org.scala-lang:scala-library` call sites

3.4 Usage Patterns of `sun.misc.Unsafe`

This section presents the patterns we have found during our study. We present them sorted by how many artifacts depend on them, as computed from the Maven dependency graph described in Section 3.2.

A summary of the patterns is shown in Table 3.1. The *Pattern* column indicates the name of the pattern. *Found in* indicates the number of artifacts in *Maven Central* that contain the pattern. *Used by* indicates the

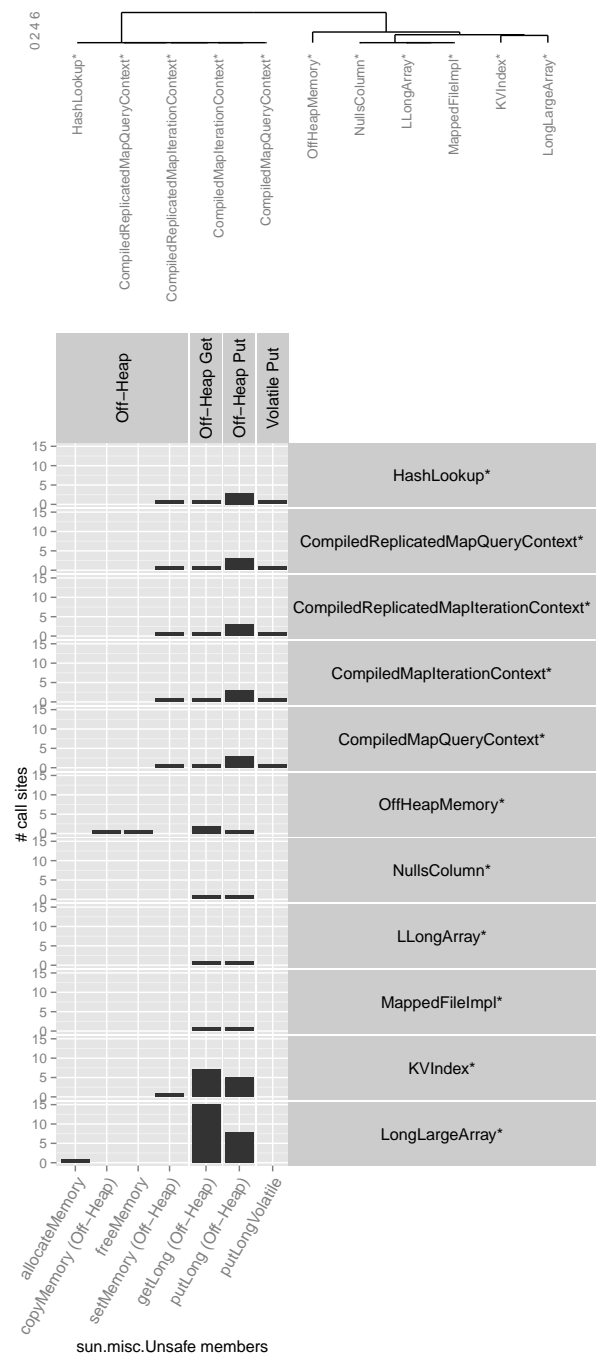


Figure 3.5. Classes using off-heap large arrays

number of artifacts that transitively depend on the artifacts with the pattern. *Most used artifacts* presents the three most used artifacts containing

Table 3.1. Patterns and their occurrences in the Maven Central repository

	Pattern	Found In	Used by	Most used artifacts
1	Allocate an Object without Invoking a Constructor	88	14794	<code>org.springframework:spring-core</code> <code>org.objenesis:objenesis</code> <code>org.mockito:mockito-all</code>
2	Process Byte Arrays in Block	44	12274	<code>com.google.guava:guava</code> <code>com.google.gwt:gwt-dev</code> <code>net.jpountz.lz4:lz4</code>
3	Atomic Operations	84	10259	<code>org.scala-lang:scala-library</code> <code>org.apache.hadoop:hadoop-hdfs</code> <code>org.glassfish.grizzly:grizzly-framework</code>
4	Strongly Consistent Shared Variables	198	9795	<code>org.scala-lang:scala-library</code> <code>org.jruby:jruby-core</code> <code>com.hazelcast:hazelcast-all</code>
5	Park/Unpark Threads	62	7330	<code>org.scala-lang:scala-library</code> <code>org.codehaus.jsr166-mirror:jsr166y</code> <code>com.netflix.servo:servo-internal</code>
6	Update Final Fields	11	7281	<code>org.codehaus.groovy:groovy-all</code> <code>org.jodd:jodd-core</code> <code>com.lmax:disruptor</code>
7	Non-Lexically-Scoped Monitors	14	7015	<code>org.jboss.modules:jboss-modules</code> <code>org.apache.cassandra:cassandra-all</code> <code>org.gridgain:gridgain-core</code>
8	Serialization/Deserialization	32	5689	<code>com.hazelcast:hazelcast-all</code> <code>com.esotericsoftware.kryo:kryo</code> <code>com.thoughtworks.xstream:xstream</code>
9	Foreign Data Access and Object Marshaling	8	3690	<code>eu.stratosphere:stratosphere-core</code> <code>com.github.jnr:jffi</code> <code>org.python:jython</code>
10	Throw Checked Exceptions without Being Declared	59	3566	<code>io.netty:netty-all</code> <code>net.openhft:lang</code> <code>ai.h2o:h2o-core</code>
11	Get the Size of an Object or an Array	4	3003	<code>net.sf.ehcache:ehcache</code> <code>com.github.jbellis:jamm</code> <code>org.openjdk.jol:jol-core</code>
12	Large Arrays and Off-Heap Data Structures	12	487	<code>org.neo4j:neo4j-primitive-collections</code> <code>com.orienttechnologies:orientdb-core</code> <code>org.mapdb:mapdb</code>
13	Get Memory Page Size	11	359	<code>org.apache.hadoop:hadoop-common</code> <code>net.openhft:lang</code> <code>org.xerial.larray:larray-mmap</code>
14	Load Class without Security Checks	21	294	<code>org.elasticsearch:elasticsearch</code> <code>org.apache.geronimo.ext.openejb:core</code> <code>net.openhft:lang</code>

the pattern, that is, the artifact with the most other artifacts that transitively depend upon it. Artifacts are shown using their Maven identifier, *i.e.* `<groupId>:<artifactId>`.

We present each pattern using the following template.

Description. What is the purpose of the pattern? What does it do?

Rationale. What problem is the pattern trying to solve? In what contexts is it used?

Implementation. How is the pattern implemented using `sun.misc.Unsafe`?
Issues. Issues to consider when using the pattern. In addition, we present the problems discussed in the Stack Overflow question/answer database based on our previous work [Mastrangelo et al., 2015].

3.4.1 Allocate an Object without Invoking a Constructor

Description. With this pattern an object can be allocated on the heap without executing its constructor.

Rationale. This pattern is useful for creating mock objects for testing and in deserializing serialized objects.

Implementation. The `allocateInstance` method takes a `java.lang.Class` object as parameter, and returns a new instance of that class. Unlike allocating an object directly, or through the reflection API, the object's constructor is not invoked.

Issues. If the constructor is not invoked, the object might be left uninitialized and its invariants might not hold. Users of `allocateInstance` must take care to properly initialize the object before it is used by other code. This is often done in conjunction with other methods of *Unsafe*, for instance those in the *Heap Put* group, or by using the Java reflection API.

3.4.2 Process Byte Arrays in Block

Description. When processing the elements of a byte array, better performance can be achieved by processing the elements 8 bytes at a time, treating it as a long array, rather than one byte at a time.

Rationale. The pattern is used for fast byte array processing, for instance, when comparing two byte arrays lexicographically.

Implementation. The `arrayBaseOffset` method is invoked to get the base offset of the byte array. Then `getLong` is used to fetch and process 8 bytes of the array at a time.

Issues. The pattern assumes that bytes in an array are stored contiguously. This may not be true for some VMs, e.g. those implementing large arrays using discontinuous arrays or arraylets Siebert [2000]; Bacon et al. [2003]. Users of the pattern should be aware of the endianness of the underlying hardware. In one Stack Overflow discussion, this pattern is discouraged since it is non-portable and, on many JVMs, results in slower code.¹²

¹²<http://stackoverflow.com/questions/12226123>

3.4.3 Atomic Operations

Description. This pattern is used to implement non-blocking concurrent data structures and synchronization primitives. Hardware-specific atomic operations provided by `sun.misc.Unsafe` are used.

Rationale. Non-blocking algorithms often scale better than algorithms that use locking.

Implementation. To get the offset of a Java variable either `objectFieldOffset` or `arrayBaseOffset/arrayIndexScale` can be used. With this offset, the methods from the *CAS* or *Fetch & Add* groups are used to perform atomic operations on the variable. Other methods of *Unsafe* are often used in the implementation of concurrent data structures, including *Volatile Get/Put*, *Ordered Put*, and *Fence* methods.

Issues. Non-blocking algorithms can be difficult to implement correctly. Programmers must understand the Java memory model and how the *Unsafe* methods interact with the memory model.

3.4.4 Strongly Consistent Shared Variables

Description. Because of Java's weak memory model, when implementing concurrent code, it is often necessary to ensure that writes to a shared variable by one thread become visible to other threads, or to prevent reordering of loads and stores. Volatile variables can be used for this purpose, but `sun.misc.Unsafe` can be used instead with better performance. Additionally, because Java does not allow array elements to be declared volatile, there is no possibility other than to use *Unsafe* to ensure visibility of array stores. The methods of the *Ordered Put* groups and the *Volatile Get/Put* groups can be used for these purposes. In addition, the *Fence* methods were introduced in Java 8 expressly to provide greater flexibility for this use case.

Rationale. This pattern is useful for implementing concurrent algorithms or shared variables in concurrent settings. For instance, JRuby uses a *fullFence* to ensure visibility of writes to object fields.

Implementation. To ensure a write is visible to another thread, *Volatile Put* methods or *Ordered Put* methods can be used, even on non-volatile variables. Alternatively, a *storeFence* or *fullFence* can be used. *Volatile Get* methods ensure other loads and stores are not reordered across the load. A *loadFence* could also be used before a read of a shared variable.

Issues. Fences can replace volatile variables in some situations, offering better performance. Most of the uses of the pattern use the *Ordered Put* and

Volatile Put methods. Since they were added to Java only recently, there are currently few instances of the pattern that use the *Fence* methods.

3.4.5 Park/Unpark Threads

Description. To implement locks and other blocking synchronization constructs, the *park* and *unpark* methods are used. With these methods, the developer can block and unblock threads.

Rationale. The alternative to parking a thread is to busy-wait, which uses CPU resources and does not allow other threads to proceed.

Implementation. The *park* method blocks the current thread while *unpark* unblocks a thread given as an argument.

Issues. Users of *park* must be careful to avoid deadlock.

3.4.6 Update Final Fields

Description. This pattern is used to update a final field.

Rationale. Although it is possible to use reflection to implement the same behavior, updating a final field is easier and more efficient using `sun.misc.Unsafe`. Some applications update final fields when cloning objects or when deserializing objects.

Implementation. The *objectFieldOffset* methods and one of the *Put* methods work in conjunction to directly modify the memory where a final field resides.

Issues. There are numerous security and safety issues with modifying final fields. The update should be done only on newly created objects (perhaps also using *allocateInstance* to avoid invoking the constructor) before the object becomes visible to other threads. The Java Language Specification (Section 17.5.3) Gosling et al. [2013] recommends that final fields not be read until all updates are complete. In addition, the language permits compiler optimizations with final fields that can prevent updates to the field from being observed. Since final fields can be cached by other threads, one instance of the pattern uses *putObjectVolatile* to update the field rather than simply *putObject*. Using this method ensures that any cached copy in other threads is invalidated.

3.4.7 Non-Lexically-Scoped Monitors

Description. In this pattern, monitors are explicitly acquired and released without using synchronized blocks.

Rationale. The pattern is used in some situations to avoid deadlock, releasing a monitor temporarily, then reacquiring it.

Implementation. One usage of the pattern is to temporarily release monitor locks acquired in client code (e.g., through a synchronized block or method) and then to reenter the monitor before returning to the client. The `monitorExit` method is used to exit the synchronized block. Because monitors are reentrant, the pattern uses the method `Thread.holdsLock` to implement a loop that repeatedly exits the monitor until the lock is no longer held. When reentering the monitor, `monitorEnter` is called the same number of times as `monitorExit` was called to release the lock.

Issues. Care must be taken to balance calls to `monitorEnter` and `monitorExit`, or else the lock might not be released or an `IllegalMonitorStateException` might be thrown.

3.4.8 Serialization/Deserialization

Description. In this pattern, `sun.misc.Unsafe` is used to persist and subsequently load objects to and from secondary memory dynamically. Serialization in Java is so important that it has a `Serializable` interface to automatically serialize objects that implement it. Although this kind of serialization is easy to use, it does not offer good performance and is inflexible. It is possible to implement serialization using the reflection API. This is also expensive in terms of performance. Therefore, fast serialization frameworks often use `Unsafe` to get and set fields of objects. Some of these projects use reflection to check if `sun.misc.Unsafe` is available, falling back on a slower implementation if not.

Rationale. De/serialization requires reading and writing fields to save and restore objects. Some of these fields may be final or private.

Implementation. Methods of `Heap Get` and `Heap Put` are used to read and write fields and array elements. Deserialization may use `allocateInstance` to create objects without invoking the constructor.

Issues. Using `Unsafe` for serialization and deserialization has many of the same issues as using `Unsafe` for updating final fields (Section 3.4.6) and for creating objects without invoking a constructor (Section 3.4.1). Objects must not escape before being completely deserialized. Type safety can be

violated by using methods of the *Heap Put* group. In addition, care must be taken when deserializing some data structures. For instance, data structures that use `System.identityHashCode` or `Object.hashCode` may need to re-hash objects on deserialization because the deserialized object might have a different hash code than the original serialized object.

3.4.9 Foreign Data Access and Object Marshaling

Description. In this pattern `sun.misc.Unsafe` is used to share data between Java code and code written in another language, usually C or C++.

Rationale. This pattern is needed to efficiently pass data, especially structures and arrays, back and forth between Java and native code. Using this pattern can be more efficient than using native methods and JNI.

Implementation. The methods of the *Off-Heap* group are used to access memory off the Java heap. Often a buffer is allocated using `allocateMemory`, which is then passed to the other language using JNI. Alternatively, the native code can allocate a buffer in a JNI method. The *Off-Heap Get* and *Off-Heap Put* methods are used to access the buffer.

Issues. Use of *Unsafe* here is inherently not type-safe. Care must be taken especially with native pointers, which are represented as long values in Java code.

3.4.10 Throw Checked Exceptions without Being Declared

Description. This pattern allows the programmer to throw checked exceptions without being declared in the method's throws clause.

Rationale. In testing and mocking frameworks, the pattern is used to circumvent declaring the exception to be thrown, which is often unknown. It is used in the Java Fork/Join framework to save the generic exception of a thread to be re-thrown later.

Implementation. This pattern is implemented using the `throwException` method.

Issues. This method can violate Java's subtyping relation, because it is not expected for a method that does not declare an exception to actually throw it. At run time, this can manifest as an uncaught exception.

3.4.11 Get the Size of an Object or an Array

Description. This pattern uses `sun.misc.Unsafe` to estimate the size of an object or an array in memory.

Rationale. The object size can be useful for making manual memory management decisions. For instance, when implementing a cache, object sizes can be used to implement code to limit the cache size.

Implementation. To compute the size of an array, add `arrayBaseOffset` and `arrayIndexScale` (for the given array base type) times the array length. For objects, use `objectFieldOffset` to compute the offset of the last instance field. In both cases, a VM-dependent fudge factor is added to account for the object header and for object alignment and padding.

Issues. Object size is very implementation dependent. Accounting for the object header and alignment requires adding VM-dependent constants for these parameters.

3.4.12 Large Arrays and Off-Heap Data Structures

Description. This pattern uses off-heap memory to create large arrays or data structures with manual memory management.

Rationale. Java's arrays are indexed by `int` and are thus limited to 2^{31} elements. Using *Unsafe*, larger buffers can be allocated outside the heap.

Implementation. A block of memory is allocated with `allocateMemory` and then accessed using *Off-Heap Get* and *Off-Heap Put* methods. The block is freed with `freeMemory`.

Issues. This pattern has all the issues of manual memory management: memory leaks, dangling pointers, double free, etc. One issue, mentioned on Stack Overflow, is that the memory returned by `allocateMemory` is uninitialized and may contain garbage.¹³ Therefore, care must be taken to initialize allocated memory before use. The *Unsafe* method `setMemory` can be used for this purpose.

3.4.13 Get Memory Page Size

Description. `sun.misc.Unsafe` is used to determine the size of a page in memory.

Rationale. The page size is needed to allocate buffers or access memory by page. A common use case is to round up a buffer size, typically a

¹³<http://stackoverflow.com/questions/16723244>

java.nio.ByteBuffer, to the nearest page size. Hadoop uses the page size to track memory usage of cache files mapped directly into memory using *java.nio.MappedByteBuffer*. Another use is to process a buffer page-by-page. Some native libraries require or recommend allocating buffers on page-size boundaries.¹⁴

Implementation. Call *pageSize*.

Issues. Some platforms on which the JVM runs do not have virtual memory, so requesting the page size is non-portable.

3.4.14 Load Class without Security Checks

Description. *sun.misc.Unsafe* is used to load a class from an array containing its bytecode. Unlike with the *ClassLoader* API, security checks are not performed.

Rationale. This pattern is useful for implementing lambdas, dynamic class generation, and dynamic class rewriting. It is also useful in application frameworks that do not interact well with user-defined class loaders.

Implementation. The pattern is implemented using the *defineClass* method, which takes a byte array containing the bytecode of the class to load.

Issues. The pattern violates the Java security model. Untrusted code could be introduced into the same protection domain as trusted code.

3.5 What is the Unsafe API Used for?

In response to *RQ/U2 (How and when are Unsafe features used?)*, many of the patterns we found indicate that *Unsafe* is used to achieve better performance or to implement functionality not otherwise available in the Java language or standard library.

However, many of the patterns described can be implemented using APIs already provided in the Java standard library. In addition, there are several existing proposals to improve the situation with *Unsafe* already under development within the Java community. Oracle software engineer Paul Sandoz [2014] performed a survey on the OpenJDK mailing list to study how *Unsafe* is used¹⁵ and describes several of these proposals.

A summary of the patterns with existing and proposed alternatives to *Unsafe* is shown in Table 3.2. The table consists of the following columns:

¹⁴<http://stackoverflow.com/questions/19047584>

¹⁵<http://www.infoq.com/news/2014/02/Unsafe-Survey>

Table 3.2. Patterns and their alternatives. A bullet (●) indicates that an alternative exists in the Java language or API. A check mark (✓) indicates that there is a proposed alternative for Java.

#	Pattern	Lang	VM	Lib	Ref
1	Allocate an Object without Invoking a Constructor	✓			
2	Process Byte Arrays in Block		✓		
3	Atomic Operations			●	
4	Strongly Consistent Shared Variables			✓	
5	Park/Unpark Threads			●	
6	Update Final Fields				●
7	Non-Lexically-Scoped Monitors	✓			
8	Serialization/Deserialization	✓		●	●
9	Foreign Data Access and Object Marshaling	✓		●	
10	Throw Checked Exceptions without Being Declared	✓			
11	Get the Size of an Object or an Array	✓		✓	
12	Large Arrays and Off-Heap Data Structures	✓		✓	
13	Get Memory Page Size	✓		✓	
14	Load Class without Security Checks	✓		✓	

The *Pattern* column indicates the name of the pattern. The next three columns indicate whether the pattern could be implemented either as a language feature (*Lang*), virtual machine extension (*VM*), or library extension (*Lib*). The *Ref* column indicates that the pattern can be implemented using reflection. A bullet (●) indicates that an alternative exists in the Java language or API. A check mark (✓) indicates that there is a proposed alternative for Java.

Many APIs already exist that provide functionality similar to *Unsafe*. Indeed, these APIs are often implemented using *Unsafe* under the hood, but they are designed to be used safely. They maintain invariants or perform runtime checks to ensure that their use of *Unsafe* is safe. Because of this overhead, using *Unsafe* directly should in principle provide better performance at the cost of safety.

For example, the *java.util.concurrent* package provides classes for safely performing atomic operations on fields and array elements, as well as sev-

eral synchronizer classes. These classes can be used instead of *Unsafe* to implement atomic operations or strongly consistent shared variables. The standard library class *java.util.concurrent.locks.LockSupport* provides *park* and *unpark* methods to be used for implementing locks. These methods are just thin wrappers around the *sun.misc.Unsafe* methods of the same name and could be used to implement the park pattern. Java already supports serialization of objects using the *java.lang.Serializable* and *java.io.ObjectOutputStream* API. The now-deleted JEP 187 Serialization 2.0 proposal¹⁶ ¹⁷ addresses some of the issues with Java serialization.

Because volatile variable accesses compile to code that issues memory fences, strongly consistent variables can be implemented by accessing volatile variables. However, the fences generated for volatile variables may be stronger (and therefore less performant) than are needed for a given application. Indeed, the *Unsafe Put Ordered* and *Fence* methods were likely introduced to improve performance versus volatile variables. The accepted proposal JEP 193 (Variable Handles¹⁸ [Lea, 2014]) introduces *variable handles*, which allow atomic operations on fields and array elements.

Many of the patterns can be implemented using the reflection API, albeit with lower performance than with *Unsafe* [Korland et al., 2010]. For example, reflection can be used for accessing object fields to implement serialization. Similarly, reflection can be used in combination with *java.nio.ByteBuffer* and related classes for data marshaling. The reflection API can also be used to write to final fields. However, this feature of the reflection API makes sense only during deserialization or during object construction and may have unpredictable behavior in other cases.

Writing a final field through reflection may not ensure the write becomes visible to other threads that might have cached the final field, and it may not work correctly at all if the VM performs compiler optimizations such as constant propagation on final fields.

Many patterns use *Unsafe* to use memory more efficiently. Using structs or packed objects can reduce memory overhead by eliminating object headers and other per-object overhead. Java has no native support for structs, but they can be implemented with byte buffers or with JNI.¹⁹

¹⁶<http://mail.openjdk.java.net/pipermail/core-libs-dev/2014-January/024589.html>

¹⁷<http://web.archive.org/web/20140702193924/http://openjdk.java.net/jeps/187>

¹⁸<https://openjdk.java.net/jeps/193>

¹⁹<http://www.oracle.com/technetwork/java/jvms2013sciam-2013525.pdf>

The Arrays 2.0 [Rose, 2012b] and the value types^{20,21}[Rose et al., 2014; Rose, 2012a] proposals address the large arrays pattern. Project Sumatra²² [OpenJDK, 2013] proposes features for accessing GPUs and other accelerators, one of the use cases for foreign data access. Related proposals include JEP 191²³ [Nutter, 2014], which proposes a new foreign function interface for Java, and Project Panama²⁴ [Rose, 2014], which supports native data access from the JVM.

A *sizeof* feature could be introduced into the language or into the standard library. A use case for this feature includes cache management implementations. A higher-level alternative might be to provide an API for memory usage tracking in the JVM. A page size method could be added to the standard library, perhaps in the *java.nio* package, which already includes *MappedByteBuffer* to access memory-mapped storage.

Other patterns may require Java language changes. For instance, the language could be changed to not require methods to declare the exceptions they throw, obviating the need for *Unsafe* in this case. Indeed, there is a long-running debate²⁵ about the software-engineering benefits of checked exceptions. C#, for instance, does not require that exceptions be declared in method signatures at all. One alternative not requiring a language change, proposed in a Stack Overflow discussion, is to use Java generics instead.²⁶ Because of type erasure, a checked exception can be coerced unsafely into an unchecked exception and thrown.

Changing the Java language to support allocation without constructors or non-lexically-scoped monitors is feasible. However, implementation of these features must be done carefully to ensure object invariants are properly maintained. In particular, supporting arbitrary unconstructed objects can require type system changes to prevent usage of the object before initialization [Qi and Myers, 2009]. Limiting the scope of this feature to support deserialization only may be a good compromise and has been suggested in the JEP 187 Serialization 2.0 proposal.

Since *Unsafe* is often used simply for performance reasons, virtual machine optimizations can reduce the need for *Unsafe*. For example, the JVM's

²⁰<https://openjdk.java.net/jeps/169>

²¹<http://cr.openjdk.java.net/~jrose/values/values-0.html>

²²<https://openjdk.java.net/projects/sumatra/>

²³<https://openjdk.java.net/jeps/191>

²⁴<https://cr.openjdk.java.net/~jrose/panama/isthmus-in-the-vm-2014.html>

²⁵<http://www.ibm.com/developerworks/library/j-jtp05254/>

²⁶<http://stackoverflow.com/questions/11410042>

runtime compiler can be extended with optimizations for vectorizing byte array accesses, eliminating the motivation to use *Unsafe* to process byte arrays. Many patterns use *Unsafe* to use memory more efficiently. This could be ameliorated with lower GC overhead. There are proposals for this, for instance JEP 189 Shenandoah: Low Pause GC²⁷ [Christine H. Flood, 2014].

3.6 Conclusions

`sun.misc.Unsafe` is an API that was designed for limited use in system-level runtime library code. The *Unsafe* API is powerful, but dangerous. The improper use of *Unsafe* undermines Java's safety guarantees. We studied to what degree *Unsafe* usage has spread into third-party libraries, to what degree such third-party usage of *Unsafe* can impact existing Java code, and which *Unsafe* API features such third-party libraries actually use.

We thereby provided a basis for evolving the *Unsafe* API, the Java language, and the JVM by eliminating unused or abused unsafe features, and by providing safer alternatives for features that are used in meaningful ways. We hope this will help to make *Unsafe* safer.

²⁷<https://openjdk.java.net/jeps/189>

Chapter 4

Empirical Study on the Cast Operator

A common mechanism for relaxing the static typing constraints in object-oriented languages is *casting*. In programming languages with subtyping—or *subtype polymorphism* [Cardelli and Wegner, 1985]—such as Java, C# or C++, casting allows an expression to be viewed at a different type than the one at which it was defined. Casts are checked dynamically, *i.e.*, at run-time, to ensure that the object being cast is an instance of the desired type.

We aim to understand why developers use casts. Why is the static type system insufficient, requiring an escape hatch into dynamic type checking? Specifically, we attempt to answer the following three research questions:

RQ/C1 : How frequently is casting used in common application code?

To what extent does application code actually use casting operations?

RQ/C2 : How and when casts are used? If casts are used in application code, how and when do developers use them?

RQ/C3 : How recurrent are the patterns for which casts are used? In addition to understand how and when casts are used, we want to measure how often developers need to resort to certain idioms to solve a particular problem.

To answer these research questions, we devise *usage patterns*. Usage patterns are *recurrent programming idioms* used by developers to solve a specific issue. Usage patterns enable the categorization of different kinds of cast usages and thus provide insights into how the language is being used by

developers in real-world applications. Our cast usage patterns can be: (1) a reference for current and future language designers to make more informed decisions about programming languages, *e.g.*, the addition of *smart casts* in Kotlin,¹ (2) a reference for tool builders, *e.g.*, by providing more precise or new refactoring or code smell analyses, (3) a guide for researchers to test new language features, *e.g.*, Winther [2011] or to carry out controlled experiments about programming, *e.g.*, Stuchlik and Hanenberg [2011], and (4) a guide for developers for best or better practices. To answer our research questions, we empirically study how casts are used by developers. The results of this study have been submitted for publication to OOPSLA'19.

Outline

Section 4.1 provides an introduction to casts in Java, while Section 4.2 illustrates the sort of problems developers have when applying casting conversions. In Section 4.3 we introduce the methodology we used to analyse casts and to devise cast usage patterns. Sections 4.4 and 4.5 present the cast usage patterns and answers our research questions. Finally, Section 4.6 discusses the patterns we found, while Section 4.7 concludes.

4.1 Casts in Java

While casts should be familiar to most developers of object-oriented languages, because casts have different semantics in different programming languages, we briefly summarize the meaning of casts in Java and the terminology used in the rest of this chapter.

One common extension of type systems is *subtyping*, usually seen in *object-oriented* programming languages like Java. The subtype mechanism allows the interoperability of two different but related types. As Pierce [2002] states, "[...] S is a subtype of T , written $S <: T$, to mean that any term of type S can safely be used in a context where a term of type T is expected. This view of subtyping is often called the *principle of safe substitution*." Conversely, if S is a subtype of T , we say that T is a supertype of S .

A cast operation, written $(T) e$ in Java consists of a *target type* T and an *operand* e . The operand evaluates to a *source value* which has a run-time

¹<https://kotlinlang.org/docs/reference/typecasts.html#smart-casts>

source type. In Java, a source reference type is always a class type. For a particular cast evaluated at run time, the *source* of the cast is the expression in the program that created the source value. For reference casts, the source is an object allocation. The source may or may not be known statically.

An *upcast* occurs when the cast is from a source reference type S to a target reference type T , where T is a supertype of S . In our terminology, upcasts include identity casts where the target type is the same as the type of the operand. An upcast does not require a run-time check.

A *downcast*, on the other hand, occurs when converting from a source reference type S to a target reference type T , where T is a proper subtype of S . Listing 4.1 shows how to use the cast operator (line 2) to treat a reference (the variable `o`) as a different type (`String`) as it was defined (`Object`).

```
1 Object o = "foo";  
2 String s = (String)o;
```

Listing 4.1. Variable `o` (defined as `Object`) cast to `String`.

In type-safe object-oriented languages, downcasts require a run-time check to ensure that the source value is an instance of the target type. The above snippet is compiled into the Java bytecode shown in listing 4.2. The `aload_1` instruction (line 3) pushes the local variable `o` into the operand stack. The `checkcast` instruction (line 4) then checks at run-time that the top of the stack has the specified type (`java.lang.String` in this example).

```
1 ldc      "foo"  
2 astore_1  
3 aload_1  
4 checkcast java.lang.String  
5 astore_2
```

Bytecode

Listing 4.2. Compiled bytecode to the `checkcast` instruction.

This run-time check can either *succeed* or *fail*. A `ClassCastException` is thrown when a downcast fails. This exception is an unchecked exception, *i.e.*, the programmer is neither required to handle nor to specify the exception in the method signature. Listing 4.3 shows how to detect whether a cast failed by catching this exception.

```
1 try {  
2     Object x = new Integer(0);  
3     System.out.println((String)x);  
4 } catch (ClassCastException e) {  
5     System.out.println("");  
6 }
```

Listing 4.3. Catch `ClassCastException` when a cast fails.

A *guard* is a conditional expression on which a cast, usually a downcast, is control-dependent and that ensures that the cast is evaluated only if it will succeed. Guards are often implemented using the `instanceof` operator, which tests if an expression is an instance of a given reference type. If an `instanceof` guard returns true, the guarded cast should not throw a `ClassCastException`. Listing 4.4 shows a usage of the `instanceof` operator together with a cast expression.

```
1 if (x instanceof Foo) {  
2     ((Foo)x).doFoo();  
3 }
```

Listing 4.4. Runtime type test using `instanceof` before applying a cast.

An object's type can also be checked using reflection: the `getClass` method² returns the run-time class of an object. This `Class` object can be then compared against a class literal, e.g., `x.getClass() == C.class`. This test is more precise than an `x instanceof C` test since it succeeds only when the operand's class is exactly `C`, rather than any subclass of `C`. Listing 4.5 shows how to use the `getClass` method to test for an object's type.

```
1 if (x.getClass() == Foo.class) {  
2     ((Foo)x).doFoo();  
3 }
```

Listing 4.5. Runtime type test using `getClass` before applying a cast.

Because they can fail, downcasts pose potential threats. Unguarded downcasts in particular are worrisome because the developer is essentially

²<https://docs.oracle.com/javase/8/docs/api/java/lang/Object.html#getClass-->

telling the compiler “*Trust me, I know what I’m doing.*” Because downcasts are an escape-hatch from the static type system—they permit dynamic type errors—a cast is often seen as a design flaw or code smell in an object-oriented system [Tufano et al., 2015].

A cast can also fail at compile time if the cast operand and the target type are incompatible. For instance, in the expression `(String) new Integer(1)` a value of type `Integer` can *never* be converted to `String`, so the compiler rejects the cast expression.

Another form of casts in Java are *primitive conversions*, or more specifically *numeric conversions*. These are conversions from one primitive (non-reference) type, usually a numeric type, to another. These conversions can result in loss of precision of the numeric value, although they do not fail with a run-time exception.

Boxing and *unboxing* occur when casting from a primitive type to the corresponding reference type or vice versa, *e.g.*, `(Integer) 3` converts the primitive `int 3` into a boxed `java.lang.Integer`. Unlike downcasts, unboxing casts never throw a `ClassCastException`. However, an unboxing conversion throws a `NullPointerException` when the cast operand is null, *e.g.*, `(double) (Double) null`.³ Java supports *autoboxing* and *autounboxing* between primitives and their corresponding boxed type in the `java.lang` package.

Generics were introduced into Java to provide more static type safety. For instance, the type `List<T>` contains only elements of type `T`. The underlying implementation of generics, however, erases the actual type arguments when compiling to bytecode. To ensure type safety in the generated bytecode, the compiler inserts cast instructions into the generated code. Improper use of generic types or mixing of generic and raw types can lead to dynamic type errors—*i.e.*, `ClassCastException`. Our study, however, does not consider these compiler-inserted casts. Moreover, upcasts inserted by the developer in the source code are completely removed from the generated bytecode by the compiler. Our first attempt to conduct this study was to use our bytecode library analysis [Mastrangelo and Hauswirth, 2014] described in Appendix B. Nevertheless, unlike our *Unsafe* study—which targeted Java bytecode level—in this chapter we are only concerned with programmer-inserted casts in the source code, not in the generated bytecode—given the discrepancy between them.

³<https://docs.oracle.com/javase/specs/jls/se8/html/jls-5.html#jls-5.1.8>

4.2 Issues Developers have using Cast Operations

Do cast operations pose a problem for developers? Several studies [Kechagia and Spinellis, 2014; Coelho et al., 2015; Zhitnitsky, 2016] suggest that in Java, the `ClassCastException` is in the top ten of exceptions being thrown when analysing stack traces. These studies have analysed the exceptions thrown in stack traces. The exceptions come from third-party libraries and the Android API, indicating a misuse of such APIs. `ClassCastException` is in the top 10 of exceptions thrown, thus it represents a problem for developers.

To illustrate the sort of problems developers have when applying casting conversions, we performed a search for commits and issues including the term `ClassCastException` within projects using Java on *GitHub*, the largest host of source code in the world [Gousios et al., 2014]. Our searches returned about 171K commits⁴ and 73K issues,⁵ respectively, at the time of this writing. At first glance, these results indicate that `ClassCastException` indeed represents a source for problems for developers.

Typical classes of bugs encountered when using a cast are using the wrong cast target type, or using the wrong operand, or failing to guard a cast. We present a few examples we found. Each example presented here contains the link to the commit in *GitHub*. Instead of presenting long *GitHub* URLs, we have used the URL shortening service *Bitly* for easier reading. Each *Bitly* link was customized to include the project name.

The following snippet shows a cast applied to the variable `job` (in line 3) that throws `ClassCastException` because the developer forgot to include a guard. In this case, the developer fixed the error by introducing an `instanceof` guard to the cast (lines 1 and 2).

```
1 if(! (job instanceof AbstractProject<?, ?>))
2     return "";
3 AbstractProject<?, ?> project = (AbstractProject<?, ?>) job;
4                                     http://bit.ly/jenkinsci_extra_columns_plugin_2vviBuc
```

Listing 4.6. Cast throws `ClassCastException` because of a forgotten guard.

In the next example the developer made a mistake by choosing a wrong class for the cast target, *i.e.*, `JCustomFileChooser` instead of `CustomFileFilter` (line 9). The `CustomFileFilter` is an inner static class inside the

⁴<https://github.com/search?l=Java&q=ClassCastException&type=Commits>

⁵<https://github.com/search?l=Java&q=ClassCastException&type=Issues>

JCustomFileFilter class. There is no subclass relationship between these two classes. The cast happens inside an equals method—where this idiom is well known—within the CustomFileFilter class. But the developer picked the wrong class, the outer class (JCustomFileFilter), instead of the inner class (CustomFileFilter).

```

1 public final class JCustomFileChooser extends JFileChooser {
2     /* [...] */
3     public static class CustomFileFilter extends FileFilter {
4         /* [...] */
5         public boolean equals(Object obj) {
6             if (getClass() != obj.getClass()) {
7                 return false;
8             }
9             final JCustomFileChooser other = (JCustomFileChooser) obj;
10            if (!Objects.equals(this.extensions, other.extensions)) {
11                return false;
12            }
13        }
14    }
15 }

```

http://bit.ly/GoldenGnu_jeveassets_2vsLbMr

Listing 4.7. Cast throws ClassCastException because of wrong cast target.

More subtle, however, is the interaction between casting and generics. For example, the following call to the getProperty method (line 1), throws a ClassCastException. The method definition is shown in line 3.⁶

```

1 config.getProperty("peer.p2p.pingInterval", 5L)
2
3 public <T> T getProperty(String propName, T defaultValue) {
4     if (!config.hasPath(propName)) return defaultValue;
5     String string = config.getString(propName);
6     if (string.trim().isEmpty()) return defaultValue;
7     return (T) config.getAnyRef(propName);
8 }

```

http://bit.ly/ethereum_ethereumj_2vw4If8

Listing 4.8. Cast throws ClassCastException because of generic inference.

The first argument to the method is the name of a property, used to lookup a value in a table. The second argument is a default value to

⁶http://bit.ly/ethereum_ethereumj_getProperty_2vwQIBH

use if the property is not in the table. If the lookup is successful, the method casts the value found to type `T`. In the call, the given property `"peer.p2p.pingInterval"` is in the table and mapped to an `Integer`. However, Java uses the type of the `defaultValue` argument to instantiate the type parameter `T`. In this case, `Long`—autoboxed from `5L` of type `long`—is used as the type parameter `T`.

Note, however, that the cast inside `getProperty`, which in this context should cast from `Integer` to `Long`, *does not fail*. This is because the Java compiler erases the type parameters like `T` and so dynamic type tests are not performed on them. Instead, the compiler inserts a cast where the return value of `getProperty` is used later with type `Long`. It is this cast that fails at run time and that is reported at run time.

The fix for this bug is to change the default value argument from `5L` to just `5`. This causes the call's return type is inferred to be `Integer`, and the compiler-inserted cast succeeds.

As these examples show, problems with casts are not always obvious. In this thesis we aim to uncover the many different ways in which developers use casts by manually analysing a large sample of cast usages in open source software.

4.3 Finding Cast Usage Patterns

Since casts represent a problem for developers, we aim to provide an answer for our research questions, *RQ/C1 (How frequently is casting used in common application code?)*, *RQ/C2 (How and when casts are used?)* and *RQ/C3 (How recurrent are the patterns for which casts are used?)*. To answer them several elements are needed. We need a corpus of representative “real world” code and we need to perform source code analysis to identify cast operations and to help classify these operations into usage patterns.

Corpus Analysis

We gathered cast usage data using the QL query language, “a declarative, object-oriented logic programming language for querying complex, potentially recursive data structures encoded in a relational data model” [Augustinov et al., 2016]. QL allows us to analyse programs at the source code level. QL extracts the source code of a project into a Datalog model. Besides providing structural data for programs, *i.e.*, ASTs, QL has the ability to

query static types and perform data-flow analysis. To test our QL queries, we have used the *lgtm* service provided by Semmle,⁷ the developers of QL. To gather all cast expressions used in this study, we asked Semmle developers to run a query essentially like Query A.1 on their entire database.

The *lgtm* database includes—at the time of writing—7,559 Java projects imported from open-source projects hosted in *GitHub*. The *lgtm* database was constructed by importing popular open-source projects, *e.g.*, Apache Maven,⁸ Neo4j,⁹ and Hibernate¹⁰. Additionally it includes projects exported by developers to *lgtm* to query them for bug finding, smell detection, and other analyses. We argue that this project selection provides a wide coverage over realistic Java applications, excluding uninteresting projects, *e.g.*, student projects.

Manual Detection of Cast Patterns

We initially sought to describe patterns precisely as QL queries so that detection and categorization was repeatable, but we found this was infeasible because of the complexity of the reasoning involved in identifying some patterns. Often determining to which pattern a cast belongs requires reasoning about the runtime source of the cast, which might be non-local and might depend on external application frameworks or generated code. Thus, we resorted to manually inspect casts in order to devise cast usage patterns.

Nevertheless, whenever possible, we provide a QL query that *approximates* the detection of some patterns. That is, cast expressions returned as the result of a QL query for a pattern, most often belong to that pattern. However, there could be other cast expressions not returned by the query, which are instances of the pattern. Our QL queries used for pattern detection can be found online.¹¹

Unfortunately, we do not possess the *lgtm* database. It is not feasible for us to run our queries on their database. Therefore, it is impractical to gather partial statistics using our queries.

⁷<https://lgtm.com/>

⁸<https://lgtm.com/projects/g/apache/maven>

⁹<https://lgtm.com/projects/g/neo4j/neo4j/>

¹⁰<https://lgtm.com/projects/g/hibernate/hibernate-orm/>

¹¹<https://gitlab.com/acuarica/java-cast-queries>

Methodology

To identify patterns of cast usage, we analysed all Java projects in the *lgtm* database, 7,559 projects with a total 10,193,435 casts, at the time of writing. There are 215 projects in the database for which we could not retrieve the source code. In total, these 215 projects contain 1,162,583 casts. Moreover, there are also 516 projects that do not contain any cast. Therefore the total cast population to be analysed consists of 9,030,852 casts in 6,840 projects.

Because the number of cast instances is large, it is not feasible to *manually* analyse all of them. Therefore we have opted to perform random sampling to get a subset of cast instances to analyse. To choose a sample size such that the probability of missing the least frequent pattern is extremely low, we assume a hypergeometric distribution of the data. The hypergeometric distribution is a discrete probability distribution used with a finite population of N subjects. It is used to calculate the probability of drawing k subjects with a given feature—provided that there are K subjects with that feature in the population—in n draws, without replacement.

Returning to our problem of finding an appropriate sample size, we model our question as follows: We assume there are K casts that are members of the least frequently occurring pattern. We want to know the probability of not finding this pattern, *i.e.*, sampling exactly $k = 0$. Our population consists of $N = 9,030,852$ cast instances. For our study, we assume that a pattern is irrelevant if it represents less than 0.1% of the population, or $K = 9,031$ cast instances. Plugging-in these parameters using the hypergeometric distribution formula,¹² we found that with a sample size of $n = 5,000$ the probability of not sampling the least frequently occurring pattern is 0.67%.

The manual categorization file can be found online.¹³ This file is a comma-separated values (CSV) table. Each row represents a cast instance. This table contains 6 columns. The *castid* and *reloid* columns represent internal IDs to uniquely identify each cast instance and each project. The *target* and *source* columns indicate the source and target types used in the cast. The last two columns—*link* and *value*—are the link to the source code file in *lgtm* and the result of the manual inspection. The script to process the results of the manual inspection is available online as well.¹⁴ We had

¹²The reader can use any hypergeometric distribution calculator, *e.g.*, <https://keisan.casio.com/exec/system/1180573201>

¹³<https://gitlab.com/acuarica/phd-thesis/blob/master/analysis/casts.csv>

¹⁴<https://gitlab.com/acuarica/phd-thesis/blob/master/analysis/analysis.r>

to sample more than 5,000 casts. The CSV table mentioned above contains 5,530 casts (rows). This is because we found 526 links that were not accessible during our analysis, making manual code inspection impossible. Inaccessible links can be found because some projects were removed from the *lgtm* platform. We also found 1 cast that was clearly a bug, a downcast using the wrong cast operand. Thus, we had to resample the cast instances until we reach 5,000 manually inspected casts. When resampling, we took care of inspecting *different* cast instances, *i.e.*, we have discarded duplicated casts. We found 3 duplicated casts when resampling.

4.4 Overview of the Sampled Casts

The casts we sampled are summarized in Table 4.1. Our sample of casts spans 1,299 different projects (19%, out of 6,840 projects). In our sample of 5,000 casts, we found 1,043 (20.86%) primitive conversions. The remaining 3,957 (79.14%) casts are either reference upcasts, downcasts, boxing casts, or unboxing casts.

Table 4.1. Statistics on Sampled Casts

All sampled casts	5,000	100%
Reference casts	3,957	79.14%
Primitive casts	1,043	20.86%
Upcasts	106	2.12%
Downcasts	3,851	77.02%
Boxing casts	11	0.22%
Unboxing casts	18	0.36%
Guarded by instanceof	881	17.62%
Guarded by getClass	64	1.28%
Guarded by type tag	237	4.74%
Unguarded or possibly unguarded	2,499	49.98%
In application/library code	3,436	68.72%
In test code	560	11.20%
In generated code	1,058	21.16%

Castes can be classified as either *guarded* or *unguarded* casts. A guard is a conditional expression on which the cast is control dependent, which, if successful, ensures the cast will not fail. Guards are typically implemented using the `instanceof` operator or using a test of the source value's class (retrieved using the `Object::getClass` method) against a subtype of the cast target type. Guards can also be implemented in an application-specific

manner, for instance by associating a “type tag” with the source value that can be used to distinguish the run-time type.

Of the 3,957 analysed reference casts, we found that 1,458 (29.16%) were guarded by a guard in the same method as the cast and 2,499 (49.98%) were either unguarded or had a guard in another method. In the latter case, which we refer to as *possibly unguarded*, determining by manual inspection if a guard is actually present is often infeasible. The possibly unguarded casts are cases where the application developer has some reason for believing the cast will succeed, but it is not immediately apparent in the source code.

As with any expression, casts can appear in either application/library code, test code, or generated code. As expected, most casts appear in application or library code (68.72%). However, casts in test and generated code are not negligible (11.20% and 21.16% respectively).

As we describe in the next section, nearly all guarded casts fit into just a few patterns. Unguarded or possibly unguarded casts account for most of the patterns.

4.5 Cast Usage Patterns

Using the methodology described in the above section, we have devised 25 cast usage patterns. Table 4.2 presents our patterns and their occurrences sorted by frequency.

The patterns were arrived at by an iterative process. Each sampled cast was assigned a pattern. If no pattern fit the given cast, a new pattern was invented and described. My advisors and I then discussed the patterns and their instances, refining, merging, or splitting them into new patterns. This process was repeated until consensus among us was reached. The particular categorization here is therefore subjective.

We do not claim that our list of patterns is exhaustive, although our methodology should ensure that any pattern that occurs more than 0.1% of the time has a small probability of being excluded.

Moreover, we are interested in the scope of the cast instance, *i.e.*, *does it appear in application/library code, test code, or generated code?* Figure 4.1 shows our patterns and their occurrences grouped by scope and sorted by frequency.

Each pattern is described using the following template:

- *Description.* Tells what the pattern is about, gives a general overview of its structure, and briefly describes the rationale behind how this

Table 4.2. Cast Usage Patterns

Pattern	Description	# Casts	%
TYPECASE	A cast is guarded with an instance of expression, class literal, or application-specific tag.	1,182	23.64%
STASH	A cast to an heterogeneous collection element.	561	11.22%
FACTORY	A cast used to convert a newly created objects.	381	7.62%
FAMILY	A cast applied in a family of classes.	344	6.88%
USERAWTYPE	A cast used instead of the declared generic type.	335	6.70%
EQUALS	A cast used in the implementation of the well-known equals method.	247	4.94%
REDUNDANT	A cast that is not necessary for compilation.	122	2.44%
COVARIANTRETURNTYPE	A cast when the return type of a method is covariant.	106	2.12%
SELECTOVERLOAD	A cast to disambiguate between overloaded methods.	99	1.98%
KNOWNRETURNTYPE	The client of an API knows the exact return type of a method invocation.	89	1.78%
DESERIALIZATION	A cast used to convert newly created objects in deserialization.	71	1.42%
VARIABLESUPERTYPE	A cast to a variable that could be declared to be more specific.	64	1.28%
SOLESUBCLASSIMPLEMENTATION	A cast to the only subclass implementation.	61	1.22%
NEWDYNAMICINSTANCE	Cast the result of newInstance in Class, Constructor, or Array.	59	1.18%
OBJECTASARRAY	A cast to a constant array slot used as a field of an object.	47	0.94%
IMPLICITINTERSECTIONTYPE	A cast to implicitly use an intersection type.	45	0.90%
REMOVEWILDCARD	A cast used instead of the declared generic type.	34	0.68%
OPERANDSTACK	A cast to an heterogeneous stack.	29	0.58%
REFLECTIVEACCESSIBILITY	Cast the result of the Method::invoke, or Field::get.	27	0.54%
FLUENTAPI	Cast to permit a fluent API through method chaining.	23	0.46%
COVARIANTGENERIC	Remove type parameter or an upcast to permit covariant generics.	22	0.44%
COMPOSITE	A composite cast.	21	0.42%
GENERICARRAY	A cast to create a generic array.	7	0.14%
ACCESSSUPERCLASSFIELD	A cast to access a private field in a superclass.	4	0.08%
UNOCCUPIEDTYPEPARAMETER	A cast to a raw type or to remove the wildcard in a generic type.	1	0.02%

pattern was characterized as such. A few patterns can have distinct *variants*, *i.e.*, different ways of implementing the pattern. Whenever a pattern has variants, we state how they differ from each other.

- *Instances.* Gives one or more concrete examples found in real code. The code snippets presented here were modified for formatting purposes. Each example contains a highlighted line which shows the cast instance being inspected. Moreover, to facilitate some snippet presentations, we remove irrelevant code and replace it with the comment `// [...] or /* [...] */` whenever convenient. For each instance presented here, we provide the link to the source code repository in *lgtm*. We provide the link in case the reader wants to do further in-

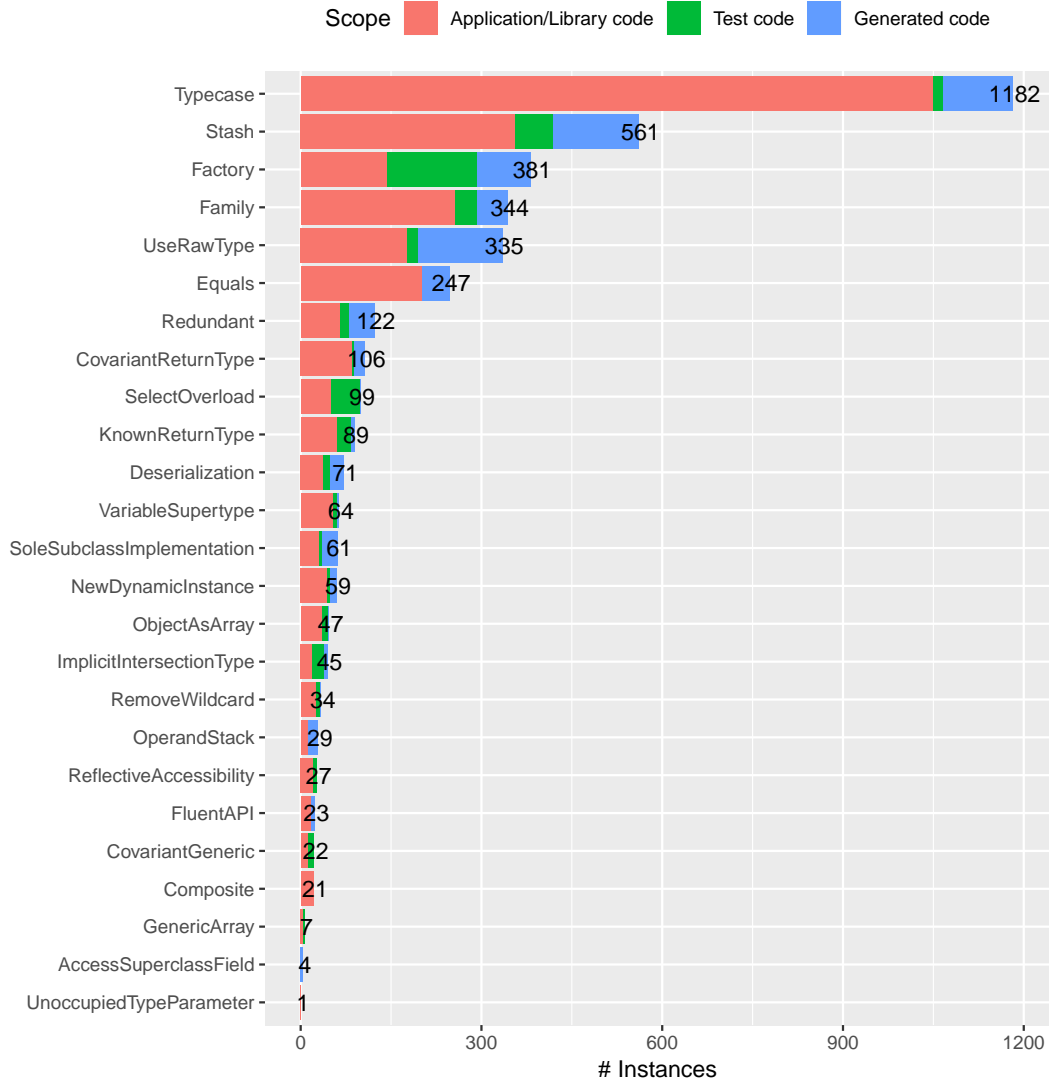


Figure 4.1. Cast Usage Pattern Occurrences

spection of the presented snippet. Instead of presenting long *lgtm* URLs, we have used the URL shortening service *Bitly* for easier reading. Each *Bitly* link was customized to include the project name. As we mentioned above, projects can be removed from the *lgtm* service, thus some links may not work.

- *Detection.* For some patterns, we provide a QL query that approximates their automatic detection, as describes in Section 4.3. Whenever

a pattern is too complex to describe in terms of QL, we explain the reasons why is so. Our additional QL classes and predicates used in detection queries can be found in Appendix A.

- *Issues.* Discusses the issues with the pattern, flaws, and alternatives that achieve the same goal without casting.

4.5.1 Typecase

Description. The TYPECASE pattern consists of dispatching to different cases depending on the run-time type of the source value. The run-time type is tested against known subtypes of the operand type, with each test followed by a cast to that type. The guard may be implemented using one of three variants: an instanceof operator (*GuardByInstanceOf*), a comparison of the runtime class against a class literal (*GuardByClassLiteral*), or an application-specific type tag (*GuardByTypeTag*).

Instances: 1,182 (23.64%). We found 1,050 in application code, 17 in test code, and 115 in generated code. TYPECASE is by far the most common pattern. Figure 4.2 shows the different variants of the pattern. The *GuardByInstanceOf* is the most used variant. Often there is just one case and the default case, *i.e.*, when the guard fails, performs a no-op or reports an error.

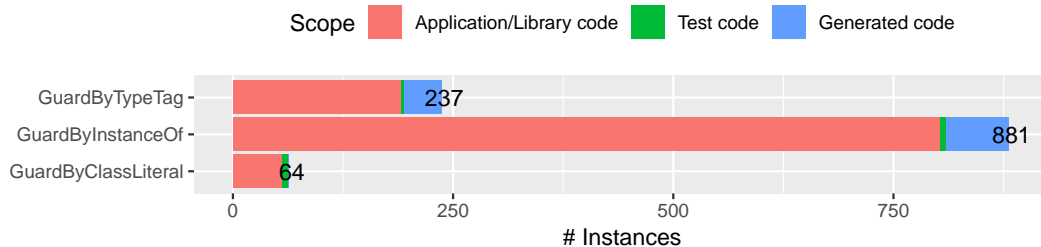


Figure 4.2. TYPECASE Variant Occurrences

The following listing shows an example of the TYPECASE pattern, using the *GuardByInstanceOf* variant.

```

1  if (object instanceof Item) {
2      return getStringFromStack(new ItemStack((Item) object));
3  } else if (object instanceof Block) {
4      return getStringFromStack(new ItemStack((Block) object));
5  } else if (object instanceof ItemStack) {

```

```

6     return getStringFromStack((ItemStack) object);
7 } else if (object instanceof String) {
8     return (String) object;
9 } else if (object instanceof List) {
10    return getStringFromStack((ItemStack) ((List) object).get(0));
11 } else return "";

```

http://bit.ly/PenguinSquad_Enchiridion_2HnNwB7

In the next case a type test is performed—through a method call—before actually applying the cast to the variable props (line 3). Note that the type test is internally using the instanceof operator (line 8).

```

1  @Override
2  public CTSolidColorFillProperties getSolidFill() {
3      return isSetSolidFill() ? (CTSolidColorFillProperties) props : null;
4  }
5  @Override
6  public boolean isSetSolidFill() {
7      return (props instanceof CTSolidColorFillProperties);
8  }

```

http://bit.ly/apache_poi_2FW5SXU

Another common scenario is when several cases are used to re-throw an exception of the right type, as shown below. The cast instance is applied to a variable of type Throwable (line 13). Nevertheless, the enclosing method is only allowed to throw NamingException by the throws declaration (line 3). Since an exception of type Throwable is checked, a cast to VirtualMachineError (subclass of Error) is needed.

```

1  protected Object wrapDataSource(
2      Object datasource, String username, String password)
3      throws NamingException {
4      try {
5          // [...]
6      } catch (Exception x) {
7          if (x instanceof InvocationTargetException) {
8              Throwable cause = x.getCause();
9              if (cause instanceof ThreadDeath) {
10                 throw (ThreadDeath) cause;
11             }
12             if (cause instanceof VirtualMachineError) {
13                 throw (VirtualMachineError) cause;
14             }
15             if (cause instanceof Exception) {
16                 x = (Exception) cause;
17             }
18         }
19     }

```

```

19         if (x instanceof NamingException) throw (NamingException)x;
20         else {
21             // [...]
22         }
23     }
24 }

```

http://bit.ly/codefollower_Tomcat_Research_2SGDUG5

The next example shows that TYPECASE can also be used to filter elements by type within a stream. The cast is applied to stream operations (line 1) over the `caseAssignments` collection. The `instanceof` guard is tested in line 2.

```

1 user = (User) caseAssignments
2     .stream().filter(oe -> oe instanceof User)
3     .findFirst()
4     .orElseThrow(() -> new IllegalArgumentException());
5

```

http://bit.ly/kiegroup_jbpm_2ENCL8a

Rather than using an `instanceof` guard, in the following example the target type of the parameter reference is determined by the value of the parameter `referenceType`, which acts as a *type tag* for reference.

```

1 switch (referenceType) {
2     case ReferenceType.FIELD:
3         return fieldSection.getItemIndex((FieldRefKey) reference);
4     case ReferenceType.METHOD:
5         return methodSection.getItemIndex((MethodRefKey) reference);
6     case ReferenceType.STRING:
7         return stringSection.getItemIndex((StringRef) reference);
8     case ReferenceType.TYPE:
9         return typeSection.getItemIndex((TypeRef) reference);
10    case ReferenceType.METHOD_PROTO:
11        return protoSection.getItemIndex((ProtoRefKey) reference);
12    default:
13        throw new ExceptionWithContext("/ * [...] */", referenceType);
14 }

```

http://bit.ly/JesusFreke_smali_2Ho8bVL

In some cases, the target types of the casts are the same in every branch. In the following snippet, the cast is applied to the `message.obj` field to (line 11), according to the value of the tag `message.what` field (line 1). However, a similar cast is applied in the first branch (line 3). In both branches `message.obj` is of type `Object[]`, but with different lengths. The casts in the calls to `onSuccess` and `onFailure` (lines 5, 13–14) are instances of the `OBJECTASARRAY` pattern.

```

1  switch (message.what) {
2      case SUCCESS_MESSAGE:
3          response = (Object[]) message.obj;
4          if (response != null && response.length >= 3) {
5              onSuccess((Integer) response[0], (Header[]) response[1],
6                      (byte[]) response[2]);
7          } else { /* [...] */ }
8          break;
9      case FAILURE_MESSAGE:
10         response = (Object[]) message.obj;
11         if (response != null && response.length >= 4) {
12             onFailure((Integer) response[0], (Header[]) response[1],
13                     (byte[]) response[2], (Throwable) response[3]);
14         } else { /* [...] */ }
15         break;
16         // [...]
17     }

```

http://bit.ly/loopj_android_async_http_2IpIULk

In the next example, instead of a switch, an if statement is used to guard the cast (in line 6).

```

1  for (final IEnrolment enrolment : dismissal.getSourceIEnrolments()) {
2      if (enrolment.isExternalEnrolment()) {
3          generateExternalEnrolmentRow(mainTable, (ExternalEnrolment) enrolment,
4                                      level + 1, true);
5      } else {
6          generateEnrolmentRow(mainTable, (Enrolment) enrolment,
7                              level + 1, false, true, true);
8      }
9  }

```

http://bit.ly/FenixEdu_fenixedu_academic_2SUNOUJ

In the next example, the parameter args is cast to Object[] (line 13). The “type tag” is given by the fact that the cast is executed in a catch block, and that value is an instance of Closure (line 9). The args parameter flows into two methods, invokeMethod(String name, Object args) and call(Object... args). Thus, args is treated as an Object or Object[] depending on the type tag, resembling a union type.

```

1  public Object invokeMethod(String name, Object args) {
2      try {
3          return super.invokeMethod(name, args);
4      }
5      catch (GroovyRuntimeException e) {
6          // br should get a "native" property match first.
7          // getProperty includes such fall-back logic
8          Object value = this.getProperty(name);

```

```

9         if (value instanceof Closure) {
10             Closure closure = (Closure) value;
11             closure = (Closure) closure.clone();
12             closure.setDelegate(this);
13             return closure.call((Object[]) args);
14         } else {
15             throw e;
16         }
17     }
18 }

```

http://bit.ly/groovy_groovy_core_2SGzK16

In the *GuardByClassLiteral* variant, a cast uses an application-specific guard, but the guard depends on a class literal. In the following example, a cast is performed to the field variable (line 22), based whether the runtime class of the variable is actually `Short.class`.

```

1 Class type = field.getClass();
2 if (type == String.class) {
3     out.writeByte((byte) 1);
4     out.writeString((String) field);
5 } else if (type == Integer.class) {
6     out.writeByte((byte) 2);
7     out.writeInt((Integer) field);
8 } else if (type == Long.class) {
9     out.writeByte((byte) 3);
10    out.writeLong((Long) field);
11 } else if (type == Float.class) {
12    out.writeByte((byte) 4);
13    out.writeFloat((Float) field);
14 } else if (type == Double.class) {
15    out.writeByte((byte) 5);
16    out.writeDouble((Double) field);
17 } else if (type == Byte.class) {
18    out.writeByte((byte) 6);
19    out.writeByte((Byte) field);
20 } else if (type == Short.class) {
21    out.writeByte((byte) 7);
22    out.writeShort((Short) field);
23 } else if (type == Boolean.class) {
24    out.writeByte((byte) 8);
25    out.writeBoolean((Boolean) field);
26 } else if (type == BytesRef.class) {
27    out.writeByte((byte) 9);
28    out.writeBytesRef((BytesRef) field);
29 } else {
30    throw new IOException("Can't handle sort field value of type [" + type + "]");
31 }

```

http://bit.ly/elastic_elasticsearch_2SSgsFV

Similar to the previous example, the next snippet contains several type cases. Each type case is guarded by an equals comparison between a class literal and the `clazz` parameter. The cast is applied to the type parameter `T` only if the guard succeeds.

```

1  @Override
2  @SuppressWarnings("unchecked")
3  public <T> T get(String fieldName, Class<T> clazz) throws DecodingException {
4      if (clazz.equals(Boolean.class)) {
5          return (T) getBoolean(fieldName);
6      }
7      // [...]
8      if (clazz.equals(ExtensionObject.class)) {
9          return (T) getExtensionObject(fieldName);
10     }
11     // [...]
12 }

```

http://bit.ly/OPCFoundation_UA_Java_Legacy_2Fb2xmZ

In the following listing, a cast is applied to the result of the `getObject` method (line 2). The target type of the cast, `MyKey`, corresponds to the class literal argument, `MyKey.class`. Essentially, `getObject` is using the `isInstance` method¹⁵ of the class `java.lang.Class` to check whether an object is from a certain type.

```

1  public MyKey getMyKey() {
2      return (MyKey) getObject(MyKey.class, KEY_MY_KEY);
3  }

```

http://bit.ly/smartdevicelink_sdl_android_2EjJiaq

The following snippet shows an instance of the *GuardByClassLiteral* variant. In this case, the cast is guaranteed to succeed because the class literal used as argument to the recursive call (`Integer.class`) determines that the method returns an `int` value.

```

1  public Object convertToNumber(Number value, Class toType) throws Exception {
2      toType = unwrap(toType);
3      if (AtomicInteger.class == toType) {
4          return new AtomicInteger((Integer)convertToNumber(value, Integer.class));
5      } else if (AtomicLong.class == toType) {
6          return new AtomicLong((Long) convertToNumber(value, Long.class));
7      } else if (Integer.class == toType) {
8          return value.intValue();
9      } else if (Short.class == toType) {

```

¹⁵<https://docs.oracle.com/javase/8/docs/api/java/lang/Class.html#isInstance-java.lang.Object->

```
10     return value.shortValue();
11 } else if (Long.class == toType) {
12     return value.longValue();
13 } else if (Float.class == toType) {
14     return value.floatValue();
15 } else if (Double.class == toType) {
16     return value.doubleValue();
17 } else if (Byte.class == toType) {
18     return value.byteValue();
19 } else if (BigInteger.class == toType) {
20     return new BigInteger(value.toString());
21 } else if (BigDecimal.class == toType) {
22     return new BigDecimal(value.toString());
23 } else {
24     throw new Exception("Unable to convert number "+value+" to "+toType);
25 }
26 }
```

http://bit.ly/apache_karaf_2HE55gE

Detection. When implementing the pattern, care must be taken with complex operands that the value of the operand is not changed between the guard and the cast, possibly even by another thread. For instance, in some situations the operand expression is a method invocation. The value returned by the method should be the same for both the instanceof and the cast, thus the method should be a pure method. Typically, this problem is avoided by using an effectively final local variable in both the guard and the cast operand.

The Query 4.9 detects the *GuardByInstanceOf* variant. It is decoupled in two QL classes. The *ControlByInstanceOfCast* class checks that the cast—to a variable—is control-dependant on an instanceof on the same variable. Then, the *GuardByInstanceOfCast* class checks that the value tested by the instanceof is the same to be cast. That is, it checks that there is no assignment to the variable between the instanceof and the cast.

```

1  class ControlByInstanceOfCast extends VarCast {
2      InstanceOfExpr iof;
3      private ConditionBlock cb;
4      ControlByInstanceOfCast() {
5          iof = cb.getCondition() and
6          cb.controls(getBasicBlock(), true) and
7          var.getAnAccess() = iof.getExpr()
8      }
9      InstanceOfExpr getIof() { result = iof }
10 }
11
12 class GuardByInstanceOfCast extends ControlByInstanceOfCast {
13     GuardByInstanceOfCast() {
14         forall (VariableUpdate def | defUsePair(def, getExpr()) |
15             defUsePair(def, iof.getExpr())
16         )
17     }
18 }

```

Listing 4.9. Query for the *GuardByInstanceOf* variant.

The implementation of *GuardByTypeTag* variant is application-specific, and thus its automatic detection in QL is impractical. Nevertheless, the Query 4.10 detect the special case when a cast is applied to a field in an object inside a switch statement. The expression to be switched is another field in the same object.

```

1  class SwitchFieldTypeTagCast extends Cast {
2      FieldAccess tagAccess;
3      FieldAccess castAccess;
4      Variable v;
5      SwitchFieldTypeTagCast() {
6          tagAccess = this.(SwitchedExpr).getSwitchStmt().getExpr() and
7          castAccess = getExprOrDef() and
8          v.getAnAccess() = tagAccess.getQualifier() and
9          v.getAnAccess() = castAccess.getQualifier()
10     }
11 }

```

Listing 4.10. Detection of a cast inside a switch statement

Similar to the previous case, the Query 4.11 detects when a cast is guarded by a call to the `Class.isArray` method. This query detects *only*

the case when the variable to be cast and the getClass invocation are in the same method.

```

1  class ControlByIsArrayCast extends VarCast {
2      ConditionBlock cb;
3      MethodAccess iama;
4      ControlByIsArrayCast() {
5          exists (VariableAssign def, GetClassMethodAccess gcls |
6              gcls.getQualifier() = var.getAnAccess() and
7              def.getSource() = gcls and
8              defUsePair(def, iama.getQualifier().(VarAccess) )
9          ) and
10         iama.getMethod() instanceof IsArrayClassMethod and
11         (
12             (cb.getCondition()=iama and cb.controls(getBasicBlock(), true)) or
13             (cb.getCondition().(LogNotExpr).getExpr() = iama and
14                 cb.controls(getBasicBlock(), false)
15         )
16     )
17 }
18 }

```

Listing 4.11. Detection of a cast guarded by the Class.isArray method.

The following query detects the *GuardByClassLiteral* variant. Similar to the previous case, this query *does not* detect the case when the variable to be cast and the Class object are passed as parameters. To detect such case would require an inter-procedural (global) data flow analysis. Such analysis does not scale easily.

```

1  class GuardByClassLiteral extends VarCast {
2      TypeLiteral tl;
3      GetClassMethodAccess gcma;
4      GuardByClassLiteral() {
5          gcma.getQualifier() = getVar().getAnAccess() and
6          isSubtype(tl.getTypeName().getType(), getTargetType()) and (
7              controlByEqualityTest(tl, gcma, this) or
8              controlByEqualsMethod(tl, gcma, this)
9          )
10     }
11 }

```

Listing 4.12. Query for the *GuardByClassLiteral* variant.

Issues. Having only a single case—that is, a single guard and cast—is common. In the 742 instances of TYPECASE that used instanceof, 511 (69%) had only one case.

The TYPECASE pattern can be seen as an *ad-hoc* alternative to a typecase or pattern matching [Milner, 1984] as a language construct. In Kotlin, flow-sensitive typing is used so that immutable values can be used at a subtype when a type guard on the value is successful.¹⁶ This feature eliminates much of the need for the guarded casts. Pattern matching can be seen in several other languages, *e.g.*, SML, Scala, C#, and Haskell. For instance, in Scala the pattern matching construct is achieved using the match keyword. In this example,¹⁷ a different action is taken according to the runtime type of the parameter notification (line 10).

```
1  abstract class Notification
2  case class Email(sender: String, title: String, body: String)
3      extends Notification
4  case class SMS(caller: String, message: String)
5      extends Notification
6  case class VoiceRecording(contactName: String, link: String)
7      extends Notification
8
9  def showNotification(notification: Notification): String = {
10     notification match {
11         case Email(email, title, _) => s"Email from $email titled: $title"
12         case SMS(number, message) => s"SMS from $number! Message: $message"
13         case VoiceRecording(name, link) => s"Recording from $name! Link: $link"
14     }
15 }
```

Scala

Alternatives to the TYPECASE pattern would be to use the visitor pattern or to use virtual dispatch on the match scrutinee. However, both of these alternatives might be difficult to implement when the scrutinee is defined in a library or in third-party code. There is an ongoing proposal^{18,19} [Goetz, 2017a] to add pattern matching to the Java language. The proposal explores changing the instanceof operator in order to support pattern matching. Java 12 extends the switch statement to be used as either a statement or an expression^{20,21} [Goetz, 2017b; Bierman, 2019]. This enhancement aims to

¹⁶<https://kotlinlang.org/docs/reference/typecasts.html#smart-casts>

¹⁷Adapted from <https://docs.scala-lang.org/tour/pattern-matching.html>

¹⁸<https://openjdk.java.net/jeps/305>

¹⁹<https://cr.openjdk.java.net/~briangoetz/amber/pattern-match.html>

²⁰<https://openjdk.java.net/jeps/325>

²¹<https://openjdk.java.net/jeps/354>

ease the transition to a switch expression that supports pattern matching.

The *GuardByClassLiteral* variant may be used instead of the *instanceof* operator when the developer wants to match exactly the runtime class of an object. The *instanceof* operator²² returns true if the expression could be cast to the specified type, whereas using a class literal comparison returns true if the expression is exactly the runtime class.

In some cases, the *GuardByTypeTag* variant can be replaced by *GuardByInstanceOf*. However, if the application-specific tag is a numeric value, the *GuardByTypeTag* could perform better than the *GuardByInstanceOf* using *instanceof*. Moreover, there are situation where the *instanceof* operator cannot be avoid since the types to be cast are the same.

4.5.2 Stash

Description. This pattern is used to stash an application-specific value. It has three variants. The *LookupById* and *StaticResource* variants are used to extract values from a heterogenous container. The *Tag* variant is used to extract a “tag” value, typically in a GUI object or message payload. They look up an object by a compile-time constant identifier, tag, or name and casts the result to an appropriate type. They access a collection that holds values of different types (usually implemented as `Collection<Object>` or as `Map<K, Object>`). The actual run-time type returned from the lookup is determined by the value of the identifier.

The *StaticResource* variant of is more specific, it is used to retrieve a value instantiated from static resource file, *e.g.*, an XML, HTML or Java properties file. The file contents are (in theory) known at compile-time and the file is included in the binary distribution of the application. These files are often built using tools such as GUI builders.

Instances: 561 (11.22%). We found 356 in application code, 63 in test code, and 142 in generated code. Figure 4.3 shows different variants of the pattern. The *LookupById* is the most used variant.

In the *LookupById* variant example shown below, the return type of the *getAttribute* method is `Object`. The variable context is of type `BasicHttpContext`, which is implemented with `HashMap`.

²²<https://docs.oracle.com/javase/specs/jls/se8/html/jls-15.html#jls-15.20.2>

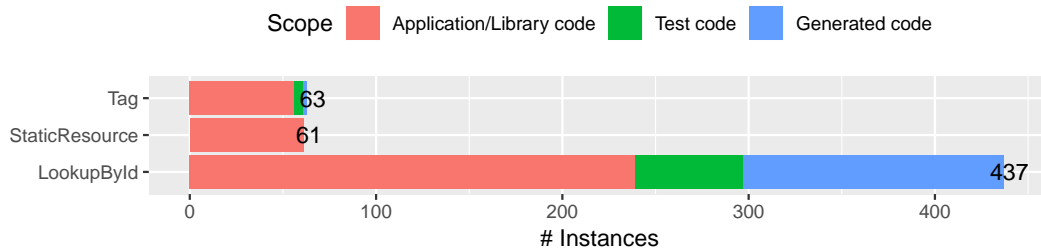


Figure 4.3. STASH Variant Occurrences

```

1 AuthState authState = (AuthState) context.getAttribute(
2     ClientContext.TARGET_AUTH_STATE);
3                                     http://bit.ly/loopj\_android\_async\_http\_2SUzY4E

```

The next snippet shows a call site to the `getComponent` method cast to the `ActiveListManager` class (line 14). The `getComponent` method in this cast instance uses as argument the `PROP_ACTIVE_LIST_MANAGER` constant. Looking at the definition of this constant (line 3), we can see there is a companion attribute (`@S4Component`) whose argument is the `ActiveListManager` class, the target of the cast instance.

```

1 /** The property that defines the type of active list to use */
2 @S4Component(type = ActiveListManager.class)
3 public final static String PROP_ACTIVE_LIST_MANAGER = "activeListManager";
4
5 @Override
6 public void newProperties(PropertySheet ps) throws PropertyException {
7     super.newProperties(ps);
8     logMath = (LogMath) ps.getComponent(PROP_LOG_MATH);
9     logger = ps.getLogger();
10    linguist = (Linguist) ps.getComponent(PROP_LINGUIST);
11    pruner = (Pruner) ps.getComponent(PROP_PRUNER);
12    scorer = (AcousticScorer) ps.getComponent(PROP_SCORER);
13    activeListManager =
14        (ActiveListManager) ps.getComponent(PROP_ACTIVE_LIST_MANAGER);
15    // [...]
16 }

```

http://bit.ly/skerit_cmusphinx_2HGgL1D

In the following example, a cast is applied to the result of looking up by index in the `iContexts` map (line 9). In case there is no value for the given index, a value of the corresponding type is stored using the same index (line 13), thus guaranteeing the success of the cast.

```

1  protected Map<Integer, AssignmentContext> iContexts =
2      new HashMap<Integer, AssignmentContext>();
3
4  @Override
5  @SuppressWarnings("unchecked")
6  public <U extends AssignmentContext> U getAssignmentContext(
7      Assignment<V, T> assignment,
8      AssignmentContextReference<V, T, U> reference) {
9      U context = (U) iContexts.get(reference.getIndex());
10     if (context != null) return context;
11
12     context = reference.getParent().createAssignmentContext(assignment);
13     iContexts.put(reference.getIndex(), context);
14     return context;
15 }

```

http://bit.ly/UniTime_cpsolver_2HUmGki

The following *StaticResource* example is from an Android application. A cast is applied to the `findViewById` method invocation. View classes are instantiated by the application framework using an XML resource file. The `findViewById` method looks up the view by its ID.

```

1  @Override
2  protected void onCreate(Bundle savedInstanceState) {
3      super.onCreate(savedInstanceState);
4      setContentView(R.layout.activity_main);
5      connectivityStatus = (TextView) findViewById(R.id.connectivity_status);
6      mobileNetworkType = (TextView) findViewById(R.id.mobile_network_type);
7      accessPoints = (ListView) findViewById(R.id.access_points);
8      busWrapper = getOttoBusWrapper(new Bus());
9      networkEvents = new NetworkEvents(getApplicationContext(), busWrapper)
10         .enableInternetCheck()
11         .enableWifiScan();
12 }

```

http://bit.ly/pwittchen_NetworkEvents_2HGbrMq

The next listing, shows a cast to a GUI component (`XulListbox`) using the `getElementById` method (lines 12 and 13). In this case the developer is using the XUL language.²³

```

1  private void createBindings() {
2      loginDialog = (XulDialog) document
3          .getElementById( "repository-login-dialog" );
4      repositoryEditButton = (XulButton) document
5          .getElementById( "repository-edit" );
6      repositoryRemoveButton = (XulButton) document

```

²³<https://developer.mozilla.org/en-US/docs/Mozilla/Tech/XUL>

```

7         .getElementById( "repository-remove" );
8     username = (XulTextbox) document
9         .getElementById( "user-name" );
10    userPassword = (XulTextbox) document
11        .getElementById( "user-password" );
12    availableRepositories = (XulListbox) document
13        .getElementById( "available-repository-list" );
14    // [...]
15 }

```

http://bit.ly/pentaho_pentaho_kettle_2TswNSf

In the following snippet of the *Tag* variant, a cast is performed to a *getSerializable* invocation (lines 9 and 10). This method gets a *Serializable* value given the specified key, *TAG_CUR_DIR* in this case. To set a value with a specified key, the *putSerializable* method is used. The mentioned cast succeeds because a value of the appropriate type is set in line 18 using the *putSerializable* method.

```

1 private TorrentContentFileTree curDir;
2
3 @Override
4 public void onActivityCreated(@Nullable Bundle savedInstanceState) {
5     super.onActivityCreated(savedInstanceState);
6     if (activity == null)
7         activity = (AppCompatActivity) getActivity();
8     if (savedInstanceState != null) {
9         curDir = (TorrentContentFileTree) savedInstanceState
10             .getSerializable(TAG_CUR_DIR);
11     } else {
12         makeFileTree();
13     }
14 }
15
16 @Override
17 public void onSaveInstanceState(Bundle outState) {
18     outState.putSerializable(TAG_CUR_DIR, curDir);
19 }

```

http://bit.ly/proninyaroslav_libretorrent_2TxpZCM

In the last example, the cast is applied to a *getModel* invocation on the *matchTable* field (line 16). Looking how *matchTable* is initialized (line 7), the *model* variable (line 5) is used as an argument to the constructor. This argument is the value returned by *getModel*, and since they are both of the same type, the mentioned cast is guaranteed to succeed.

```

1 public final class MatchPanel extends JPanel implements Observer {
2     private final JZebraTable matchTable;
3     public MatchPanel() {
4         super(new GridBagLayout());
5         DefaultTableModel model = new DefaultTableModel();
6         // [...]
7         matchTable = new JZebraTable(model) {
8             @Override
9             public boolean isCellEditable(int rowIndex, int colIndex) {
10                 return false;
11             }
12         };
13     }
14     // [...]
15     private void observe(GamerCompletedMatchEvent event) {
16         DefaultTableModel model = (DefaultTableModel) matchTable.getModel();
17         model.setValueAt("Inactive", model.getRowCount() - 1, 4);
18     }
19 }

```

http://bit.ly/ggp_org_ggp_base_2SAEXHu

Detection. The implementation of the two variants, *StaticResource* and *Tag*, is application-specific. Thus, detecting them is often impractical. However, if the methods that perform the specified patterns are known, e.g., `findViewById`, then the automatic detection becomes trivial. On the other hand, the following query detects the *LookupById* variant.

```

1 class LookupByIdCast extends Cast {
2     MethodAccess ma;
3     Method getterMethod;
4     FieldAccess constant;
5     LookupByIdCast() {
6         ma = getExprOrDef() and getterMethod = ma.getMethod() and
7         not getterMethod.isVarargs() and getterMethod.isPublic() and
8         getterMethod.getNumberOfParameters() = 1 and
9         getterMethod.getParameterType(0) instanceof TypeString and
10        getterMethod.getReturnType() instanceof TypeObject and
11        ma.getArgument(0) = constant and
12        constant.getField().isFinal() and constant.getField().isStatic()
13    }
14 }

```

Listing 4.13. Detection of the *LookupById* variant

Issues. This pattern suggests a heterogeneous dictionary. In our manual inspection, all dictionary keys and the resulting types are known at compile time, however a cast is needed because the dictionary type does not encode the relationship between key values and the result type. Casts in this pattern are typically not guarded indicating that the programmer knows the source of the cast based on the value of the key. The *LookupById* variant could be replaced by strongly typed heterogeneous collections [Kiselyov et al., 2004] although implementing it in Java would be more verbose.

The *StaticResource* variant is often seen in Android applications. The Butter Knife framework²⁴ uses annotations to avoid the “manual” casting. Instead, code is generated that casts the result of `findViewById` to the appropriate type. These casts could be solved by using code generation, or partial classes like in C#. Since the contents of the resource file are known at compile-time, code generation could be used to generate the corresponding Java code. In our sample, however, this variant only appears in application code.

The *Tag* variant can also be used to fetch a value from a collection (as in *LookupById*). The main difference is “locality”. That is, in the *Tag* variant the cast value is set “locally”, i.e., in the same method or class, whereas the cast value in the *LookupById* variant is usually set in another class.

Since this pattern casts a value to a known type from a method invocation, it can be seen as a kind of `KNOWNRETURN TYPE` pattern.

4.5.3 Factory

Description. Creates an object based on some arguments to a method call. Since the arguments are known at compile-time, cast to the specific type. In this pattern, the arguments resemble a “type tag” descriptor (cf. `TYPECASE`).

This pattern is characterized by a cast to a method call passing one or more arguments. The method call needs to create an object based on those arguments. Usually the arguments that determine the run-time type to be returned are known at compile-time.

Instances: 381 (7.62%). We found 144 in application code, 149 in test code, and 88 in generated code. The following snippet shows an instance of the `FACTORY` pattern. The cast is applied to the result of invoking `keyPair.getPrivate` (line 6). The variable `keyPair` is assigned the result of `pairGen.-`

²⁴<http://jakewharton.github.io/butterknife/>

generateKeyPair (line 3). At the same time, the pairGen variable is assigned the value returned by KeyPairGenerator.getInstance("RSA"). The argument "RSA" indicates the algorithm to use. The method²⁵ will return a reference to the private key component, and this is determined by the algorithm argument described above.

```

1 KeyPairGenerator pairGen = KeyPairGenerator.getInstance("RSA");
2 pairGen.initialize(1024);
3 KeyPair keyPair = pairGen.generateKeyPair();
4 // [...]
5 RSAKey rsaJWK2 = new RSAKey.Builder((RSAPublicKey) keyPair.getPublic())
6     .privateKey((RSAPrivateKey) keyPair.getPrivate())
7     .keyID("2")
8     .build();

```

http://bit.ly/connect2id_oauth_2_0_sdk_with_2HvRIUX

Similar to the above snippet, the next example shows an instance of this pattern where a cast is performed on the result of the openConnection method²⁶ (line 2). The method is declared to return URLConnection but can return a more specific type based on the URL string. The openConnection method is applied to the url variable, which is assigned in line 1 using the URL constructor. The argument to the constructor is an http URL, thus the result is cast to HttpURLConnection.

```

1 URL url = new URL("http://localhost:8088/ws/v1/cluster/apps");
2 HttpURLConnection conn = (HttpURLConnection) url.openConnection();
3

```

http://bit.ly/apache_hadoop_2E6KY6T

The following example shows how a cast (line 3) is being determined by the argument to the CertificateFactory.getInstance method (line 1). The argument is the string "X.509", therefore the method generateCRL will return a value of type X509CRL.

```

1 CertificateFactory cf = CertificateFactory.getInstance("X.509", "BC");
2 // [...]
3 X509CRL crl = (X509CRL)cf.generateCRL(new ByteArrayInputStream(directCRL));
4

```

http://bit.ly/bcgit_bc_java_2TEVScM

In our last example the cast instance (line 2) is applied to the result of parse method. The return type of parse is of type Statement, but, since the

²⁵[https://docs.oracle.com/javase/8/docs/api/java/security/KeyPair.html#getPrivate\(\)](https://docs.oracle.com/javase/8/docs/api/java/security/KeyPair.html#getPrivate())

²⁶<https://docs.oracle.com/javase/8/docs/api/java/net/URL.html#openConnection-->

statement is a SELECT statement, the value returned by the parse method is known to be of type Select and the cast should succeed.

```

1 statement = "SELECT * FROM mytable WHERE mytable.col = 9 LIMIT :param_name";
2 select = (Select) parserManager.parse(new StringReader(statement));
3 public class Select implements Statement {
4     // [...]
5 }
6 public class CCJSqlParserManager implements JSqlParser {
7     @Override
8     public Statement parse(Reader statementReader) throws JSqlParserException {
9         // [...]
10    }
11 }

```

http://bit.ly/JSQLParser_JSqlParser_2TecMyB

In some cases of this pattern, a cast is applied to a method invocation where one of its arguments is a class literal. The target type of the cast is determined by this class literal, like in the following snippets.

```

1 final ILiferayServerBehavior liferayServerBehavior =
2     (ILiferayServerBehavior) moduleServer.getServer()
3     .loadAdapter( ILiferayServerBehavior.class, null );
4

```

http://bit.ly/liferay_liferay_ide_2FMG0f6

```

1 CFArray o = (CFArray) CType.Marshaler.toObject(CFArray.class, handle, flags);
2

```

http://bit.ly/robovm_robovm_2FMFWvS

Detection. The detection of this pattern requires to analyse the factory method being called. This is not always possible in QL, since QL does not analyse project dependencies.

In several instances, to manually determine when a cast belongs to this pattern, we had to look-up the method implementation in external source code repositories.

Issues. In some situations, the use of this pattern can be seen as breaking the contract API between the caller and the callee. This happens because the caller needs to know how the method is implemented in order to determine the run-time return type. In FACTORY, there is a known type hierarchy below the return type and the caller casts to a known subtype in that hierarchy based on the arguments passed into the factory method.

The `KNOWNRETURNType` pattern is similar to `FACTORY`, since both depend on the knowledge that a method returns a more specific type.

This pattern is prevalent in test code 39.11%. This is because when testing, known parameters are given to factory methods. In these situations, a test method needs to know a more specific type—by using a cast—to properly check for a test condition.

4.5.4 Family

Description. The `FAMILY` pattern implements casts to provide a sort of *family polymorphism* [Ernst, 2001]. A “family” consists of multiple mutually-dependent types designed to collaborate with each other. Each type has a role in the family. Deriving from a base family to form another family requires subclassing all the members of the base family, with the subclasses in the new family retaining their roles in the new family.

Because method parameter types are invariant in Java and because covariant parameter types are unsound in general, the method parameter types in the derived family are the same as in the base family. Casts are therefore necessary for one member of a derived family to access another member using its derived family type rather than its base family type.

Instances: 344 (6.88%). We found 257 in application code, 37 in test code, and 50 in generated code. The following example shows an instance of this pattern. In this case, the interfaces `StepInterface`, `StepMetaInterface`, and `StepDataInterface` are part of a base family and the `stopRunning` method has parameters of these types. In the derived family the roles of these three interfaces are implemented by the classes `DynamicSQLRow`, `DynamicSQLRowMeta`, and `DynamicSQLRowData`. A cast is applied to the parameter `smi` of `stopRunning` in `DynamicSQLRow` (line 12). This cast is necessary to convert the method parameter, of the base family type `StepDataInterface`, into the derived family type with the same role.

```
1 public interface StepInterface extends VariableSpace, HasLogChannelInterface {
2     // [...]
3     public void stopRunning( StepMetaInterface stepMetaInterface,
4                             StepDataInterface stepDataInterface ) throws KettleException;
5 }
6 public class DynamicSQLRow extends BaseStep implements StepInterface {
7     private DynamicSQLRowMeta meta;
8     private DynamicSQLRowData data;
```

```

9 // [...]
10 public void stopRunning( StepMetaInterface smi, StepDataInterface sdi )
11     throws KettleException {
12     meta = (DynamicSQLRowMeta) smi;
13     data = (DynamicSQLRowData) sdi;
14     // [...]
15 }
16 }

```

http://bit.ly/pentaho_pentaho_kettle_2FN59J8

The next example is similar to the previous one. The masked parameter is cast to `DoubleColumnVector` (line 5). It is so because the masked variable is expected to hold an instance of `DoubleColumnVector` when the `maskData` method is applied to an object of type `DoubleIdentity`.

```

1 public class DoubleIdentity implements DataMask {
2     @Override
3     public void maskData(ColumnVector original, ColumnVector masked, int start,
4                           int length) {
5         DoubleColumnVector target = (DoubleColumnVector) masked;
6         DoubleColumnVector source = (DoubleColumnVector) original;
7         // [...]
8     }
9 }
10 public interface DataMask {
11     // [...]
12     void maskData(ColumnVector original, ColumnVector masked,
13                   int start, int length);
14 }

```

http://bit.ly/apache_orc_2SE4C2m

In both previous examples, casts were applied to a parameter in an overriding method. In the next example, the cast instance is applied to super class field (line 12). The field is declared in the `BaseExchange` class (line 20). However, the field is initialized with a `BitflyerMarketDataService` value in line 5.

```

1 public class BitflyerExchange extends BaseExchange implements Exchange {
2     // [...]
3     @Override
4     protected void initServices() {
5         this.marketDataService = new BitflyerMarketDataService(this);
6         // [...]
7     }
8     // [...]
9     @Override
10    public void remoteInit() throws IOException, ExchangeException {
11        BitflyerMarketDataServiceRaw dataService =

```

```

12     (BitflyerMarketDataServiceRaw) this.marketDataService;
13     List<BitflyerMarket> markets = dataService.getMarkets();
14     exchangeMetaData = BitflyerAdapters.adaptMetaData(markets);
15 }
16 }
17 public abstract class BaseExchange implements Exchange {
18     // [...]
19     protected MarketDataService marketDataService;
20     // [...]
21 }

```

http://bit.ly/knowm_XChange_2UPPDj9

Detection. To detect this pattern, the cast needs to be applied to a family. A family is distinguished by a covariant usage of a field or parameter in an overriding method. Since this pattern has many small variations, writing a QL query would be impractical.

Issues. Java itself does not support statically type-safe family polymorphism directly and so casts are often necessary. Various proposals have been made to better support family polymorphism (and the related “expression problem” [Wadler, 1998]) in object-oriented languages, including the use of design patterns [Wang and Oliveira, 2016; Oliveira and Cook, 2012; Nystrom et al., 2003], and type systems [Ernst, 2000; Odersky and Zenger, 2005; Myers, 2006; Oliveira et al., 2016; Kiselyov et al., 2009] that permit some restricted form of covariant method parameters.

4.5.5 UseRawType

Description. A cast is in the USERAWTYPE pattern when a *raw type* is used rather than a generic type. Methods of raw types typically return `Object` rather than a more specific type.

Instances: 335 (6.70%). We found 176 in application code, 18 in test code, and 141 in generated code. For example, in the following code, the collection `c` and iterator `it` are declared to be of the raw types `Collection` and `Iterator` rather than as parameterized types. The call to `next` on line 4 must be cast to a more specific type because static type information was lost by the use of raw types.

```

1 Collection c = recipients.getRecipients();
2 assertTrue(c.size() >= 1 && c.size() <= 2);
3 Iterator it = c.iterator();
4 verifyRecipient((RecipientInformation)it.next(), privKey);
5

```

http://bit.ly/bcgit_bc_java_2SD2HLM

The following example uses the Comparable interface (line 1). This interface is generic,²⁷ but in this case the developer is using its raw type. Therefore a cast is needed in line 5.

```

1 public class McpSettlementDetailDto implements Comparable {
2     // [...]
3     @Override
4     public int compareTo(Object o){
5         McpSettlementDetailDto mcpSettlementDetailDto=(McpSettlementDetailDto)o;
6         Integer newConsume=((int)mcpSettlementDetailDto.getConsume());
7         Integer temp=((int)this.consume);
8         return temp.compareTo(newConsume);
9     }
10 }

```

http://bit.ly/fangjie008_tixue_mcp_parent_2FSZKzm

In the following snippet, a cast is applied to the result of the doPrivileged method in lines 3 and 4. This method takes a PrivilegedAction<T>, but the cast is needed because it is invoked with a raw type, *e.g.*, new PrivilegedAction(). Inspecting further the source code application, we found that it might be a requirement to be compatible with the JDK 1.2. Generics were added to Java 5. Thus, this cast might be still necessary.

```

1 class SecuritySupport12 extends SecuritySupport {
2     ClassLoader getSystemClassLoader() {
3         return (ClassLoader)
4             AccessController.doPrivileged(new PrivilegedAction() {
5                 public Object run() {
6                     ClassLoader cl = null;
7                     try {
8                         cl = ClassLoader.getSystemClassLoader();
9                     } catch (SecurityException ex) {}
10                    return cl;
11                }
12            });
13     }
14 }
15 public final class AccessController {
16     public static <T> T doPrivileged(PrivilegedAction<T> action) {

```

²⁷<https://docs.oracle.com/javase/8/docs/api/java/lang/Comparable.html>

```

17         return action.run();
18     }
19 }
20 public interface PrivilegedAction<T> {
21     public T run();
22 }

```

http://bit.ly/robovm_robovm_2FAI5x5

Detection. The Query 4.14 detects a variation of the `USERAWTYPE` pattern, *e.g., only* the first example shown above. That is, when a cast is applied to a method *declared* as returning a generic type, but the method is invoked on an object defined as a raw type.

```

1 class UseRawTypeCast extends Cast {
2     MethodAccess ma;
3     RawType rt;
4     UseRawTypeCast() {
5         ma = getExpr() and
6         rt = ma.getQualifier().getType() and
7         ma.getMethod().getSourceDeclaration().getReturnType()
8             instanceof TypeVariable
9     }
10    MethodAccess getMethodAccess() {
11        result = ma
12    }
13 }

```

Listing 4.14. Detection of the `USERAWTYPE` pattern.

Issues. Raw types exist in Java to support legacy code. Best practice would be to rewrite the code to use generics, but this is not always feasible or cost effective.

This pattern is prevalent in generated code (42.09% of generated instances). Since these casts will not be seen by a developer, code generators make less effort to avoid them.

Casts among generic types and between raw types and generic types are unchecked at run time, although other casts are typically inserted by the compiler to ensure type safety dynamically. When these inserted casts fail, the reported location of the failure may not match the programmer's expectation. Indeed, this is similar to the problem of *blame* in gradually

typed languages [Wadler and Findler, 2009]. In this setting, when a runtime cast fails the blame should be put on the appropriate programmer-inserted cast, not on a compiler-inserted cast.

4.5.6 Equals

Description. This pattern is a common pattern to implement the well-known equals method (declared in `java.lang.Object`). It is a particularly instance of guarded casts. A cast expression is guarded by either an instanceof test—*InstanceOf* variant—or a getClass comparison—*GetClass* variant—usually to the same target type as the cast; in an equals²⁸ method implementation. This is done to check if the argument has same type as the receiver (this argument). Notice that a cast in an equals method is needed because it receives an `Object` as a parameter.

To detect this pattern, a cast must be applied to the parameter of the equals method. The result value of the cast must be then used in an equality comparison. We relax the constraint that the target type of the cast must be the enclosing class.

Instances: 247 (4.94%). We found 202 in application code, 0 in test code, and 45 in generated code. This pattern accounts for 16.94% of guarded casts, 247 instances out of 1,458. Figure 4.4 shows the different variants of the EQUALS pattern and their occurrences. The *InstanceOfSupertype*, *AutoValue*, and *InstanceOfSwitch* variants are explained below.

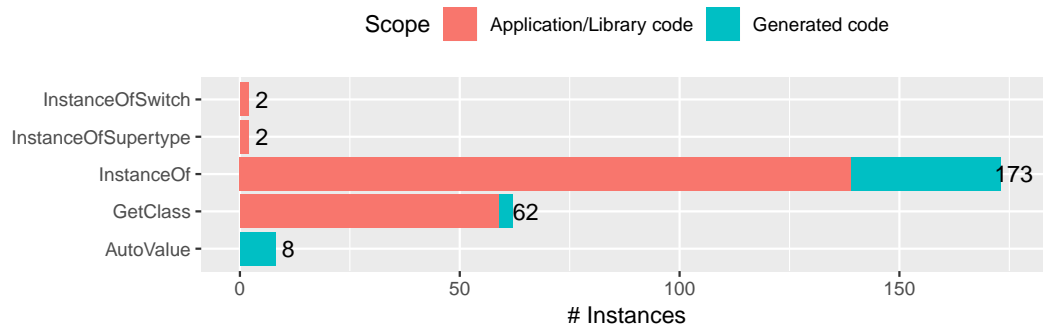


Figure 4.4. EQUALS Variant Occurrences

²⁸<https://docs.oracle.com/javase/8/docs/api/java/lang/Object.html#equals-java.lang.Object->

The following listing shows an example of the EQUALS pattern. In this case, an instanceof guards for the same type as the receiver (*InstanceOf* variant).

```

1  @Override
2  public boolean equals(Object obj) {
3      if ( this == obj ) {
4          return true;
5      }
6      if ( (obj instanceof Difference) ) {
7          Difference that = (Difference) obj;
8          return actualFirst == that.actualFirst
9              && expectedFirst == that.expectedFirst
10             && actualSecond == that.actualSecond
11             && expectedSecond == that.expectedSecond
12             && key.equals( that.key );
13      }
14      return false;
15  }
```

http://bit.ly/neo4j_neo4j_2vJw94J

Alternatively, in the following listing we show another example of the EQUALS pattern. But in this case, a getClass comparison is used to guard for the same type as the receiver in line 4 (*GetClass* variant).

```

1  @Override
2  public boolean equals( Object o ) {
3      if ( this == o ) return true;
4      if ( o == null || getClass() != o.getClass() )
5          return false;
6
7      ValuePath that = (ValuePath) o;
8      return nodes.equals(that.nodes) &&
9             relationships.equals(that.relationships);
10 }
```

http://bit.ly/neo4j_neo4j_2vKP0MW

In some situations, the type cast is not the same as the enclosing class. Instead, the target type of the cast is the super class or a super interface of the enclosing class (*InstanceOfSupertype* variant). The following example shows this scenario. The cast is performed in the WildcardTypeImpl enclosing class, but the target type is java.lang.reflect.WildcardType.

```

1  public static class WildcardTypeImpl implements WildcardType, CompositeType {
2      @Override
3      public boolean equals(Object other) {
4          return other instanceof WildcardType
5              && MoreTypes.equals(this, (WildcardType) other);
6      }
```

```

6     }
7 }

```

http://bit.ly/elastic_elasticsearch_2GHYPp5

Similar to the previous example, the *AutoValue* variant casts the equals parameter to a super class of the enclosing class. However, this happens when the Google AutoValue library²⁹ is used. AutoValue is a code generator for value classes.

```

1 @AutoValue
2 abstract class ListsItem implements Parcelable { /* [...] */ }
3
4 abstract class $AutoValue_ListsItem extends ListsItem {
5     @Override
6     public boolean equals(Object o) {
7         if (o == this) {
8             return true;
9         }
10        if (o instanceof ListsItem) {
11            ListsItem that = (ListsItem) o;
12            return (this.id == that.id())
13                && (this.name.equals(that.name()))
14                && (this.itemCount == that.itemCount());
15        }
16        return false;
17    }
18 }

```

http://bit.ly/square_sqlbrite_2HmHMYE

The following snippet shows a non-trivial implementation of equals. The enclosing class of the equals method is CapReq (line 1). However, the cast instance (line 13) is not against the enclosing class, it is against to the Requirement class (*InstanceOfSwitch* variant). Note that the cast using the enclosing class as target type is in line 9.

```

1 class CapReq {
2     @Override
3     public boolean equals(Object obj) {
4         if (this == obj)
5             return true;
6         if (obj == null)
7             return false;
8         if (obj instanceof CapReq)
9             return equalsNative((CapReq) obj);
10        if ((mode == MODE.Capability) && (obj instanceof Capability))
11            return equalsCap((Capability) obj);

```

²⁹<https://github.com/google/auto/tree/master/value>

```

12         if ((mode == MODE.Requirement) && (obj instanceof Requirement))
13             return equalsReq((Requirement) obj);
14         return false;
15     }
16 }

```

http://bit.ly/bndtools_bnd_2SM5pOw

Detection. This pattern contains several variants. The Query 4.15 detects three of them, *i.e.*, *InstanceOf*, *GetClass*, and *AutoValue*. It is not difficult to extend this query to detect the other variants.

```

1  class EqualsCast extends VarCast {
2      EqualsCast() {
3          getVar() instanceof Parameter and
4          getEnclosingCallable() instanceof EqualsMethod and (
5              this instanceof GuardByInstanceOfCast or
6              this instanceof GetClassGuardsVarCast
7          ) and (
8              getTargetType() = getEnclosingCallable().getDeclaringType() or
9              (
10                 getTargetType() = getEnclosingCallable().getDeclaringType().
11                     getASupertype+() and
12                 getEnclosingCallable().getDeclaringType()
13                     instanceof AutoValueGenerated
14             )
15          )
16      }
17 }

```

QL

Listing 4.15. Detection of the EQUALS pattern.

Issues. The pattern for an equals method implementation is well-known. Most equals methods in our sample are implemented with the same boilerplate structure: that is, first checking if the parameter is another reference to this, then checking if the argument is not null, and finally, checking if the argument is of the right class (with either an instanceof test or a getClass comparison). Once all checks are performed, a cast follows, and a field-by-field comparison is made.

To avoid this boilerplate, other languages bake in deep equality comparisons, at least for some types (*e.g.*, Scala case classes), or provide mechanisms to generate the boilerplate code (*e.g.*, deriving Eq in Haskell or `#[derive(Eq)]` in Rust). Vaziri et al. [2007] propose a declarative approach

to avoid boilerplate code when implementing both the `equals` and `hashCode` methods. They manually analysed several applications, and found there are many issues while implementing `equals()` and `hashCode()` methods. It would be interesting to check whether these issues happen in real application code.

There is an exploratory document³⁰ by Brian Goetz, Java Language Architect, addressing these issues from a more general perspective. It is definitely a starting point towards improving the Java language.

This pattern can be seen as a special instance of the `TYPECASE` pattern when the guard is an `instanceof` test or a `getClass` comparison.

4.5.7 Redundant

Description. A redundant cast is a cast that is not necessary for compilation. The cast could be removed from source code without affecting the application.

To detect the `REDUNDANT` pattern, the expression being cast needs to be of the same type as the type being cast to.

Instances: 122 (2.44%). We found 66 in application code, 14 in test code, and 42 in generated code. The following listing exhibits an instance of the `REDUNDANT` pattern. A redundant cast is applied to a lambda expression (line 8). This cast is not needed a Java compiler can infer that the lambda expression is of type `TransactionCallback<Void>` (defined in line 22).

```
1 public class FlywayTest {
2     private TransactionTemplate transactionTemplate;
3     @Test
4     public void test() {
5         // [...]
6         transactionTemplate.execute(
7             (TransactionCallback<Void>) transactionStatus -> {
8                 Post post = new Post();
9                 entityManager.persist(post);
10                return null;
11            });
12     }
13 }
14 public interface TransactionStatus { /* [...] */ }
15 @FunctionalInterface
```

³⁰<http://cr.openjdk.java.net/~briangoetz/amber/datum.html>

```

16 public interface TransactionCallback<T> {
17     T doInTransaction(TransactionStatus status);
18 }
19 public class TransactionTemplate {
20     <T> T execute(TransactionCallback<T> action) { /* [...] */ }
21 }

```

http://bit.ly/vladmihalcea_high_performance_java_persistence_2FWXw2e

The next cast instance is trivially redundant: both the target type and the static type of the operand `count(b)`, are `BigDecimal`.

```

1 @Override
2 public void accumulate(Tuple b) throws IOException {
3     // [...]
4     BigDecimal count = (BigDecimal)count(b);
5     // [...]
6 }
7 static protected BigDecimal count(Tuple input) throws ExecException {
8     // [...]
9 }

```

http://bit.ly/sigmoidanalytics_spork_2SIqWYq

In the following cast instance, a cast is applied to the `node.right` field (line 9). Nevertheless, the right field of the `Node` class is already defined as `Node<T>`, rendering the cast redundant.

```

1 public class ImplicitKeyTreap<T> implements IList<T> {
2     protected Node<T> root = null;
3     // [...]
4     private int getIndexByValue(T value) {
5         final Node<T> node = (Node<T>)root;
6         if (value == null || node == null)
7             return Integer.MIN_VALUE;
8         final Node<T> l = (Node<T>)node.left;
9         final Node<T> r = (Node<T>)node.right;
10        // [...]
11        return i;
12    }
13    public static class Node<T> {
14        private T value = null;
15        private int priority;
16        private int size;
17        private Node<T> parent = null;
18        private Node<T> left = null;
19        private Node<T> right = null;
20        // [...]
21    }
22 }

```

http://bit.ly/phishman3579_java_algorithms_implementation_2SGcH6w

There are cases when code generators insert superfluous casts to null. The following cast instance could be removed since in this case the cast to null is not needed.

```
1 public groovy.lang.MetaClass getMetaClass() {
2     return (groovy.lang.MetaClass) null;
3 }
```

http://bit.ly/togglz_togglz_2SGncXB

Detection. The following query returns casts where the static type of the cast expression is the exactly the same as the target type. That is, casts (T) e where e is declared as T. The query also detects a redundant upcast, *i.e.*, an upcast that is not used for neither the `SELECTOVERLOAD` nor the `COVARIANTGENERIC` patterns.

```
1 class RedundantCast extends Cast {
2     RedundantCast() {
3         getExpr().getType() = getTargetType() or (
4             this instanceof Upcast and
5             not this instanceof SelectOverloadCast and
6             not this instanceof CovariantGenericCast
7         )
8     }
9 }
```

QL

Listing 4.16. Detection query for the REDUNDANT pattern

Issues. Redundant casts are generally upcasts or casts involving erased type parameters. This pattern arises often in generated code. It may also appear due to code refactoring that change a type and therefore make the cast redundant.

4.5.8 CovariantReturnType

Description. The `COVARIANTRETURN TYPE` pattern is used to cast a call to a method that returns an instance of a type that is covariant with the receiver type. Commonly the method returns a instance of the receiver type itself.

Instances: 106 (2.12%). We found 85 in application code, 3 in test code, and 18 in generated code. A common instance of this pattern is for calls

to the clone method of java.lang.Object (70 instances), which returns an object of the same type as the receiver, but whose static type is Object. The following snippet shows a cast to the clone method.

```

1  @Override
2  public ListTagsForResourceResult clone() {
3      try {
4          return (ListTagsForResourceResult) super.clone();
5      } catch (CloneNotSupportedException e) {
6          throw new IllegalStateException(/* [...] */);
7      }
8  }

```

http://bit.ly/aws_aws_sdk_java_2GvHhYt

In the following example, the unmarshall method overrides a superclass method with a covariant return type. A cast is used on the call to the superclass method to change the type of the return value to match the more precise return type.

```

1  public class ResourceContentionExceptionUnmarshaller
2      extends StandardErrorUnmarshaller {
3      public ResourceContentionExceptionUnmarshaller() {
4          super(ResourceContentionException.class);
5      }
6      public AmazonServiceException unmarshall(Node node) throws Exception {
7          // Bail out if this isn't the right error code that this
8          // marshaller understands.
9          String errorCode = parseErrorCode(node);
10         if (errorCode == null || !errorCode.equals("ResourceContention"))
11             return null;
12         ResourceContentionException e =
13             (ResourceContentionException) super.unmarshall(node);
14         return e;
15     }
16 }

```

http://bit.ly/aws_amplify_aws_sdk_android_2FVWl13

The initCause method—from the java.lang.Throwable class—has return type Throwable. Nevertheless, this method returns the receiver (after setting the cause exception). Therefore a cast is needed to recover the original exception type, as shown in the following example. This use case resembles the FLUENTAPI pattern.

```

1  throw (IllegalArgumentException)
2      new IllegalArgumentException("Invalid broker URI: " + brokerURL)
3      .initCause(e);
4

```

http://bit.ly/apache_activemq_2EnSivc

Detection. The Query 4.17 approximates the detection of the COVARIANTRETURNTYPE pattern when a cast is applied to a method in a super class, *e.g.*, in the first two examples shown above.

```

1 class CovariantReturnTypeCast extends Cast {
2     Method m;
3     MethodAccess ma;
4     CovariantReturnTypeCast() {
5         getExpr() = ma and ma.isOwnMethodAccess() and
6         getEnclosingCallable() = m and m.overrides(ma.getMethod())
7     }
8 }

```

QL

Listing 4.17. COVARIANTRETURNTYPE detection query.

Issues. The situation of returning this could be avoided if Java supported self types [Bruce, 2003]. More generally, associated types [Chakravarty et al., 2005] can provide a statically typed solution, for instance in the second example above.

4.5.9 SelectOverload

Description. This pattern is used to select the appropriate version of an overloaded method³¹ where two or more of its implementations differ *only* in some argument type.

A cast to null is often used to select against different versions of a method, *i.e.*, to resolve method overloading ambiguity. Whenever a null value needs to be an argument, a cast is needed to select the appropriate implementation. This is because the type of null has the special type *null*³² which can be treated as any reference type. In this case, the compiler cannot determine which method implementation to select.

Another use case is to select the appropriate the right argument when calling a method with variable arguments.

Instances: 99 (1.98%). We found 51 in application code, 47 in test code, and 1 in generated code. The following listing shows an example of the SELECTOVERLOAD pattern. In this example, there are three versions of the

³¹Using ad-hoc polymorphism [Strachey, 2000].

³²<https://docs.oracle.com/javase/specs/jls/se8/html/jls-4.html#jls-4.1>

onSuccess method. The cast (String) null is used to select the appropriate version (line 7), based on the third parameter, between overloaded methods that differ only in their argument type (the third one).

```

1 onSuccess(statusCode, headers, (String) null);
2 public void onSuccess(
3     int statusCode, Header[] headers, JSONObject response) { /* [...] */ }
4 public void onSuccess(
5     int statusCode, Header[] headers, JSONArray response) { /* [...] */ }
6 public void onSuccess(
7     int statusCode, Header[] headers, String responseString) { /* [...] */ }
8                                     http://bit.ly/loopj_android_async_http_2FENovD

```

In the following example actual.data() returns a boxed Long. Because implicit upcasts have precedence over implicit unboxing conversions, the call is needed to invoke the method that takes a long (line 3) rather than the method that takes an Object (line 2).

```

1 assertEquals(expected, (long) actual.data());
2 public static void assertEquals(Object expected, Object actual) { /* [...] */ }
3 public static void assertEquals(long expected, long actual) { /* [...] */ }
4                                     http://bit.ly/spullara_redis_protocol_2FC9Llb

```

The following snippet is similar to the previous example, but notice how that the cast is applied to a primitive—*non-reference*—type.

```

assertEquals((byte) 0x1, record.getSpacing());
                                     http://bit.ly/apache_poi_2StrlOn

```

In the last example of SELECTOVERLOAD, an upcast of a generic type is performed to select the appropriate overload of the max method.

```

1 public static <T> T max(Iterator<T> self, Comparator<T> comparator) {
2     return max((Iterable<T>)toList(self), comparator);
3 }
4 public static <T> List<T> toList(Iterator<T> self) {
5     // [...]
6 }
7 @Deprecated
8 public static <T> T max(Collection<T> self, Comparator<T> comparator) {
9     // [...]
10 }
11 public static <T> T max(Iterable<T> self, Comparator<T> comparator) {
12     // [...]
13 }
                                     http://bit.ly/groovy_groovy_core_2HDAkbF

```

Detection. The Query 4.18 detects when a cast is used as an argument of an overloaded method. A cast returned by this query needs to be either a cast to null or an upcast. This is an approximation because the query does not check whether the overloaded method differs only on the type of the argument that is cast.

```

1  class SelectOverloadCast extends Cast {
2      SelectOverloadCast() {
3          (getExpr() instanceof NullLiteral or this instanceof Upcast) and
4          this instanceof OverloadedArgument
5      }
6      Callable getOverload() {
7          result = this.(OverloadedArgument).getAnOverload()
8      }
9  }

```

QL

Listing 4.18. Query to detect the SELECTOVERLOAD pattern.

Issues. Casting the null constant seems rather artificial. This pattern shows either a lack of expressiveness in Java or a bad API design. Passing null to a method might better be handled by using overloading with fewer parameters or by using default parameters. Several other languages support default parameters, *e.g.*, Scala, C# and C++. Adding default parameters might be a partial solution.

In addition, a pure object-oriented language would not distinguish between primitives and objects, avoiding the need for autoboxing to be visible at the type level.

Oostvogels et al. [2018] propose an extension to TypeScript to express constraints between properties, which can then be mapped onto optional parameters.

Both the ACCESSSUPERCLASSFIELD and this pattern are used to select class members. While this pattern is used to select the appropriate overloaded method, the ACCESSSUPERCLASSFIELD is used to select a field in a superclass.

4.5.10 KnownReturnType

Description. There are cases when a method's return type is less specific than the actual return type value. This is often to hide implementation

details, but may also be because the method overrides another method with a less-specific type and the return type is not changed covariantly.

This pattern is used to cast from the method's return type to the *known* actual return type. This pattern is characterized by a method that always returns a value of the same type, a subtype of the declared return type, regardless of the context or the arguments to the method call.

Instances: 89 (1.78%). We found 61 in application code, 23 in test code, and 5 in generated code. In the following example, a cast is performed to a call to the `getRealization` method (line 1). Its implementation returns a value of type `CubeInstance` (line 9).

```

1  final List<CubeSegment> mergingSegments = ((CubeInstance) seg.getRealization())
2      .getMergingSegments((CubeSegment) seg);
3  public class CubeSegment implements IBuildable, ISegment, Serializable {
4      // [...]
5      private CubeInstance cubeInstance;
6      // [...]
7      public IRealization getRealization() {
8          return cubeInstance;
9      }
10 }
11 public class CubeInstance
12     extends RootPersistentEntity implements IRealization, IBuildable {
13     // [...]
14 }
```

http://bit.ly/apache_kylin_2SljooO

In the following example, a cast is applied to the result of an invocation to the `createDebugTarget` method. This method is known to return a value of type `PHPDebugTarget`, which implements `IPHPDebugTarget`.

```

1  debugTarget = (PHPDebugTarget) createDebugTarget(/* [...] */);
2
3  protected IDebugTarget createDebugTarget(/* [...] */) throws CoreException {
4      return new PHPDebugTarget(/* [...] */);
5  }
```

http://bit.ly/eclipse_pdt_2Ekeu9v

In some situations, an API method is designed to return an abstract class or interface. This API allows the developer to then choose which implementation use at run-time. The following example shows this situation. The cast is applied to the `getLogger` method—with return type `org.slf4j.Logger`—in line 4. But the developer set up the application to use `ch.qos.logback.classic.Logger` instead.

```

1 import ch.qos.logback.classic.Logger;
2 import org.slf4j.LoggerFactory;
3
4 Logger rootLogger = (Logger) LoggerFactory.getLogger(Logger.ROOT_LOGGER_NAME);
5                                     http://bit.ly/skylot\_jadx\_2HIoR9X

```

Detection. Similar to the FACTORY pattern, KNOWNRETURN_TYPE requires analysis of the method implementation called in the cast expression. Expressing this kind of analysis in QL becomes impractical.

Issues. This pattern usually indicates an abstraction violation: the caller needs to know the method implementation to know the correct target type.

The COVARIANTRETURN_TYPE pattern can be considered a special case of this pattern where the return type is known to vary with the receiver type. Like that pattern, associated types [Chakravarty et al., 2005] in languages like Haskell or Rust could be used to avoid the cast.

4.5.11 Deserialization

Description. This pattern is used to deserialize an object at run-time. In its more common form, this pattern is characterized for a cast to the readObject method on an ObjectInputStream object.

Instances: 71 (1.42%). We found 37 in application code, 12 in test code, and 22 in generated code. The following example shows how the DESERIALIZATION pattern is used to create objects from a file system (line 9).

```

1 FileInputStream fis = new FileInputStream(serialize);
2 ObjectInputStream ois = new ObjectInputStream(fis);
3 CrawlURI deserializedCuri = (CrawlURI)ois.readObject();
4 deserializedCuri = (CrawlURI)ois.readObject();
5 deserializedCuri = (CrawlURI)ois.readObject();
6 assertEquals("...", this.seed.toString(), deserializedCuri.toString());
7                                     http://bit.ly/internetarchive\_heritrix3\_2SF4j7k

```

Detection. The following query detects DESERIALIZATION with the fact that the readObject method family is used to deserialize objects. For other

deserialization frameworks, it would require to analyse external dependencies.

```

1  class DeserializationCast extends Cast {
2      DeserializationCast() {
3          getExprOrDef().(MethodAccess).getMethod() instanceof ReadObjectMethod
4      }
5  }
```

Listing 4.19. Detection of the DESERIALIZATION pattern.

Issues. The serialization API dates back to Java 1.1 in 1997. Since then, newer serialization APIs have been developed. For instance, Apache Avro³³ uses generics and class literals to specify the expected type of an object read. In some languages, type-safe serialization and deserialization boilerplate code can be automatically generated, for instance in Rust, the Serde library³⁴ can generate code to serialize most data types in a variety of formats.

Both this pattern and the NEWDYNAMICINSTANCE pattern create objects by using reflection. While it might be considered a special case of KNOWN-RETURN-TYPE, DESERIALIZATION differs in that the run-time result type of the readObject depends on the state of the input stream and can change depending on context.

4.5.12 VariableSupertype

Description. This pattern occurs when a cast is applied to a variable (local variable, parameter, or field), that has usually been assigned just once and is declared with a proper supertype of the value assigned into it. The type of the value being assigned to can be determined locally either within the enclosing method or class.

To detect this pattern, a cast needs to be applied to a variable whose value can be determined simply by looking at the enclosing method or class.

Instances: 64 (1.28%). We found 53 in application code, 8 in test code, and 3 in generated code. The following snippet shows an example of

³³<https://avro.apache.org/docs/current/>

³⁴<https://serde.rs/>

the `VARIABLESUPERTYPE` pattern (line 4). The `samlTokenRenewer` variable is being cast to the `SAMLTokenRenewer` class. The variable is declared with type `TokenRenewer` (superclass of `SAMLTokenRenewer`) in line 1. However, the variable is being initialized with the expression `new SAMLTokenRenewer()`. Thus, the cast instance could be trivially avoided by changing the declaration of the `samlTokenRenewer` variable to `SAMLTokenRenewer` instead of `TokenRenewer`.

```

1 TokenRenewer samlTokenRenewer = new SAMLTokenRenewer();
2 samlTokenRenewer.setVerifyProofOfPossession(false);
3 samlTokenRenewer.setAllowRenewalAfterExpiry(true);
4 ((SAMLTokenRenewer)samlTokenRenewer).setMaxExpiry(1L);
5
http://bit.ly/apache\_cxf\_2SNoUXj

```

The following listing shows an example of the `VARIABLESUPERTYPE` pattern. The field `uncompressedDirectBuf` is being cast to the `java.nio.ByteBuffer` class (line 13) but it is declared as `java.nio.Buffer` (line 3). Nevertheless, the field is assigned only once in the constructor (line 7) with a value of type `java.nio.ByteBuffer`. The value assigned is returned by the method `allocateDirect` from the `ByteBuffer` class.³⁵ Inspecting the enclosing class, there is no other assignment to the `uncompressedDirectBuf` field, thus making possible to declare it as `final`. Therefore, the cast pattern in line 13 will always succeed. Any other similar use of the `uncompressedDirectBuf` field needs to be cast as well.

```

1 public class SnappyCompressor implements Compressor {
2     // [...]
3     private Buffer uncompressedDirectBuf = null;
4     // [...]
5     public SnappyCompressor(int directBufferSize) {
6         // [...]
7         uncompressedDirectBuf = ByteBuffer.allocateDirect(directBufferSize);
8         // [...]
9     }
10    // [...]
11    synchronized void setInputFromSavedData() {
12        // [...]
13        ((ByteBuffer) uncompressedDirectBuf).put(userBuf, userBufOff,
14            uncompressedDirectBufLen);
15        // [...]
16    }

```

³⁵[https://docs.oracle.com/javase/7/docs/api/java/nio/ByteBuffer.html#allocateDirect\(int\)](https://docs.oracle.com/javase/7/docs/api/java/nio/ByteBuffer.html#allocateDirect(int))

```

17     // [...]
18 }

```

http://bit.ly/facebookarchive_hadoop_20_2FuDeO7

In the next cast instance, the parameter `k1` is cast to the `Comparable` class (line 7). `k1` is declared as `E` (line 5), an unbounded type parameter (line 1). The developer likely designed the class so that `E` must be `Comparable` only if `comparator` is `null`, providing an API with two ways to compare list elements.

```

1 public class SortedArrayList<E> extends ArrayList<E> {
2     protected final Comparator<E> comparator;
3     // [...]
4     @SuppressWarnings( {"unchecked"})
5     protected int compare(final E k1, final E k2) {
6         if (comparator == null) {
7             return ((Comparable) k1).compareTo(k2);
8         }
9         return comparator.compare(k1, k2);
10    }
11 }

```

http://bit.ly/oblaac_jodd_2UKxm6H

In the next example, the `ir` field is cast to `DirectoryReader` (line 11). The `ir` field is declared as `IndexReader` (superclass of `DirectoryReader`) in line 1. The cast to `ir` is performed using the value of the expression `readers.get(0)` (line 10). But `readers` is defined as `ArrayList<DirectoryReader>` (line 3), making the cast superfluous if an extra variable of type `DirectoryReader` had been used.

```

1 private IndexReader ir = null;
2 // [...]
3 ArrayList<DirectoryReader> readers = new ArrayList<DirectoryReader>();
4 for (Directory dd : dirs) {
5     DirectoryReader reader;
6     reader = DirectoryReader.open(dd);
7     readers.add(reader);
8 }
9 if (readers.size() == 1) {
10     ir = readers.get(0);
11     dir = ((DirectoryReader)ir).directory();
12 } else {
13     ir = new MultiReader(
14         (IndexReader[])readers.toArray(new IndexReader[readers.size()]);
15     }

```

http://bit.ly/tarzanek_luke_2OhDT6O

Detection. This pattern contains many variations that require manual inspection. To detect this pattern an interprocedural data-flow analysis would be required, since the value being cast could be assigned in another method.

However, for some cases (*e.g.*, the first example), some query approximations would be possible, since the cast and the assignment are in the same method. The following query shows how to approximate the detection in this case. The `forex` quantifier asserts that for every `Type t = getADef().getType()`, then `t = getTargetType()`, and that at least exists one `t` satisfying the condition, *i.e.*, `t = getTargetType()`.

```

1 class VariableSupertypeCast extends VarCast {
2     VariableSupertypeCast() {
3         forex (Type t | t = getADef().getType() | t = getTargetType()) and
4             isSubtype(getTargetType(), getExpr().getType())
5     }
6 }

```

QL

Listing 4.20. Detection of the VARIABLESUPERTYPE pattern.

Issues. In most the cases this can be considered as a bad practice or code smell. This is because by only changing the declaration of the variable to a more specific type, the cast can be simply eliminated.

This pattern sometimes related to the REDUNDANT pattern. Although VARIABLESUPERTYPE is not redundant, by only changing the declaration of the variable to a more specific type, the cast becomes redundant.

4.5.13 SoleSubclassImplementation

Description. The SOLESUBCLASSIMPLEMENTATION occurs when an interface or abstract class has only one implementing subclass. Casting the interface to this class must succeed because it cannot possibly be an instance of another class.

Instances: 61 (1.22%). We found 30 in application code, 6 in test code, and 25 in generated code. In the following example the `jobId` variable is cast to the sole implementation (`JobIdImpl`).


```
return Longs.compare(id, ((JobIdImpl) jobId).id);
```

http://bit.ly/ow2_proactive_scheduling_2Ulqfs

Similar to the previous example, the variable `user` is cast to the known implementation (`UserImpl`).

```
1 for (User user : api.getUsers()) {
2   if (channelId.equals(((ImplUser) user).getUserChannelId())) {
3     return user;
4   }
5 }
```

http://bit.ly/Javacord_Javacord_2GwGjuV

Detection. The following query returns all casts such that the type—class or interface—of the expression being cast has only one subtype. The *transitive closure* symbol `+` indicates that `getASubtype` may be followed one or more times.

```
1 class SoleSubclassImplementation extends Cast {
2   SoleSubclassImplementation() {
3     count(RefType rt |
4       rt = getExpr().getType() and rt.fromSource() |
5       rt.getASubtype+() ) = 1
6   }
7 }
```

QL

Listing 4.21. Detection of the SOLESUBCLASSIMPLEMENTATION pattern.

Issues. This pattern occurs when there is high cohesion between super and subclass. In some cases, the cast instance appears in a generated class. This mechanism allows the developer to extend this class to add custom code. Therefore, this high cohesion is acceptable. The developer assumes that there is no other implementation of the base class, otherwise the cast instance fails.

4.5.14 NewDynamicInstance

Description. In the NEWDYNAMICINSTANCE pattern, a new object or array is created by means of reflection. The type of the object being created is determined at run time, and the new object is cast to some statically known supertype of the run time type.

The `newInstance` method family declared in the `Class`,³⁶ `Array`^{37, 38} and `Constructor`³⁹ classes creates an object or array dynamically by means of reflection, *i.e.*, the type of object being created is not known at compile-time. This pattern consists of casting the result of these methods to the appropriate target type.

Instances: 59 (1.18%). We found 44 in application code, 5 in test code, and 10 in generated code. The following example shows a cast of the result of the `Class.newInstance()` method.

```
logger = (AuditLogger) Class.forName(className).newInstance();
```

http://bit.ly/apache_hadoop_2HC3IPg

The following example shows how to dynamically create an array, using the `Array` class.

```
1 return list.toArray( (T[]) Array.newInstance( componentType, list.size()));
```

2 http://bit.ly/neo4j_neo4j_2Hp5Hqc

Whenever a constructor other than the default constructor is needed, the `newInstance` method declared in the `Constructor` class should be used to select the appropriate constructor, as shown in the following example.

```
1 return (Exception) Class.forName(className)
```

```
2     .getConstructor(String.class)
```

```
3     .newInstance(message);
```

4 http://bit.ly/gradle_gradle_2HsUgOo

The following example shows a guarded instance of the `NEW_DYNAMIC_INSTANCE` pattern. This seems rather unusual, as this pattern is not guarded.

```
1 private static List<String> getMapperMethodNames(final Class clazz) {
```

```
2     try {
```

```
3         if (clazz != null) {
```

³⁶<https://docs.oracle.com/javase/8/docs/api/java/lang/Class.html#newInstance-->

³⁷<https://docs.oracle.com/javase/8/docs/api/java/lang/reflect/Array.html#newInstance-java.lang.Class-int->

³⁸<https://docs.oracle.com/javase/8/docs/api/java/lang/reflect/Array.html#newInstance-java.lang.Class-int...->

³⁹<https://docs.oracle.com/javase/8/docs/api/java/lang/reflect/Constructor.html#newInstance-java.lang.Object...->

```

4         Object obj = clazz.newInstance();
5         if (obj instanceof BaseMethodMapper) {
6             return ((BaseMethodMapper) obj).getAllFunctionNames();
7         }
8     }
9 } catch (Exception e) {
10     e.printStackTrace();
11 }
12 return null;
13 }
```

http://bit.ly/alibaba_LuaViewSDK_2HC33xg

There are cases when the cast is not directly applied to the result of the `newInstance` method. The following snippet shows such a case. The cast is used to convert from `Class<?>` to `Class<ConfigFactory>` (line 4). The invocation to the `newInstance` method then does not need a direct cast (line 8) given the definition of the `clazz` variable (line 2). Nevertheless, the cast is unchecked, and a `checkcast` instruction is going to be emitted anyway for the result of the `newInstance` invocation.

```

1 ClassLoader tccl = Thread.currentThread().getContextClassLoader();
2 final Class<ConfigFactory> clazz;
3 if (tccl == null) {
4     clazz = (Class<ConfigFactory>) Class.forName(factoryName);
5 } else {
6     clazz = (Class<ConfigFactory>) Class.forName(factoryName, true, tccl);
7 }
8 final ConfigFactory factory = clazz.newInstance();
9
```

http://bit.ly/pac4j_pac4j_2HJtXUn

Detection. The Query 4.22 detects the `NEWDYNAMICINSTANCE` pattern. The QL `NewDynamicInstanceAccess` class checks where the cast expression is the `newInstance` method on the aforementioned classes.

```

1 class NewDynamicInstanceCast extends Cast {
2     NewDynamicInstanceCast() {
3         getExprOrDef() instanceof NewDynamicInstanceAccess
4     }
5 }
```

QL

Listing 4.22. Detection of the `NEWDYNAMICINSTANCE` pattern.

Issues. The cast here is needed because of the dynamic nature of reflection. This pattern is usually unguarded; that is, the programmer knows what target type is being created.

Generics could be used to avoid the cast on `newInstance`, assuming the `Class` instance is not a raw type or a `Class<?>`. However, the usual API for getting a class instance `Class.forName` returns such a type. Indeed, the following two snippets:

```
Class<?> c = Class.forName("java.lang.String");
String pf = (String) c.newInstance();
```

```
Class<String> c = (Class<String>) Class.forName("java.lang.String");
String pf = c.newInstance();
```

compile to the same bytecode below.

```
ldc          "java.lang.String"
invokestatic java.lang.Class.forName()
astore_1
aload_1
invokevirtual java.lang.Class.newInstance()
checkcast    java.lang.String
```

Bytecode

In the first case, the cast is to the `newInstance` method, an instance of the `NEW_DYNAMIC_INSTANCE` pattern. In the second case, the cast is to the call to `Class.forName`, an instance of the `FACTORY` pattern.

This pattern is related to `DESERIALIZATION`, since both create an object dynamically. It is also related to `REFLECTIVE_ACCESSIBILITY`, where both retrieve objects by using reflection.

4.5.15 ObjectAsArray

Description. In this pattern an array is used as an untyped object. A cast is applied to a constant array slot, *e.g.*, `(String) array[1]`.

Instances: 47 (0.94%). We found 36 in application code, 10 in test code, and 1 in generated code. The following example shows an instance of this pattern. The variable `currentState` contains an `Object[]` with a fixed

schema.⁴⁰ Then, a cast is performed of a constant array slot, (BitSet) currentState[3] on line 5.

```

1      BitSet theLoadedFields = (BitSet)currentState[2];
2      for (int i = 0; i < this.loadedFields.length; i++) {
3          this.loadedFields[i] = theLoadedFields.get(i);
4      }
5      BitSet theModifiedFields = (BitSet)currentState[3];
6      for (int i = 0; i < dirtyFields.length; i++) {
7          dirtyFields[i] = theModifiedFields.get(i);
8      }
9      setVersion(currentState[1]);
10

```

http://bit.ly/datanucleus_datanucleus_core_2S1L5Zf

Detection. The following query detects when the cast expression is an array access, and that access is indexed with a compile-time constant.

```

1  class ObjectToArrayCast extends Cast {
2      ArrayAccess arr;
3      ObjectToArrayCast() {
4          arr = getExprOrDef() and
5          arr.getIndexExpr() instanceof CompileTimeConstantExpr and
6          arr.getArray().getType().(Array).
7              getElementType() instanceof TypeObject
8      }
9  }

```

QL

Listing 4.23. Detection of the OBJECTARRAY pattern.

Issues. This pattern usually suggests an abuse of the type system. Using an object with statically typed fields might be a better alternative.

4.5.16 ImplicitIntersectionType

Description. This pattern occurs when there is a downcast of reference v of type T to a target interface type I . Although T does not implement I , the cast succeeds because all possible run-time types of v do implement I .

⁴⁰<http://www.datanucleus.org/javadocs/core/5.0/org/datanucleus/enhancement/Detachable.html>

Instances: 45 (0.90%). We found 19 in application code, 21 in test code, and 5 in generated code. For instance, in the following example the method call returns a `Number`, which does not implement `Comparable`; however, all values that could be returned by the method are subclasses of `Number` in `java.lang` that do implement `Comparable`.

```
final Comparable max = (Comparable) properties.getMaxValue();
```

http://bit.ly/senbox_org_snap_desktop_2FQOt4v

This pattern can be used to implement a dynamic proxy. In the following example, `pyObjectValue` is a proxy to `PyObjectValue`. Nevertheless, a cast to `Proxy` is needed to invoke the `setHandler`.

```
1 PyObjectValueProxyClass proxyClass = getProxyClass(pyObject);
2 PyObjectValue pyObjectValue = (PyObjectValue) proxyClass.getConstructor()
3     .newInstance(proxyClass.getParams());
4 ((Proxy) pyObjectValue).setHandler(
5     new PyObjectValueMethodHandler(content, sensitive, pyObject));
6
```

http://bit.ly/CloudSlang_cloud_slang_2Ekp4l

Detection. The Query 4.24 detects this pattern. Usually, in a downcast (T) e , the class or interface T is a subtype of e 's class or interface. This query essentially detects whether T has no subtyping relation with the type of e .

```
1 class ImplicitIntersectionTypeCast extends Cast {
2     ImplicitIntersectionTypeCast() {
3         getTargetType() instanceof Interface and
4         not isSubtype(getTargetType(), getExpr().getType()) and
5         not this instanceof Upcast and
6         not getExpr() instanceof NullLiteral and
7         notGenericRelated(getTargetType()) and
8         notGenericRelated(getExpr().getType())
9     }
10 }
```

QL

Listing 4.24. Detection of the `IMPLICITINTERSECTIONTYPE` pattern.

Issues. The cast could be avoided by having the operand type implement the target type interface or by introducing a more precise interface. In the

first example, one could imagine an interface `ComparableNumber` that extends both `Number` and `Comparable`. Scala supports interface types, allowing the type `Number with Comparable` to be used directly.

Fourtounis et al. [2018] propose a static analysis of dynamic proxies, which are a special case of this pattern. To implement their analysis, they have used Doop [Bravenboer and Smaragdakis, a].

4.5.17 RemoveWildcard

Description. A cast is in the `REMOVEWILDCARD` pattern when a *wildcard type* is used rather than a generic type.

Instances: 34 (0.68%). We found 26 in application code, 7 in test code, and 1 in generated code. In the following example, `unit` is declared as `Unit<?>`, but to actually be able to use it a cast to a concrete type is needed.

```
copy.setUnitOfMeasure( (Unit<Length>) unit );
```

http://bit.ly/eclipse_jetty_project_2WMI0Ld

Detection. The following query detects the `REMOVEWILDCARD` pattern. The query checks that the type of the cast operand is a wildcard, or a parameterized type containing a wildcard.

```

1 predicate containsWildcard(Type t) {
2     t instanceof Wildcard or
3     containsWildcard( t.(ParameterizedType).getATypeArgument() )
4 }
5
6 class RemoveWildcardCast extends Cast {
7     RemoveWildcardCast() {
8         containsWildcard(getExpr().getType())
9     }
10 }
```

QL

Listing 4.25. Detection of the `REMOVEWILDCARD` pattern.

Issues. Wildcard types are a form of existential type and consequently can limit access to members of a generic type. Casts are used to restore access at a particular type.

Since this pattern is an unchecked cast, the discussion about compiler-inserted casts and *blame* is similar to the `USERAWTYPE` pattern.

4.5.18 OperandStack

Description. The `OPERANDSTACK` pattern consists of multiple cases, dispatched depending on some application-specific control state, with casts of the top elements of stack-like collection in each case. An application invariant ensures that if the application is in a given state then the top elements of the stack should be of known run-time types.

Instances: 29 (0.58%). We found 13 in application code, 0 in test code, and 16 in generated code. The following example, shows a cast whose value is on top of a stack (line 2). In this case, the code is transforming a parse tree into an abstract syntax tree. The casts in the switch case are guarded by the parse tree node type and its arity.

```

1  case JJTASSERT_STMT:
2      exprType msg = arity == 2 ? ((exprType) stack.popNode()) : null;
3      test = (exprType) stack.popNode();
4      return new Assert(test, msg);
5
http://bit.ly/fabioz\_Pydev\_2HF6nrF

```

Similar to the previous example, in this case a guarded cast is performed on a stack of grammar symbols. The code was generated using an LR parser generator. The guard ensures that the parser has already matched a given prefix of the input and so the top of the stack should contain the expected symbols.

```

1  case 40: // qualified_name_decl = name_decl.n DOT.DOT IDENTIFIER.i
2  {
3      final Symbol _symbol_n = _symbols[offset + 1];
4      final IdUse n = (IdUse) _symbol_n.value;
5      final Symbol DOT = _symbols[offset + 2];
6      final Symbol i = _symbols[offset + 3];
7      return new IdUse(n.getID() + "." + ((String)i.value));
8  }
http://bit.ly/Sable\_soot\_2MZLZ3m

```

Detection. To manually detect this pattern, we look for methods that pop up an element from a stack, and then cast to it. Automatic detection for

this pattern becomes impractical, since a query would need to detect such a method, and when a class is implementing a stack-like structure.

Issues. In our sample, this pattern is always seen when implementing grammar-related operations, such as parsers or interpreters. In some situations, similar to the STASH pattern, this pattern could be replaced with a strongly typed heterogeneous collection [Kiselyov et al., 2004].

Similar to TYPECASE, multiple cases are evaluated with casts to different types, depending on application-specific guards. However, unlike TYPECASE, the success of the casts is ensured not by a type-tag-like value, but by application-specific state (*e.g.*, the current parser state or the state of an evaluator) and proper use of the stack.

4.5.19 ReflectiveAccessibility

Description. This pattern accesses a field of an object by means of reflection. Typically reflection is used because the field is private and therefore inaccessible at compile time and the developer cannot change the field declaration itself. In this case, the method `Field::setAccessible(true)` is invoked on the field before getting the value of the field. The cast is needed because `Field::get` returns an `Object`.

Instances: 27 (0.54%). We found 22 in application code, 5 in test code, and 0 in generated code. The following two snippets show how this pattern is used:

```

1 f.setAccessible(true);
2 HttpEntity wrapped = (HttpEntity) f.get(entity);
3
    http://bit.ly/loopj_android_async_http_2SOISRr

1 Field fieldPosition=ChangesOutputter.class.getDeclaredField("changesPosition");
2 fieldPosition.setAccessible(true);
3 ChangesOutputter changesDisplayBis = output(changes);
4 PositionWithChanges<ChangesAssert, ChangeAssert> positionBis =
5     (PositionWithChanges) fieldPosition.get(changesDisplayBis);
6
    http://bit.ly/joel_costigliola_assertj_db_2Ip1Rho

```

Detection. The Query 4.26 detects this pattern. The query looks for a cast applied to a get or invoke method in an object *o* of type `Field` or `Method` respectively. Moreover, it checks that `setAccessible(true)` has been invoked in *o*. However, this query does not check that `setAccessible(true)` has been invoked *before* the cast.

```
1 class ReflectiveAccessibilityCast extends Cast {
2     Variable fieldVariable;
3     ReflectiveMethodAccess reflectiveMethodAccess;
4     SetAccessibleTrueMethodAccess satma;
5     ReflectiveAccessibilityCast() {
6         reflectiveMethodAccess = getExprOrDef() and
7         fieldVariable.getAnAccess() =
8             getExprOrDef().(MethodAccess).getQualifier().(VarAccess) and
9         fieldVariable.getAnAccess() = satma.getQualifier()
10    }
11 }
```

Listing 4.26. Detection of the REFLECTIVEACCESSIBILITY pattern.

Issues. Using reflection to access a field is a common workaround to tight access control restrictions. However, it should generally be regarded as a code smell.

As with `DESERIALIZATION`, this pattern is necessary because a library method can return values of many different types at run time, and so is declared to return `Object`.

4.5.20 FluentAPI

Description. A fluent API is an API that allows the developer to operate on the same object using method chaining. This pattern is exhibited when the receiver (`this` reference) is cast to a type parameter which is itself bounded by the self type.

Instances: 23 (0.46%). We found 18 in application code, 0 in test code, and 5 in generated code. In the following snippet, the receiver (`this`) is cast to a type parameter (`B`) (line 5). This allows subclasses to reuse the methods in the base class without overriding them just to change the return type.

```

1 public class ClockBuilder <B extends ClockBuilder<B>> {
2     // [...]
3     public final B alarms(final Alarm... ALARMS) {
4         properties.put("alarmsArray", new SimpleObjectProperty<>(ALARMS));
5         return (B) this;
6     }
7 }

```

http://bit.ly/HanSolo_Medusa_2TyBObH

The following example implements FLUENTAPI by directly casting the receiver (this) in line 3. Similarly to the addAllThrown method, the rest of methods in the enclosing class perform a cast to this. Although there is a lot of boilerplate code, this instance happens in generated code. The cast succeeds because there is a guard (line 8) in the constructor that guarantees the receiver is of the appropriate type (*cf.* TYPECASE).

```

1 public final EncodedElement.Builder addAllThrown(
2     Iterable<? extends Type> elements) {
3     this.thrown.addAll(elements);
4     return (EncodedElement.Builder) this;
5 }
6
7 public Builder() {
8     if (!(this instanceof EncodedElement.Builder)) {
9         throw new UnsupportedOperationException("/* [...] */");
10    }
11 }

```

http://bit.ly/immutableables_immutableables_2S4BoJs

Detection. The Query 4.27 detects the FLUENTAPI pattern. The query detects the case like the first example.

```

1  class FluentAPICast extends Cast {
2      TypeVariable x;
3      GenericType enclosingClass;
4      FluentAPICast() {
5          getExpr() instanceof ThisAccess and
6          getParent() instanceof ReturnStmt and
7          x = getTargetType() and
8          enclosingClass = getExpr().getType() and
9          x.hasTypeBound() and
10         x.getFirstTypeBound().getType() = enclosingClass and
11         x.getFirstTypeBound().getType().(GenericType).getATypeParameter() = x
12     }
13 }

```

Listing 4.27. Detection of the FLUENTAPI pattern.

Issues. In most cases, this pattern is concerned with a particular implementation of fluent APIs where recursive generics are used to mimic self types [Bruce, 2003]. Other implementations of fluent APIs simply return this without a cast, but these are less extensible.

4.5.21 CovariantGeneric

Description. The COVARIANTGENERIC pattern occurs when a cast is used to use an invariant generic type as if it were covariant.

Instances: 22 (0.44%). We found 13 in application code, 9 in test code, and 0 in generated code. In the following snippet, an upcast is performed to ensure that the inferred type of the call to `singletonList` is a supertype of the type that would be otherwise inferred. The `singletonList` method has the signature `<T> List<T> singletonList(T o)`.⁴¹ If `curframe` were passed in without the cast, the type of the list would be inferred to be `List<FrameBuilder>`, which is not a subtype of the method return type `List<Framedata>`, causing a compilation error. With the cast, the list type is inferred to be the same as the return type.

```

1  @Override
2  public List<Framedata> createFrames(String text, boolean mask) {
3      FrameBuilder curframe = new FramedataImpl1();

```

⁴¹<https://docs.oracle.com/javase/8/docs/api/java/util/Collections.html>

```

4      /* [...] */
5      return Collections.singletonList( (Framedata) curframe );
6  }
7  public interface FrameBuilder extends Framedata { /* [...] */ }
8  public class Collections {
9      public static <T> List<T> singletonList(T o) { /* [...] */ }
10 }

```

http://bit.ly/arpruss_raspberryjammod_2USL7Ai

Similar to the previous example, in the following case, an upcast is performed to change the return type of the `Matcher<T> equalTo(T)` method.

```

1  @Test
2  public void testUpdateReturnBoolean() throws Exception {
3      /* [...] */
4      List<Object> args = boundSql.getArgs();
5      assertThat(args.get(0), equalTo((Object) "ash"));
6  }
7  public static <T> Matcher<T> equalTo(T operand) {
8      // [...]
9  }

```

http://bit.ly/jfaster_mango_2EhXzUW

Instead of an upcast, in this example, a cast to null is performed to change the return type. This use case resembles the `SELECTOVERLOAD` pattern.

```

1  assertThat(result.queryValue(memberOne, DefaultFlag.BUILD), is((State) null));
2  public static <T> Matcher<T> is(T value) {
3      // [...]
4  }

```

http://bit.ly/EngineHub_WorldGuard_2IVUOx1

Another common version of this pattern for type `S` a subtype of `T`, is to cast a generic type like `List<S>` to a raw type (`List`), which can then be assigned to a variable of `List<T>`.

```

1  private final List<VariableExpression> dataProcessorVars = new ArrayList<>();
2  new ArrayExpression(ClassHelper.OBJECT_TYPE, (List) dataProcessorVars);
3  public class ArrayExpression extends Expression {
4      public ArrayExpression(ClassNode elementType, List<Expression> exprs) {}
5  }

```

http://bit.ly/spockframework_spock_2UYEsF5

Detection. The Query 4.28 detects when a cast is used to select the return type of a generic method. This query *does not* detect a cast belonging to this pattern when the raw type is used, *e.g.*, last example shown above.

```

1  class CovariantGenericCast extends Cast {
2      Argument arg;
3      Call call;
4      Callable m;
5      CovariantGenericCast() {
6          this = arg and
7          call = arg.getCall() and
8          arg.getCall().getCallee() = m and
9          (
10             m.getReturnType().(ParameterizedType).getATypeArgument() =
11                 m.getParameterType(arg.getPosition()).(TypeVariable) or
12             m.getReturnType().(TypeVariable) =
13                 m.getParameterType(arg.getPosition()).(TypeVariable)
14         )
15     }
16 }

```

Listing 4.28. Query to detect the COVARIANTGENERIC pattern.

Issues. In some cases, this pattern could be avoided using explicit type parameter, *e.g.*, `Collections.<Framedata>singletonList(curframe)`. From Java 8 this cast is unnecessary due to better type inference.⁴²

Altidor et al. [2011] define a type system that infers definition-site variance to Java. This can reduce the need for this pattern.

4.5.22 Composite

Description. The COMPOSITE pattern is characterized by a cast to another element of a composite data structure, typically a tree, where the target type is known because of its position in the data structure.

Instances: 21 (0.42%). We found 21 in application code, 0 in test code, and 0 in generated code. The following example shows a cast from a `Box`—as returned by the `getPreviousSibling` method—to a `TableSectionBox`. The programmer reasons that the cast will succeed because the source of the cast is a sibling of another `TableSectionBox`.

```

1  public class TableBox extends BlockBox {
2      protected TableSectionBox sectionAbove(TableSectionBox section /*[...]*/) {
3          TableSectionBox prevSection = (TableSectionBox)section.

```

⁴²<https://docs.oracle.com/javase/specs/jls/se8/html/jls-18.html#jls-18.5>

```
4         getPreviousSibling();
5     }
6 }
7 public abstract class Box implements Styleable {
8     public Box getPreviousSibling() { }
9 }
```

http://bit.ly/flyingsaucerproject_flyingsaucer_2N2nYbY

Detection. Since the COMPOSITE pattern resembles the FAMILY pattern, its detection suffers the same inconveniences.

Issues. The pattern is typical of hierarchical data structures such as abstract syntax trees, document models, or UI layouts. Based on the grammar of the data structure, the types of adjacent objects in the structure can be known. The cast succeeds if the data structure is well-formed. This pattern is only seen in application code, since it is used when designing an extensible API.

More precise typing of the links in the data structure could eliminate the need for the casts. For example, in the above example, the sibling of a `TableSectionBox` might be declared to have type `TableSectionBox`. However, this may require the programmer to override methods to refine return types covariantly. Language features available in other languages like generalized algebraic data types (GADTs) [Peyton Jones et al., 2006] or self types [Bruce, 2003; Odersky and Zenger, 2005] could also be used to provide a more precise typing.

The pattern can be thought of as a more dynamic variant of the FAMILY pattern. Rather than reasoning that the cast will succeed because of the source type's relative position in the class hierarchy, the cast will succeed because of the source value's position in a composite data structure.

4.5.23 GenericArray

Description. A cast due to the instantiation of an array with a parameterized base type. In Java these arrays cannot be instantiated, instead an `Object[]` or an array of raw types must be created. The cast is necessary to use the array at the intended type.

Instances: 7 (0.14%). We found 5 in application code, 2 in test code, and 0 in generated code. In the following snippet, a cast is required when

accessing an element in the array (line 4). The array is created using the raw type `List[][]` and assigned to a variable of using the wildcard type `List<?>[][]` (line 1). It is not possible to simply allocate a `List<byte[]>`.

```

1 List<?>[][] partialResults = new List<th>[tw];
2 for (...) {
3     partialResults[ty][tx] = build(tx, ty, order, cCompatibility);
4     layers.addAll((List<byte[]>) partialResults[y][x]);
5 }

```

http://bit.ly/ppiastucki_recast4j_2EM7zWK

Instead of casting individual elements, the following example shows a cast applied directly when the array is created.

```

T[] newArray = (T[]) new Object[growSize(currentSize)];

```

http://bit.ly/seven332_Nimingban_2UdBwIL

Detection. The following queries detect different variations of the `GENERICARRAY` pattern. The first one detects when a generic cast is applied to the array, *e.g.*, `(E[]) new Object[length]`.

```

1 class OnArrayGenericArrayCast extends Cast {
2     OnArrayGenericArrayCast() {
3         getTargetType().(Array).getComponentType() instanceof TypeVariable and
4         getExpr().getType() instanceof Array
5     }
6 }

```

QL

The following query detects the case when the target type of the cast is a type variable used to get an element from the array. For instance, `(T) events[i]`, where the `events` array is defined as `EventObject[]` and `T` is declared as `T extends EventObject`.

```

1 class TypeVariableGenericArrayCast extends Cast {
2     TypeVariableGenericArrayCast() {
3         getExprOrDef() instanceof ArrayAccess and
4         getTargetType() instanceof TypeVariable
5     }
6 }

```

QL

Our last query for this pattern is similar to the previous one. But in this case, the component type of the array is either a raw type or a wildcard type (e.g., `List<?>`). For instance, a cast `(Any<T>) entries[i]` where `entries` is defined as `Any[] entries = new Any[n]`. In this example, `Any` is a generic type, but the array is using the raw type instead.

```

1  class OnElementGenericArrayCast extends Cast {
2      OnElementGenericArrayCast() {
3          (
4              getExprOrDef().(ArrayAccess).getArray().getType().(Array)
5                  .getComponentType() instanceof RawType or
6                  containsWildcard( getExprOrDef().(ArrayAccess).getArray()
7                      .getType().(Array).getComponentType() )
8              ) and
9              getTargetType() instanceof ParameterizedType
10     }
11 }

```

Issues. This pattern occurs because generic type parameters are not reified at runtime, but array types are reified. To create a generic `T[]`, for instance, since the parameter `T` is not known statically, the compiler cannot know the runtime representation of the array. The Java specification just forbids these problematic cases and therefore requires programmers to create arrays of raw types and to use casts.

4.5.24 AccessSuperclassField

Description. Perform an upcast to access a field of a superclass of the cast operand.

Instances: 4 (0.08%). We found 0 in application code, 0 in test code, and 4 in generated code. The following snippet shows an instance of this pattern.

```

1  public abstract class StudentsPerformanceReport_Base extends QueueJobWithFile {
2      // [...]
3      public ExecutionSemester getValue(StudentsPerformanceReport o1) {
4          return ((StudentsPerformanceReport_Base)o1).executionSemester.get();
5      }
6      private OwnedVBox<ExecutionSemester> executionSemester;
7  }
8  public class StudentsPerformanceReport extends StudentsPerformanceReport_Base {

```

```

9      // [...]
10     }

```

http://bit.ly/FenixEdu_fenixedu_academic_2SQxlkC

Detection. The Query 4.29 detects the two variants of this pattern. The first variant, as shown in the example below, is when an upcast is performed to access a field when the subclass does not have access privileges to access the field. In this case, the field is declared as either private or protected in the superclass. On the other variant—not found in our manual sample—an upcast is performed to access a field declared in a superclass, when the subclass declares a field with the same name.

```

1  class AccessSuperclassFieldCast extends Cast {
2      FieldAccess fa;
3      AccessSuperclassFieldCast() {
4          this instanceof Upcast and
5          fa.getQualifier().getProperExpr() = this and (
6          getExpr().getType().(RefType).declaresField(fa.getField().getName()) or
7          ( fa.getField().isPrivate() or fa.getField().isProtected() )
8          )
9      }
10 }

```

QL

Listing 4.29. Detection of the ACCESSSUPERCLASSFIELD pattern.

As in our first example, the following snippet shows an example of the second variant mentioned above.

```

1  private SomeObject from = new SomeObject(100);
2
3  assertThat(((InheritMe) to).privateInherited)
4      .isNotEqualTo(((InheritMe) from).privateInherited);
5
6  static class InheritMe {
7      protected String protectedInherited = "protected";
8      private String privateInherited = "private";
9  }
10 public static class SomeObject extends InheritMe {
11 }

```

http://bit.ly/mockito_mockito_2vF51Em

Issues. The particular instance we encountered has a method whose parameter is a subclass of the current class. The cast is needed to access a

private field of the current class. Being an upcast, the cast is always safe. More problematic is the strong coupling between the base class and the derived class, however the base class is generated code; possibly, a manually written version would just combine the two classes.

Another use of the pattern, not found in our sample however, is to upcast a value to access a field of a superclass which is shadowed by another field of the same name in the subclass.

The REFLECTIVEACCESSIBILITY pattern is also used to access private or protected fields, albeit fields of unrelated classes that cannot be accessed simply by casting to another type. Like SOLESUBCLASSIMPLEMENTATION, this pattern occurs when there is high cohesion between super and subclass.

Both the SELECTOVERLOAD and this pattern are used to select class members. While this pattern is used to select a field in a superclass, the SELECTOVERLOAD is used to select the appropriate overloaded method.

4.5.25 UnoccupiedTypeParameter

Description. This pattern occurs when a generic type changes its type parameter, but the new type parameter holds no values.

Instances: 1 (0.02%). We found 1 in application code, 0 in test code, and 0 in generated code. This instance is used to implement an Either type. A value of type Either<L, R> can be either a value of type L or of type R. In this instance, the receiver—of type Either<L, R>—is cast to Either<U, R> (line 9). There is no subtype relation between L and U. However, the cast succeeds because the programmer ensures (using the guard isLeft in line 6) that no value of type U is accessible from this. Note that this cast does not conform to the TYPECASE pattern, despite the guard, because the target type is not a subtype of the cast operand. The cast succeeds only because of Java's type erasure implementation.

```

1 public interface Either<L, R> extends Value<R>, Serializable {
2     @SuppressWarnings("unchecked")
3     default <U> Either<U, R> mapLeft(
4         Function<? super L, ? extends U> leftMapper) {
5         Objects.requireNonNull(leftMapper, "leftMapper is null");
6         if (isLeft()) {
7             return Either.left(leftMapper.apply(getLeft()));
8         } else {
9             return (Either<U, R>) this;
10        }

```

```

11     }
12 }

```

http://bit.ly/vavr_io_vavr_2SMIfI2

Detection. To detect this pattern, application-specific knowledge is required. The developer knows that no value of a type parameter is ever being created. Thus, automatic detection of this pattern seems infeasible.

Issues. This pattern is related to the use of *phantom types* in parametrically polymorphic languages [Leijen and Erik, 1999; Cheney and Hinze, 2003]. Phantom types are type parameters used solely for type checking and are not occupied by any value.

This pattern also occurs with empty collections. For instance, the Java standard library implementation of the method `Collections.<T>emptyList` casts a private constant with raw type `List` to a `List<T>`. This is safe because the list is empty and has no elements of type `T`.

Scala has an unoccupied `Nothing` type to handle this situation. For instance, an (immutable) empty list has `List[Nothing]`, which is a subtype of `List[T]` for any type `T`.

4.6 Discussion

There are common aspects shared by several patterns. Table 4.3 presents a summary of the patterns and their different aspects. The table consists of the following columns:

- **Pattern** Indicates the name of the pattern.
- **Guarded** The patterns here are guarded casts. A guarded cast is a cast such that before the cast is applied, some condition—the *guard*—needs to be verified. The condition to be verified guarantees that the cast will not fail at runtime (unless there is a bug in the application), *i.e.*, the cast will not throw a `ClassCastException`. Some kind of guards ensure that the cast will not fail at the language-level, while others only can guarantee it at the application-level.
- **Language** These casts could be ameliorated if there is enough language support by changing the type system.

- **Tools** The casts in this group could be checked with new analysis or compiler tools.
- **Auto** These casts are related to generated or boilerplate code.
- **Refactor** The casts with this aspect can be simply removed by the developer, can be removed with little refactoring, or suggest a code smell in the source code.
- **Generics** The casts in this category are related to generics or reified generics.
- **Boxing** These casts are related with explicit boxing/unboxing operations, *i.e.*, explicit converting values of primitive types to boxed types and vice versa.
- **QL** A bullet (●) in this column indicates that a pattern is partially detected in QL; a check mark (✓) indicates that we have provided a QL query for automatic detection; and a cross mark (✗) indicates that is infeasible or impractical to detect this pattern in QL.

Many programming languages provide features to ameliorate the more common use cases of casts. For instance, Kotlin's smart casts couple together the `instanceof` operator and cast operation on value, providing direct support for the `TYPECASE` and `EQUALS` patterns. More generally, ML-style pattern matching subsumes this pattern. Smart casts do not apply directly to the `OPERANDSTACK` pattern, since it is dispatched depending on some application-specific control state.

Other language features that might at least partially obviate the need for some of the patterns are intersection types (*cf.* `IMPLICITINTERSECTIONTYPE`), and self types or associated types (*cf.* `FACTORY`, `KNOWNRETURNTYPE`, `DESERIALIZATION`, `COVARIANTRETURNTYPE`, `FLUENTAPI`). Virtual classes [Ernst, 2000; Odersky and Zenger, 2005] and languages that support family polymorphism [Ernst, 2001] would help with casts in the `FAMILY` pattern.

Some cast can be automatically generated. The *StaticResource* variant in `STASH` could be generated by a GUI editor, given that it is most seen in Android applications. The `EQUALS` pattern is composed of boilerplate code. For instance, Scala' solves this issue by introducing *case classes*, which among other features, provide equality out of the box.

Patterns like `FACTORY` are prevalent in test code, because when testing the developer calls the factory methods with known-parameters. The *StaticResource* appear only in source code. This could be because of our sample

Table 4.3. Categorization of Cast Usage Patterns

Pattern	Guarded	Language	Tools	Auto	Refactor	Generics	Boxing	QL
TYPECASE	✓	✓	✓					•
STASH			✓	✓				•
FACTORY			✓					x
FAMILY		✓						x
USERAWTYPE					✓	✓	✓	✓
EQUALS	✓			✓				✓
REDUNDANT					✓			✓
COVARIANTRETURNTYPE		✓						•
SELECTOVERLOAD		✓					✓	✓
KNOWNRETURNTYPE			✓		✓			x
DESERIALIZATION			✓					•
VARIABLESUPERTYPE					✓			•
SOLESUBCLASSIMPLEMENTATION		✓						✓
NEWDYNAMICINSTANCE			✓					✓
OBJECTASARRAY					✓			•
IMPLICITINTERSECTIONTYPE		✓						✓
REMOVESWILDCARD		✓				✓		✓
OPERANDSTACK	✓	✓						x
REFLECTIVEACCESSIBILITY		✓					✓	✓
FLUENTAPI		✓				✓		•
COVARIANTGENERIC		✓				✓	✓	•
COMPOSITE		✓						x
GENERICARRAY		✓				✓	✓	✓
ACCESSSUPERCLASSFIELD					✓			✓
UNOCCUPIEDTYPEPARAMETER		✓				✓		x

does not contain any code generation for Android, *e.g.*, Butter Knife. `USERAWTYPE` is prevalent in generated code. In those cases, code generators do not make the effort to avoid these casts. In our sample, the `ACCESSSUPERCLASSFIELD` pattern only appears in generated code. However, we found other instances using QL.

Our study also suggests analyses could be performed to improve code quality and eliminate some cast usages, for instance finding opportunities to use generics instead (*cf.* `USERAWTYPE`), removing redundant casts (*cf.* `REDUNDANT`), or locating and removing code smells (*cf.* `KNOWNRETURNTYPE`, `VARIABLESUPERTYPE`, and `OBJECTASARRAY`).

The `REMOVESWILDCARD`, `GENERICARRAY` and `COVARIANTGENERIC` patterns are used to workaround the erasure of generic type parameters in Java; while the `UNOCCUPIEDTYPEPARAMETER` pattern is used to take advantage of it. Reified generics or definition-site, rather than use-site, variance annotations [Altidor et al., 2011] would reduce the need for these patterns.

There is an ongoing proposal⁴³ [Smith, 2014] to enhance Java with this feature.

The `USERAWTYPE`, `COVARIANTGENERIC`, and `GENERICARRAY` patterns use either boxing or unboxing because of the interplay between primitive types and generics. The JEP 218 Generics over Primitive Types⁴⁴ [Goetz, 2014] could ameliorate the situation in this respect. On the other hand, the `SELECTOVERLOAD` pattern uses boxing/unboxing to select the appropriate method, while the `REFLECTIVEACCESSIBILITY` pattern uses unboxing when the field being accessed is of a primitive type.

The *QL* column shows whether a pattern can be automatically detected using *QL*. Currently we have 11 patterns for which we can automatically detect them, and 8 where at least we can partially detect them. Just 6 patterns (out of 25 or 24.00%) are impractical to automatically detect.

To detect patterns like `TYPECASE`, `EQUALS`, and `REDUNDANT`, only a local analysis (within a method) is needed. Some generic related patterns, *e.g.*, `USERAWTYPE`, `REMOVESWILDCARD`, and `GENERICARRAY`, are local. On the other hand, patterns like `VARIABLESUPERTYPE` and `UNOCCUPIEDTYPEPARAMETER` require a non-local analysis. However, the `VARIABLESUPERTYPE` pattern can be detected when is instantiated locally, *i.e.*, the cast and the variable assignment are in the same method. Generic related patterns like `FLUENTAPI` and `COVARIANTGENERIC` require a non-local analysis as well.

There are patterns that depend exclusively on known methods, *e.g.*, `NEWDYNAMICINSTANCE`, and `REFLECTIVEACCESSIBILITY`. These patterns are easily detectable. Although a pattern like `DESERIALIZATION` depends on a well-known method—`readObject`—an application could use others deserialization mechanisms.

Some other patterns are inherently complex to detect, *e.g.*, `STASH`, `FAMILY`, `OPERANDSTACK`, and `COMPOSITE`. Recognition of these patterns would require to take into account many different variants, which makes automatic detection impractical. Manually inspection would be better suited in these cases.

Detection of patterns like `FACTORY` and `KNOWNRETURNType` requires to look-up method definitions, often define in external dependencies. At the time of this writing, *QL* does not permit to analyse external dependencies.

⁴³<https://openjdk.java.net/jeps/300>

⁴⁴<https://openjdk.java.net/jeps/218>

4.7 Conclusions

The cast operator in Java bridges the gap between compile-time and run-time safety. We have discovered several cast usage patterns. We found the rationale behind some cast patterns is due to the inexpressiveness of Java's type system. On the other hand, there are patterns that abuse or misuse it.

Many of the patterns we found should be unsurprising to most object-oriented programmers. That nearly 45% of casts are (possibly) unguarded suggests that developers use application-specific knowledge that cannot be easily encoded in the type system to ensure the absence of run-time type errors.

Our study provides insight on the boundary between static and dynamic typing, which may inform research on both static and dynamic, as well as gradual type systems [Siek and Taha, 2006]. Conversely, this research can inform the design of extensions of the Java type system to reduce the need for casting.

Chapter 5

Conclusions

In this thesis I have presented the research I carried out together with my advisors to fulfil the requirements for the Ph.D. degree. We empirically studied how the two mechanisms—*Unsafe* API and casting—are used by developers. We performed qualitative analyses on source code text. In particular, we manually inspected source code text to devise usage patterns.

We have discovered common usage patterns for the Java *Unsafe* API. We discussed several current and future alternatives to improve the Java language. This work has been published in [Mastrangelo et al., 2015]. On the other hand, we complement our *Unsafe* API study with our casting study. We have submitted this study for publication to the OOPSLA’19 conference. We have discovered common usage patterns that involve the cast operator. Having a taxonomy of usage patterns—for both the *Unsafe* API and casting—can shed light on how Java developers give up static type checking.

For our *Unsafe* study, in response to *RQ/U1* (To what extent does the *Unsafe* API impact common application code?), we found that `sun.misc. - Unsafe` is used heavily either directly or indirectly. In response to *RQ/U2* (How and when are *Unsafe* features used?), the results could change if we choose another dataset, *e.g.*, *GitHub*.

For our casting study, in response to *RQ/C2* (How and when casts are used?), we do not claim that our list of patterns is exhaustive. Although our methodology should ensure that any pattern that occurs more than 0.1% of the time has a small probability of being excluded. In response to *RQ/C1* (How frequently is casting used in common application code?) and *RQ/C3* (How and when casts are used?) we assume that casts are uniformly distributed, otherwise our pattern distribution would not reflect

reality.

5.1 Java's Evolution

The Java language is evolving constantly. There are several proposals to improve different aspects of the language. The proposal JEP 193 [Lea, 2014] that introduces Variable Handles is already accepted and included in Java 9. The GC algorithm introduced in JEP 189 Shenandoah [Christine H. Flood, 2014] is included as a experimental feature in Java 12.

There is an ongoing proposal^{1,2} [Goetz, 2017a] to add pattern matching to the Java language. The proposal explores changing the `instanceof` operator in order to support pattern matching. Java 12 extends the `switch` statement to be used as either a statement or an expression^{3,4} [Goetz, 2017b; Bierman, 2019]. This enhancement aims to ease the transition to a `switch` expression that supports pattern matching.

On the other hand, JEP 191 Foreign Function Interface [Nutter, 2014], JEP 169 Value Objects [Rose, 2012a], and JEP 300 Augment Use-Site Variance with Declaration-Site Defaults [Smith, 2014] are still in draft status.

5.2 Limitations and Future Work

Both our *Unsafe* and casting studies rely on manual inspection to devise usage patterns. The main issue with manual inspection is that rely heavily on the personal experience of the authors. Different authors could devise different set of patterns.

In our casting study—whenever possible—we have used QL to automatically detect some patterns. For some other patterns, it is infeasible to perform automatic detection since QL—and the *lgtm* dataset—currently analyse a given project, not its dependencies. Furthermore, some patterns require application-specific knowledge to be detected, which cannot be expressed in QL.

A possible future work could be to run our detection queries on the entire *lgtm* database. This can open up the possibility to devise new us-

¹<https://openjdk.java.net/jeps/305>

²<https://cr.openjdk.java.net/~briangoetz/amber/pattern-match.html>

³<https://openjdk.java.net/jeps/325>

⁴<https://openjdk.java.net/jeps/354>

age patterns, or to refine existing ones. Moreover, by running our queries at large-scale we can corroborate—or refute—the distribution of patterns given in both Sections 4.4 and 4.5.

Conducting ultra-large scale studies, either on source code or compiled code, is not a trivial task. There are several factors to consider when doing these kind of studies, *e.g.*, downloading, storing, parsing, compiling, and analysing software repositories. Services like *Boa* and *lgtm* make conducting this kind of studies easier. In recent versions [Dyer et al., 2015], *Boa* added support to conduct studies on open source projects from *GitHub* and *Qualitas Corpus* [Tempero et al., 2010]. However, at the time we conducted our study on *Unsafe*, this support was not included yet.

We could recast our *Unsafe* study to use *Boa* on the *GitHub* dataset, or *lgtm* through QL queries, although as mentioned above, we will not be able to analyse project dependencies. The patterns we have already devised for the *Unsafe* study could be formalized using QL [Avgustinov et al., 2016].

To conduct our studies, we have used static analysis. Static analyses are always more conservative than dynamic analyses. Another possible future direction could be complement the static analyses with dynamic ones. For the *Unsafe* study, we found that it is used in 1% of the *Maven Central* artifacts. Using project dependencies, 25% of artifacts depend on `sun.misc.Unsafe`. A dynamic analysis could actually measure how often the *Unsafe* API is invoked at run-time, thus giving more precision about its usage. As for the casting study, using a dynamic analysis could measure how many casts fail with `ClassCastException`.

The two studies we conducted in this thesis analyse a single snapshot of a project, *i.e.*, we did not look into the evolution of a project. Some patterns could be better understood in terms of their history. Questions like *How did they solve this problem before using Unsafe?*, or *Why this cast is redundant?* could be answered by analysing the project’s history. For instance, we found that `sun.misc.Unsafe` is heavily used in only 1% of analysed artifacts (48,139 call sites). By looking into the project’s history would be possible to understand why this happened. Source code management tools, *e.g.*, *Git*, maintain a detailed track of changes, which can point out the precise moment in time when an *Unsafe* operation or cast was introduced in a project.

5.3 Lessons Learnt

Throughout my Ph.D. studies I carried out several research projects, most of them described in this thesis. Here are a few lessons I learnt in each of them.

Unsafe API. For our *Unsafe* study, we have engineered a mining software repository infrastructure from scratch. In particular, we learnt about the internals of the *Maven Central* repository. This requires a considerable amount of time to be implemented properly, and for a research project can be of little value.

Cast operator. The cast operator in static typed languages provides a bridge between compile-time and run-time checking. Developers need to resort to the cast operator due to the inexpressiveness of Java's type system.

On the other hand, we have discovered QL, a powerful query language for static analysis. Researchers can use QL in mining software repositories, while software engineers can use it to find vulnerabilities in their code.

Java bytecode instrumentation. Data-flow analysis is complex. The bytecode verification through stack map frames introduced in Java 6 requires implementing a data-flow analysis at bytecode level. Making an industrial-strength implementation is not trivial and requires a lot of careful design.

Supercompilation. Although not mentioned in this thesis, I spent a considerable time working on Supercompilation [Turchin, 1986]. "Supercompilation is a program optimisation technique that is particularly effective at eliminating unnecessary overheads." [Mitchell, 2010].

Although there are several attempts to make Supercompilation mainstream [Mitchell and Runciman, 2007; Bolingbroke and Peyton Jones, 2010], still has some drawbacks that makes it impossible to include it in a compiler of industrial-strength.

5.4 Artifacts

In each of these research projects, I have produced several artifacts. Here are the artifacts mentioned throughout this thesis.

Unsafe API.

- <https://gitlab.com/acuarica/java-unsafe-analysis>: The complete scripts and results used for this study.

- <http://dx.doi.org/10.1145/2814270.2814313>: Companion artifact that aims to reproduce the results found in this study. The companion artifact can be found in the *Source Materials* tab in the ACM Digital Library.

Cast operator.

- <https://gitlab.com/acuarica/phd-thesis/blob/master/analysis/casts.csv>: Comma-separated values (CSV) table. Each row represents a cast instance. This table contains 6 columns. The last two columns—link and value—are the link to the source code file in *lgtm* and the result of the manual inspection.
- <https://gitlab.com/acuarica/phd-thesis/blob/master/analysis/analysis.r>: The script to process the results of the manual inspection.
- <https://gitlab.com/acuarica/java-cast-queries>: QL detection queries.
- <https://gitlab.com/acuarica/java-cast-inspection>: Early prototype of our cast study using javalang and a custom javac plug-in.
- <https://gitlab.com/acuarica/java-cast-study>: Early prototype of our cast study using JavaParser.

Java bytecode instrumentation.

- <https://gitlab.com/acuarica/jnif>: JNIF source code.
- <http://acuarica.gitlab.io/jnif>: JNIF user documentation.

Supercompilation.

- <https://gitlab.com/acuarica/hsc>: Our proof of concept of a Haskell Supercompiler.

Appendix A

Automatic Detection of Patterns using QL

QL [Avgustinov et al., 2016] is “a declarative, object-oriented logic programming language for querying complex, potentially recursive data structures encoded in a relational data model”. QL allows us to analyse programs at the source code level. QL extracts the source code of a project into a Datalog model. Besides providing structural data for programs, *i.e.*, ASTs, QL has the ability to query static types and perform data-flow analysis.

In addition to gather cast usage data using QL, given its powerfulness, we have used QL to approximate the automatic detection of some cast patterns. This appendix gives an introduction to QL. Section A.2 provides the definition of several additional classes and predicates used in Chapter 4.

A.1 Introduction to QL

QL is logic query language, with a syntax that resembles both SQL and Java. It is built up of logical formulas. QL uses logical connectives, *e.g.*, or and not, quantifiers `exists` and `forall`, and logical predicates. QL is highly optimized to support recursive queries. It is possible to use aggregates, *e.g.*, count or sum, in QL as well.

The source code ASTs are modeled as QL classes, *e.g.*, the `CastExpr` class represents the table of cast expressions. Unlike Java, QL “classes are just logical properties describing sets of already existing values.”¹ For instance, Query A.1 gets all cast expressions in a given project. The `from` and `select`

¹<https://help.semmle.com/QL/learn-ql/ql/about-ql.html>

clauses have similar semantics as in SQL. The `import` clause is used to select the language to be analysed. QL supports analysis of several languages, *e.g.*, JavaScript, Python, C/C++, and C#.

```
1 import java
2
3 from CastExpr ce
4 select ce
```

QL

Listing A.1. Query to fetch all cast expressions in a project.

The `where` clause is used to constraint the results. The Query A.2 returns all unused parameters. It selects both the unused parameter and the method where it is declared. In this case, `where` is used similarly to a SQL join clause. The `getAnAccess` class predicate returns any access—read or write—to that parameter.

```
1 import java
2
3 from Parameter p, Method m
4 where not exists(p.getAnAccess())
5     and m.getAParameter() = p
6     and not m.isAbstract()
7 select p, m
```

QL

Listing A.2. Query to fetch unused parameters.

The following query demonstrates how to use the count aggregate. In this example it returns the number of methods with body, *i.e.*, neither abstract nor native.

```
1 import java
2
3 select count(Method m | exists(Block b | m.getBody() = b))
```

QL

Listing A.3. Query to count methods with implementation.

QL permits to define custom classes to refine query results. For example, the following query fetches all primitive cast expressions. A primitive cast is a cast where both the target type and the type of the operand are

primitive types. This excludes any boxed type. The class predicate defined in line 2 is called the *characteristic predicate* or *character* of a class. It is the predicate that determines which values correspond to a given class. Thus, it is similar to filter out results using the where clause.

```

1 class PrimitiveCast extends CastExpr {
2     PrimitiveCast() {
3         getExpr().getType() instanceof PrimitiveType and
4         getTypeExpr().getType() instanceof PrimitiveType
5     }
6 }
7
8 from PrimitiveCast ce
9 select ce

```

QL

The following section describes the additional classes and predicates used throughout Chapter 4.

A.2 Additional QL Classes and Predicates

All cast pattern classes inherit from this base class. As we have seen before, the CastExpr is the QL class that represents all casts. Note that this class does not provide a characteristic predicate. It just adds helper predicates to be used by detection patterns.

```

1 class Cast extends CastExpr {
2     Type getTargetType() { result = getTypeExpr().getType() }
3     Expr getExprOrDef() {
4         result = getExpr() or
5         exists (VariableAssign def |
6             defUsePair(def, getExprOrDef()) and
7             result = def.getSource()
8         )
9     }
10 }

```

QL

Listing A.4. Cast class definition.

This class represents all upcasts. An upcast (T) e happens when the type of e is a subtype of T. Since to detect an upcast is needed to look-up in the class hierarchy, the + operator—transitive closure operator—is used.

```

1  class Upcast extends Cast {
2      Upcast() {
3          getExpr().getType().(RefType).getASupertype+() = getTargetType()
4      }
5  }

```

QL

Listing A.5. Upcast class definition

The following class represents when an argument is used in an overloaded method.

```

1  class ArgumentEx extends Argument {
2      Parameter param;
3      ArgumentEx() {
4          call.getCallee().getParameter(pos) = param
5      }
6  }
7
8  class OverloadedArgument extends ArgumentEx {
9      Callable target;
10     Callable overload;
11     OverloadedArgument() {
12         target = call.getCallee() and
13         overload = target.getDeclaringType().getACallable() and
14         overload.getName() = target.getName() and
15         target != overload
16     }
17     Callable getTarget() { result = target }
18     Callable getAnOverload() { result = overload }
19 }

```

QL

Listing A.6. OverloadedArgument class definition.

This class represents a cast applied to a variable. A QL variable is a field, a local variable or a parameter. The getADef class predicate returns any definition for the variable being cast. This is used in the VARIABLESUPERTYPE pattern.

```

1  class VarCast extends Cast {
2      Variable var;
3      VarCast() { var.getAnAccess() = getExpr() }
4      Variable getVar() { result = var }
5      Expr getADef() {
6          exists (VariableAssign def |
7              defUsePair(def, getExpr()) and
8              result = def.getSource()
9          )
10     }
11 }

```

QL

Listing A.7. VarCast class definition.

The following predicate holds with expressions e , f , and c such that either `if (e == f) c;` or `if (e != f) /*...*/ else c;`. Expressions e and f are interchangeable.

```

1  predicate controlByEqualityTest(Expr e, Expr f, Expr c) {
2      exists (ConditionBlock cb, EqualityTest eqe |
3          eqe.hasOperands(e, f) and eqe = cb.getCondition() and (
4              (eqe.getOp() = " == " and cb.controls(c.getBasicBlock(), true)) or
5              (eqe.getOp() = " != " and cb.controls(c.getBasicBlock(), false))
6          )
7      )
8  }

```

QL

Listing A.8. controlByEqualityTest predicate definition.

Similar to the previous predicate, this predicate holds with expressions e , f , and c such that `if (e.equals(f)) c;`. Expressions e and f are interchangeable.

```

1 predicate controlByEqualsMethod(Expr e, Expr f, Expr c) {
2     exists (ConditionBlock cb, MethodAccess ema |
3         ema.getMethod() instanceof EqualsMethod and (
4             (ema.getQualifier() = e and ema.getArgument(0) = f) or
5             (ema.getQualifier() = f and ema.getArgument(0) = e)
6         ) and (
7             ema = cb.getCondition() and cb.controls(c.getBasicBlock(), true)
8         )
9     )
10 }

```

Listing A.9. controlByEqualsMethod predicate definition.

This predicate holds whenever sub is direct or indirect subclass of sup, or whenever both are the same type.

```

1 predicate isSubtype(RefType sub, RefType sup) {
2     sub.getASupertype*() = sup
3 }

```

Listing A.10. isSubtype predicate definition.

The following class represents all casts guarded with a getClass comparison in an equals method.

```

1 class GetClassGuardsVarCast extends Cast {
2     GetClassMethodAccess tma;
3     GetClassMethodAccess oma;
4     GetClassGuardsVarCast() {
5         tma.isOwnMethodAccess() and
6         oma.getQualifier() = this.(VarCast).getVar().getAnAccess() and
7         (
8             controlByEqualityTest(tma, oma, this) or
9             controlByEqualsMethod(tma, oma, this)
10        )
11    }
12 }

```

Listing A.11. GetClassGuardsVarCast class definition.

These classes define the AutoValue related classes used in the EQUALS pattern.

```

1  class AutoValueAnnotation extends Annotation {
2      AutoValueAnnotation() {
3          getType().hasQualifiedName("com.google.auto.value", "AutoValue")
4      }
5  }
6
7  class AutoValueClass extends Class {
8      AutoValueClass() {
9          getAnAnnotation() instanceof AutoValueAnnotation and
10         isAbstract()
11     }
12 }
13
14 class AutoValueGenerated extends Class {
15     AutoValueGenerated() {
16         count(getASupertype()) = 1 and
17         getASupertype() instanceof AutoValueClass
18     }
19 }

```

Listing A.12. AutoValueGenerated class definition.

This class represents a newInstance method, used in the detection of the NEWDYNAMICINSTANCE pattern.

```

1  class NewDynamicInstanceAccess extends MethodAccess {
2      NewDynamicInstanceAccess() {
3          getCallee().hasName("newInstance") and (
4              getCallee().getDeclaringType() instanceof TypeClass or
5              getCallee().getDeclaringType() instanceof TypeConstructor or
6              getCallee().getDeclaringType() instanceof TypeArray
7          )
8      }
9  }

```

Listing A.13. NewDynamicInstanceAccess class definition.

These classes represent either a Method.invoke or Field.get method access. It is used in the detection of the REFLECTIVEACCESSIBILITY pattern. To implement this query, we made the ReflectiveMethodAccess class *abstract*.

In QL, abstract classes allow the developer “to think of a class as being the union of its subclasses”.² The semantics of an abstract class is as fol-

²<https://help.semmle.com/QL/learn-ql/ql/advanced/abstract-classes.html>

lows: “an abstract class has one or more superclasses and a characteristic predicate. However, for a value to be in an abstract class, it must not only satisfy the character of the class itself, but it must also satisfy the character of a subclass.” Thus, for a value—here, a method access—to be in the `ReflectiveMethodAccess` class, it needs to be either a method access in the `MethodInvokeMethodAccess` or `FieldGetMethodAccess` classes.

```
1  abstract class ReflectiveMethodAccess extends MethodAccess {}
2
3  class MethodInvokeMethodAccess extends ReflectiveMethodAccess {
4      MethodInvokeMethodAccess() {
5          getMethod().hasName("invoke") and
6          getMethod().getDeclaringType()
7              .hasQualifiedName("java.lang.reflect", "Method")
8      }
9  }
10
11 class FieldGetMethodAccess extends ReflectiveMethodAccess {
12     FieldGetMethodAccess() {
13         getMethod().hasName("get") and
14         getMethod().getDeclaringType()
15             .hasQualifiedName("java.lang.reflect", "Field")
16     }
17 }
```

Listing A.14. `ReflectiveMethodAccess` class definition.

The following class represents an invocation to the `setAccessible(true)` method either on a `Field` or `Method` object. It is used in the detection of the `REFLECTIVEACCESSIBILITY` pattern.

```

1  class SetAccessibleTrueMethodAccess extends MethodAccess {
2      Argument flagArgument;
3      SetAccessibleTrueMethodAccess() {
4          getMethod().hasName("setAccessible") and
5          getMethod().getDeclaringType()
6              .hasQualifiedName("java.lang.reflect", "AccessibleObject") and
7          (
8              (getNumArgument() = 1 and flagArgument = getArgument(0)) or
9              (getNumArgument() = 2 and flagArgument = getArgument(1))
10         ) and
11         flagArgument.(BooleanLiteral).getBooleanValue() = true
12     }
13 }

```

Listing A.15. SetAccessibleTrueMethodAccess class definition.

This predicate holds whenever type is neither a raw type (e.g., List), a parameterized type (e.g., List<String>), nor a bounded type (i.e., a type parameter or a wildcard). It is used in the detection of the IMPLICITINTERSECTIONTYPE.

```

1  predicate notGenericRelated(Type type) {
2      not type instanceof RawType and
3      not type instanceof ParameterizedType and
4      not type instanceof BoundedType
5  }

```

Listing A.16. notGenericRelated predicate definition.

Appendix B

JNIF: Java Native Instrumentation

This appendix presents JNIF, our library to instrument Java applications in native code using C/C++. Although the material presented here is not directly related to this thesis, we have used JNIF in several experiments during the development of both Chapters 3 and 4. The original article have been published in Mastrangelo and Hauswirth [2014]. Moreover, the source code of JNIF can be found online.¹

B.1 Introduction

Program analysis tools are important in software engineering tasks such as comprehension, verification and validation, profiling, debugging, and optimization. They can be broadly categorized either as static or dynamic, based on the input that they take. Static analysis tools carry out their task using as input only a program in a given representation, *e.g.*, source code, abstract syntax tree, bytecode, or binary code. In contrast, dynamic analysis tools observe the program being analysed by collecting runtime information. Many dynamic analysis tools rely on instrumentation to achieve their goals.

In the context of the JVM, static analysis and instrumentation for dynamic analysis often happens on the level of Java bytecode. Analysis tools thus need to decode and analyse—and in the case of instrumentation also edit and encode—Java bytecode. Given the relative complexity of the Java class file format, a diverse set of libraries (see Section B.2) has been created for this purpose. All those libraries are implemented in Java.

¹<https://gitlab.com/acuarica/jnif>

Instrumentation at bytecode level can be done in two ways: using a Java instrumentation agent or using a native JVMTI agent.² A Java instrumentation agent is written in Java and runs in the same JVM as the application. This leads to two main problems: poor isolation and poor coverage. It provides *poor isolation* because to instrument the VM, the agent's classes must be loaded in the same VM, and this can lead to perturbation in the VM. It provides *poor coverage* because an instrumentation agent (implemented in Java) will require some runtime library classes to be loaded before it can start instrumenting, and those runtime classes thus cannot be instrumented at load time.

A native JVMTI agent can instrument every class that the VM loads, including runtime classes. The main issue when using JVMTI is that instrumentation must be done in a native language, usually C or C++. Using C/C++ as the instrumentation language can be problematic, because of the lack of a C/C++ library for Java bytecode rewriting. Therefore developers have been using an extra JVM as an “instrumentation server” in which they could use Java-based bytecode rewriting libraries. The C/C++ JVMTI agent thus only has to send code to the server, and no native bytecode rewriting library is needed. However, this approach has a drawback: it requires an additional JVM, and it causes IPC traffic between the observed JVM and the instrumentation server.

We created JNIF to overcome this problem. To the best of our knowledge, JNIF is the first native Java bytecode rewriting library. JNIF is a C++ library for decoding, analysing, editing, and encoding Java bytecode. The main benefit of JNIF is that it can be plugged into a JVMTI agent for instrumenting all classes in a JVM transparently, i.e., without connecting to another JVM and without perturbing the observed JVM.

Starting with Java 6, class files can include stack maps to simplify bytecode verification for the JVM. Java 7 made those stack maps mandatory. Thus, unless one wants to disable the JVM's verifier, code rewriting tools need to also generate stack maps. Stack maps contain, for each basic block, type information for each local variable and operand stack slot. To generate stack maps, a bytecode rewriting tool needs to perform a static analysis. Due to the fact that bytecode does not contain type declarations of variable slots and local variables, these types have to be inferred using an intra-procedural data flow analysis. For reference types, computing the least upper bound of two types in a join point of a control flow graph even requires

²<http://docs.oracle.com/javase/7/docs/technotes/guides/jvmti/index.html>

access to the class hierarchy of the program. Thus, the seemingly innocuous requirement for stack maps significantly complicates the creation of a bytecode rewriting library. JNIF solves these issues, also thanks to the fact that it can be used in-process in a JVMTI agent, and thus can determine the necessary subtyping relationships by requesting the bytes of arbitrary classes loaded or loadable at any given point in time. This works for classes loaded via user-defined class loaders as well as for classes generated dynamically on-the-fly.

Overall, the main contributions of this paper are:

- We present JNIF, a C++ library for decoding, analysing, editing, and encoding Java class files.
- JNIF includes a data-flow analysis for stack map generation, a complication necessary for any library that provides editing and encoding support for modern JVMs with split-time verification.
- We evaluate JNIF by comparing its performance against the most prevalent Java bytecode rewriting library, ASM.

The rest of this chapter is organized as follows: Section B.2 presents related work. In Section B.3 we show how to use the JNIF API. Section B.4 describes the design of JNIF. Section B.5 explains how we validated JNIF. Section B.6 evaluates JNIF's performance against the mainstream bytecode manipulator, ASM. Section B.7 discusses limitations, and Section B.8 concludes.

B.2 Related Work

We now discuss low-level Java bytecode rewriting libraries, JVM hooks for dynamic bytecode rewriting, high-level dynamic bytecode rewriting frameworks, and how they relate to JNIF.

Low-level Rewriting Libraries

JNIF certainly is not the first Java bytecode analysis and instrumentation framework. The probably earliest is BCEL³, a well-designed Java library

³<http://commons.apache.org/bcel/>

with a tree-based API. ASM⁴ [Bruneton et al., 2002b; Kuleshov, 2007] is probably the most prevalent, which aims to be more efficient, especially due to the addition of a visitor-based streaming API, but which has a somewhat less encapsulated design.

SOOT⁵ [Vallée-Rai et al., 1999] is a Java bytecode optimization framework supporting whole-program analysis with four different intermediate representations: Baf, which is simple to manipulate, Jimple, which is easy to optimize, Shimple, an SSA-based variant of Jimple, and Grimp, focused on decompilation.

WALA⁶ is a framework for static analysis, which also includes SHRIKE⁷, a library for instrumenting bytecode using a patch-based approach.

Unlike the above libraries, Javassist⁸ [Chiba and Nishizawa, 2003] provides an API for editing class files like they were Java source code, thereby enabling developers who do not understand bytecode to instrument class files.

Dynamic Instrumentation Hooks

The most limited way for dynamically rewriting Java classes at runtime is the use of a custom class loader. This requires modifications to the application, so that it uses that class loader. This can be problematic for applications, especially for large programs based on powerful frameworks, that already use their own class loaders. This limitation can be circumvented by using dynamic instrumentation hooks provided by the JVM [Lindholm et al.]. Java provides two such hooks: Java agents and JVMTI. Java agents⁹ are supported via the `-javaagent` JVM command line argument. They are implemented in Java and use the `java.lang.instrument` package to interact with the JVM. This allows them to get notified when classes are about to get loaded, and it allows them to modify the class bytecode. They can also modify and reload already loaded classes, however the kinds of transformations allowed with class reloading are severely limited. JVMTI (the *Java Virtual Machine Tool Interface*) is a native API into the JVM that, amongst

⁴<http://asm.ow2.org/>

⁵<http://www.sable.mcgill.ca/soot/>

⁶<http://wala.sourceforge.net/>

⁷http://wala.sourceforge.net/wiki/index.php/Shrike_technical_overview

⁸<http://www.javassist.org/>

⁹<http://docs.oracle.com/javase/7/docs/api/java/lang/instrument/package-summary.html>

many other things, provides hooks that allow the rewriting of bytecode. The advantage of JVMTI over Java agents is that JVMTI allows the instrumentation of *all* Java classes, including the entire runtime library. Java also provides JDI¹⁰ (the *Java Debug Interface*), a high-level interface on top of JVMTI to control a running application in a remote JVM.

High-level Dynamic Analysis Frameworks

We now discuss dynamic analysis frameworks that are built on top of the previously mentioned rewriting libraries and use the above instrumentation hooks. These frameworks do not allow arbitrary code transformations and they shield the developer from the necessary instrumentation effort. Sofya¹¹ [Kinneer et al., 2007] is a dynamic analysis framework that runs the analysis in a separate JVM from the observed application. It provides analysis developers with a set of observable events, to which the analyses can subscribe. Sofya combines bytecode instrumentation using BCEL with the use of JDI for capturing events. FERRARI [Binder et al., 2007] is a dynamic bytecode instrumentation framework that combines static instrumentation of runtime library classes with dynamic instrumentation of application classes to achieve full coverage. FERRARI hooks into the JVM using a Java agent. DiSL [Marek et al., 2012a,b] is a domain-specific aspect language for dynamic analysis. It eliminates the need for static instrumentation from FERRARI by using a separate JVM for instrumentation. It uses JVMTI to hook into the JVM and forwards loaded classes to an instrumentation server, where it performs instrumentation using the ASM rewriting library. Turbo DiSL [Zheng et al., 2012] significantly improves the performance of DiSL by partially evaluating analysis code at instrumentation time. RoadRunner¹² [Flanagan and Freund, 2010] is a high-level framework for creating dynamic analyses focusing on concurrent programs. An analysis implemented on top of RoadRunner simply consists of analysis code in the form of a class that can handle the various event types (such as method calls or field accesses) that RoadRunner can track. RoadRunner uses a custom classloader to be able to rewrite classes at load time, and it uses ASM for bytecode rewriting. Btrace¹³ is an instrumentation tool that allows developers to inject probes based on a predefined set of probe

¹⁰<https://docs.oracle.com/javase/7/docs/jdk/api/jpda/jdi/>

¹¹<http://sofya.unl.edu>

¹²<http://dept.cs.williams.edu/~freund/rr/>

¹³<https://kenai.com/projects/btrace>

types (such as method entry, or bytecode for a specific source line number). Btrace uses the Java agent hooks and builds on top of ASM for instrumentation. Chord¹⁴ [Naik, 2011] is a static analysis framework based on Datalog. It uses joeq¹⁵ to decode classes and convert bytecode into a three-address quadcode internal representation for static analysis. Chord also supports dynamic analysis, for which it instruments programs using Javassist.

How JNIF Differs

Similar to BCEL, JNIF is a low-level library that uses a clean object model to represent java class files. However, unlike all the libraries described above, JNIF is not implemented in Java, but in C++. This allows JNIF to be used directly inside a JVMTI agent. Java-based libraries do not allow dynamic instrumentation in this way: they either are limited to Java agents (which only provide limited coverage), or they require out-of-process instrumentation inside a second JVM (a so-called instrumentation server), and inter-process communication between the JVMTI agent and the instrumentation server.

JNIF simplifies the development and deployment of full-coverage dynamic analysis tools, because one does not need to run an instrumentation server in a separate JVM process. The fact that this is essential is demonstrated by the HPROF¹⁶ profiling agent coming with the JVM. HPROF does not use Java libraries for rewriting bytecode, but implements (a limited form of) class file instrumentation as native code inside a JVMTI agent.

The high-level frameworks described above all abstract away from the underlying instrumentation approach. Thus, they could make use of JNIF to provide their users with full-coverage while eliminating the need for a separate instrumentation server.

B.3 Using JNIF

This section shows common use cases of the JNIF library, such as writing instrumentation code and analysing class files, thus giving an overview of the library. Its components are explained in more detail in Section B.4. JNIF can be used both in stand-alone tools or embedded inside a JVMTI agent.

¹⁴<http://pag.gatech.edu/chord/>

¹⁵<http://joeq.sourceforge.net>

¹⁶<http://docs.oracle.com/javase/7/docs/technotes/samples/hprof.html>

The complete API documentation and more extensive examples are available online¹⁷. We present the examples in an incremental fashion, adding complexity in each example. In order to be able to work with class files, they must be parsed. Given a buffer with a class file and its length, Listing B.1 shows how to parse it.

```
1  const char* data = ...;
2  int len = ...;
3
4  jnif::ClassFile cf(data, len);
```

C++

Listing B.1. Decoding a class

The class `ClassFile` represents a Java class file and contains the definition for each method and fields. Besides providing access to all members of a class, `ClassFile` also provides access to the constant pool via methods like `getUtf8()` and `addMethodRef()`.

Once a class file is correctly parsed and loaded it can be manipulated using the methods and fields in `ClassFile`. For instance, in order to write back the parsed class file in a new buffer, the `write` method is used in conjunction with the `computeSize` method as shown in listing B.2.

```
1  const char* data = ...;
2  int len = ...;
3  jnif::ClassFile cf(data, len);
4  int newlen = cf.computeSize();
5  u1* newdata = new u1[newlen];
6  cf.write(newdata, newlen);
7
8  // Use newdata and newlen
9
10 delete [] newdata;
```

C++

Listing B.2. Encoding a class

The `ClassFile` class has a collection of fields and methods which can be used to discover the members of the class file. The listing B.3 prints in the standard output every method's name and descriptor in a class file.

¹⁷<http://acuarica.gitlab.io/jnif/>

Note that every jnif class overloads the operator<< in order send it to an `std::ostream`.

```
1  const char* data = ...;
2  int len = ...;
3  jnif::ClassFile cf(data, len);
4  for (jnif::Method* m : cf.methods) {
5      cout << "Method: ";
6      cout << cf.getUtf8(m->nameIndex);
7      cout << cf.getUtf8(m->descIndex);
8      cout << endl;
9  }
```

Listing B.3. Traversing all methods in a class

To hook every invocation of a constructor, a method named `<init>` in Java bytecode, one can traverse the method list and check whether the current method is an `<init>` method. Once detected, it is possible to add instrumentation code, like for instance call a static method in a given class. Figure B.4 shows how to add instruction to the instruction list.

```
1  ConstIndex mid = cf.addMethodRef(ci, "alloc", "(Ljava/lang/Object;)V");
2
3  for (Method* method : cf.methods) {
4      if (method->isInit()) {
5          InstList& instList = method->instList();
6
7          Inst* p = *instList.begin();
8          instList.addZero(OPCODE_aload_0, p);
9          instList.addInvoke(OPCODE_invokestatic, mid, p);
10     }
11 }
```

Listing B.4. Instrumenting constructor entries

Another common use case is to instrument every method entry and exit. In order to do so, one can add the instrumentation code at the beginning of the instruction list to detect the method entry. To detect method exit, it is necessary to look for instructions that terminate the current method execution, i.e., `xRETURN` family and `ATHROW` as showed in figure B.5.


```

1  ConstIndex sid = cf.addMethodRef(proxyClass, "enterMethod",
2      "(Ljava/lang/String;Ljava/lang/String;)V");
3  ConstIndex eid = cf.addMethodRef(proxyClass, "exitMethod",
4      "(Ljava/lang/String;Ljava/lang/String;)V");
5  ConstIndex classNameIdx = cf.addStringFromClass(cf.thisClassIndex);
6
7  ...
8
9  InstList& instList = method->instList();
10
11  ConstIndex methodIndex = cf.addString(m->nameIndex);
12
13  Inst* p = *instList.begin();
14
15  instList.addLdc(OPCODE_ldc_w, classNameIdx, p);
16  instList.addLdc(OPCODE_ldc_w, methodIndex, p);
17  instList.addInvoke(OPCODE_invokestatic, sid, p);
18
19  for (Inst* inst : instList) {
20      if (inst->isExit()) {
21          instList.addLdc(OPCODE_ldc_w, classNameIdx, inst);
22          instList.addLdc(OPCODE_ldc_w, methodIndex, inst);
23          instList.addInvoke(OPCODE_invokestatic, eid, inst);
24      }
25  }

```

Listing B.5. Instrumenting <init> methods

B.4 JNIF Design and Implementation

JNIF is written in C++11 [ISO, 2012], in an object-oriented style similar to Java-based class rewriting APIs.

Design

JNIF consists of five main modules: model, parser, writer, printer, and analysis. *Model* implements JNIF's intermediate representation. It is centered around its *ClassFile* class. It is possible to create and manipulate class files from scratch. *Parser* implements the parsing of class files from a given byte array. The parser parses a byte array and translates it to the model's IR. Once a *ClassFile* is created by the parser, it can be manipulated with the methods available in the model. *Writer* and *printer* represent two back-ends

for the model. *Writer* serializes the entire `ClassFile` into a byte array ready to be loaded inside a JVM. *Printer* instead serializes the `ClassFile` into a textual format useful for debugging. Finally, *analysis* implements the static analyses necessary for computing stack maps.

JVM-Independence

JNIF is a stand-alone C++ library that can be used outside a JVM. It does not depend on JVM TI or JNI. However, for the purpose of stack map generation, it may need to determine the common super class of two classes. For this it will need to retrieve a class file given the name of an arbitrary class. This functionality is provided by a plugin that implements JNIF's `IClassPath`. JNIF comes with such a plugin that uses JNI in case it is running inside a JVM.

Explicit Constant Pool Management

Unlike some other class rewriting libraries, JNIF exposes the constant pool instead of hiding it. Our reasons for this design decision were two-fold: (1) We wanted to fully control the structure of the class file, and for that it is necessary to expose the constant pool. To reduce the additional complexity, we provide a rich set of methods that facilitate constant pool management. (2) We wanted to preserve, whenever possible, the original structure of the class file. This means that if one parses and then writes a class file, the original bytes will be obtained. This decreases the perturbation done by the instrumentation and allows for better testing.

Memory Management

Given that JNIF is implemented in an unmanaged language, we have to worry about memory deallocation. Our API follows a simple ownership model where all IR objects are owned by their enclosing objects. This means, that the `ClassFile` object owns the complete IR of a class. Our API design enforces this ownership model by requiring IR objects to be created by their enclosing objects. For example, to create a `Method`, one has to use the `ClassFile::addMethod()` factory method instead of directly allocating a new `Method` object.

Stack Map Generation

When encoding a `ClassFile` into a byte array, JNIF needs to generate stack maps. The necessary static analyses are implemented in the analyser module. This module uses data flow analysis and abstract interpretation to determine the types of operand stack slots and local variables. The analysis module first builds a control flow graph of the method. The data flow analysis associates to each basic block an input and output stack frame, which represents the types of the local variables and operand stack slots at that point in the code. The input frame represents the type before any instruction in the basic block is executed. The output frame is computed by abstract interpretation of each instruction in the basic block. The entry basic block has an empty stack and each entry in its local variable table is set to top. Then the algorithm starts from the entry block and follows each edge. If a basic block is reachable from multiple edges, then a merge is involved.

Merging involves finding the least upper bound of multiple incoming types. While this is trivial for primitive types, it can require access to the class hierarchy for reference types. This requirement represents a severe complication for binary rewriting tools: when rewriting a single class, they may require access to many other classes in the program. JNIF solves this problem by providing the `IClassPath` interface. Different `IClassPath` implementations can provide different ways for getting access to classes. For example, a static instrumentation tool may use a user-defined class path to find classes, while a dynamic instrumentation tool may use JNI to request the bytes of a class given that class' name.

Running JNIF Inside a JVM TI Agent

When using JNIF inside a JVM TI agent, JNIF uses an `IClassPath` implementation that uses JNI to load the bytes of classes required for least upper bound computations during stack map generation.

Avoiding Premature Static Initialization

Using JNI to load a class (with `ClassLoader.loadClass()`), however, will call that class' static initializer. This is a side effect that may change the observable behavior of the program under analysis. To avoid this, one can request the bytes of the class (with `ClassLoader.getResourceAsStream()`)

instead of loading the class. It can then parse the bytes of the class into its IR to determine that class' supertypes.

Avoiding the Loading of the Class Being Instrumented

If during the instrumentation of a class X JNIF needs to perform a least upper bound computation involving type X , then using `ClassLoader.loadClass` to load class X would cause an infinite recursion. The above solution with `getResourceAsStream()` also prevents this problem.

Avoiding Premature `ClassNotFoundException`

If during the instrumentation of a class X JNIF needs to perform a least upper bound computation involving a type Y , and if class Y cannot be found, then throwing a `ClassNotFoundException` at that time would be premature (because without instrumentation, such an exception would only be thrown later). We solve that problem by assuming a least upper bound of `java.lang.Object` in that case.

B.5 Validation

We used a multitude of testing strategies to ensure JNIF is working correctly.

The JNIF parser and writer makes no extra modification to the class file, thus it is an exact representation of the class file. This property makes the parser and writer returns the identical same byte stream which can be useful for testing purposes.

Unit Tests

JNIF includes a unit test suite that tests individual features of its various modules.

Integration Tests

Our integration test suite includes six different JNIF clients we run on over 40000 different classes. Each test reads, analyses, and possibly modifies, prints, or writes classes from the Java runtime library (`rt.jar`), and all Dacapo benchmarks, Scala benchmarks, and the JRuby compiler.

testPrinter. This test parses and prints all classes. Its main goal is to cover the printing functionality. It has no explicit assertions. We consider it passed if it does not throw any exceptions.

testSize. This test covers the decoding and encoding modules. It asserts that the encoded byte array has the same length as the original byte array.

testWriter. This is similar to testSize, but it asserts that the contents of the encoded byte array is identical to the original bytes.

testNopAdderInstrPrinter. This also tests the instrumentation functionality, by injecting NOP instructions and dumping the result. It passes if it does not throw any exceptions.

testNopAdderInstrSize. This is similar to testSize, however it performs NOP injection. The resulting size must be identical to the original size plus the size of the injected NOP instructions.

testNopAdderInstrWriter. This is similar to testNopAdderInstrSize, but it asserts that the resulting array is identical except for the modified method bytecodes.

The “size” and “writer” tests work thanks to the fact that JNIF produces output identical to its input as long as classes are not modified and stack maps do not need to be re-generated.

Live Tests

Our live tests use JNIF inside a JVMTI dynamic instrumentation agent to ensure that the output of JNIF can successfully be loaded, verified, and run by a JVM. In addition to the aspects covered by the unit and integration tests, the live tests also validate that stack map generation works correctly, essentially by using the JVM’s verifier to check correctness. For the live tests, we run a set of microbenchmarks, the Dacapo benchmarks, the Scala Benchmarking Project¹⁸, and a microbenchmark using the JRuby compiler, with the goal of including InvokeDynamic bytecode instructions generated by JRuby.

¹⁸<http://www.benchmarks.scalabench.org/>

Assertions and Checks

The JNIF code is sprinkled with calls to `Error::assert` that check preconditions, postconditions, and invariants. To provide a developer experience similar to Java's, all assertion violations print out call stack traces in addition to understandable error messages.

Moreover, JNIF checks its inputs (such as class files while parsing, or instrumented code while generating stack maps), and it calls `Error::check` to throw exceptions with stack traces and helpful messages when checks fail.

B.6 Performance Evaluation

We evaluated the performance of a JNIF-based dynamic instrumentation approach versus an approach using an ASM-based instrumentation server.

Measurement Contexts

We ran our experiments on three different machines: (1) A machine with two Intel Xeon E5-2620 2 GHz CPUs, each with 6 cores and 2 threads per core, and 8 GB RAM, running Debian Linux x86 64 3.10.11-1. (2) A Dell PowerEdge M620, 2 NUMA node with 64 GB of RAM, Intel Xeon E5-2680 2.7 GHz CPU with 8 cores, CPU frequency scaling and Turbo Mode disabled, running Ubuntu Linux x86 64 3.8.0-38. For consistent memory access speed, we bound our program to a specific NUMA node using `numactl`. (3) A MacBook Pro with an Intel Core i7 2.7 GHz CPU with 4 cores and 16 GB running Mac OS X 10.8.2.

Benchmarks

We used the Dacapo benchmarks, except for `tradebeans` and `tradesoap`, which suffer from a well known issue.¹⁹ We also include the Scala benchmarks (except for the subset identical to Dacapo).

¹⁹<http://sourceforge.net/p/dacapobench/bugs/70/>

Subjects

We compare JNIF to ASM for the purpose of performing dynamic instrumentation. For JNIF we built a JVMTI agent that directly includes JNIF to instrument loaded classes. For ASM, we use a JVMTI agent that forwards loaded classes to an instrumentation server that uses ASM's streaming API (which is faster than ASM's tree API).

Results

Figure B.1 shows the results of our performance evaluation in terms of time spent instrumenting classes. The figure shows the results from our first machine. The other machines produced results similar to Figure B.1. The figure shows box plots summarizing five measurement runs. It shows one box for JNIF and two boxes for ASM. The "ASM Server" box represents the time as measured on the instrumentation server. This is equivalent to the time a static instrumentation tool would take. It excludes the time spent in the JVMTI agent and the time for the IPC between the agent and the server. The "ASM Server on Client" box represents the total time needed for instrumentation, as measured in the JVMTI agent, and thus includes the IPC and JVMTI agent time.

Each chart in the figure consists of five groups of boxes: "Empty" is the time when using a JVMTI agent that does not process bytecodes at all. "Identity" is for an agent that simply decodes and encodes each class, without any instrumentation, and without recomputing stack maps. "Compute-Frames" also includes recomputing stack maps. "Allocations" represents a useful dynamic analysis that captures all allocations. "Nop Padding" is a different dynamic analysis that injects NOPs after each bytecode instruction.

The figure shows that frame computation adds significant overhead, on ASM as well as JNIF. Moreover, it shows that except for dacapo-eclipse, dacapo-jython, and scala-scalatest, JNIF is faster even than just the ASM Server time.

Reproducibility

To run these evaluations, a Makefile script is provided in the git repository. These tasks take care of the compilation of the JNIF library and also all java

files needed. The repository is self-contained, no need to download dacapo benchmarks separately.

```
> make testapp
```

Listing B.6. Running testapp

```
> make testapp
```

Listing B.7. Running dacapo

To run a particular dacapo benchmark with default settings

```
> make dacapo BENCH=avroora
```

Listing B.8. Running dacapo

To run a full evaluation with all dacapo benchmarks in all configuration a task `-eval-` is provided. You can set how many times run each configuration with the variable `times`, like

```
> make eval times=5
```

Listing B.9. Running full eval five times

Finally, there is a task to create plots for the evaluation. This task needs R with the package `ggplot2`.

```
> make plots
```

Listing B.10. Plots

B.7 Limitations

JNIF still has some limitations.

jsr/ret. JNIF does not support stack map generation for jsr and ret. Class files requiring stack maps do not include jsr/ret. The jsr/ret instructions make the control flow graph generation difficult, because a ret instruction can jump to multiple targets instead of a predefined one.

invokedynamic. JNIF's support for invokedynamic is not yet fully tested but our initial tests with JRuby have been successful. Dacapo bach was released in 2009, before the creation of Java 7, which introduces the invoke-dynamic instruction. Thus it does not contain any benchmark with invoke-dynamic. Instead we use JRuby 1.7 in order to create a self-contained jar file. This jar file does not contain any invokedynamic instruction, but it does contain the JRuby compiler, that when specified via `-Djruby.compile.invoke-dynamic=true` it will generate class files with invokedynamic. We tested our parser and writer with this settings with successful results.

Stack map generation with full coverage. When the JVM loads the first few runtime library classes, and calls the JVMTI agent to have those classes instrumented, it is still too early to use JNI for loading classes needed for computing least upper bounds for stack map generation. For this reason, we do not generate stack maps for runtime library classes. This no problem, because the JVM does not verify the runtime library classes by default, and thus it does not need stack maps for those classes. However, should developers decide to explicitly turn on the verification of runtime library classes (with `-Xverify:all`), the verifier would complain because JNIF would not have generated stack maps.

To get full coverage for the instrumentation inside a JVMTI agent, it is necessary to instrument every class, even the whole java class library. If the instrumentation needs to change or add branch targets, the compute frames option must be used, but it cannot be used against the class library, because to compute frames, the class hierarchy must be known, and this imposes a dependency with a classloader which is not yet available.

Luckily, by default the Java library classes are not verified, because they are trusted. Thus the instrumentation only needs to compute frames on classes not belonging to the Java library.

B.8 Conclusions

Until now, full-coverage dynamic instrumentation in production JVMs required performing the code rewriting in a separate JVM, because of the lack of a native bytecode rewriting framework. This paper introduces JNIF, the

first full-coverage in-process dynamic instrumentation framework for Java. It discusses the key issues of creating such a framework for Java—such as stack-map generation—and it evaluates the performance of JNIF against the most prevalent Java-level framework: ASM. We find that JNIF is faster than using out-of-process ASM in most cases. We hope that thanks to JNIF, and this paper, a broader number of researchers and developers will be enabled to develop native JVM agents that analyse and rewrite Java bytecode without limitations.

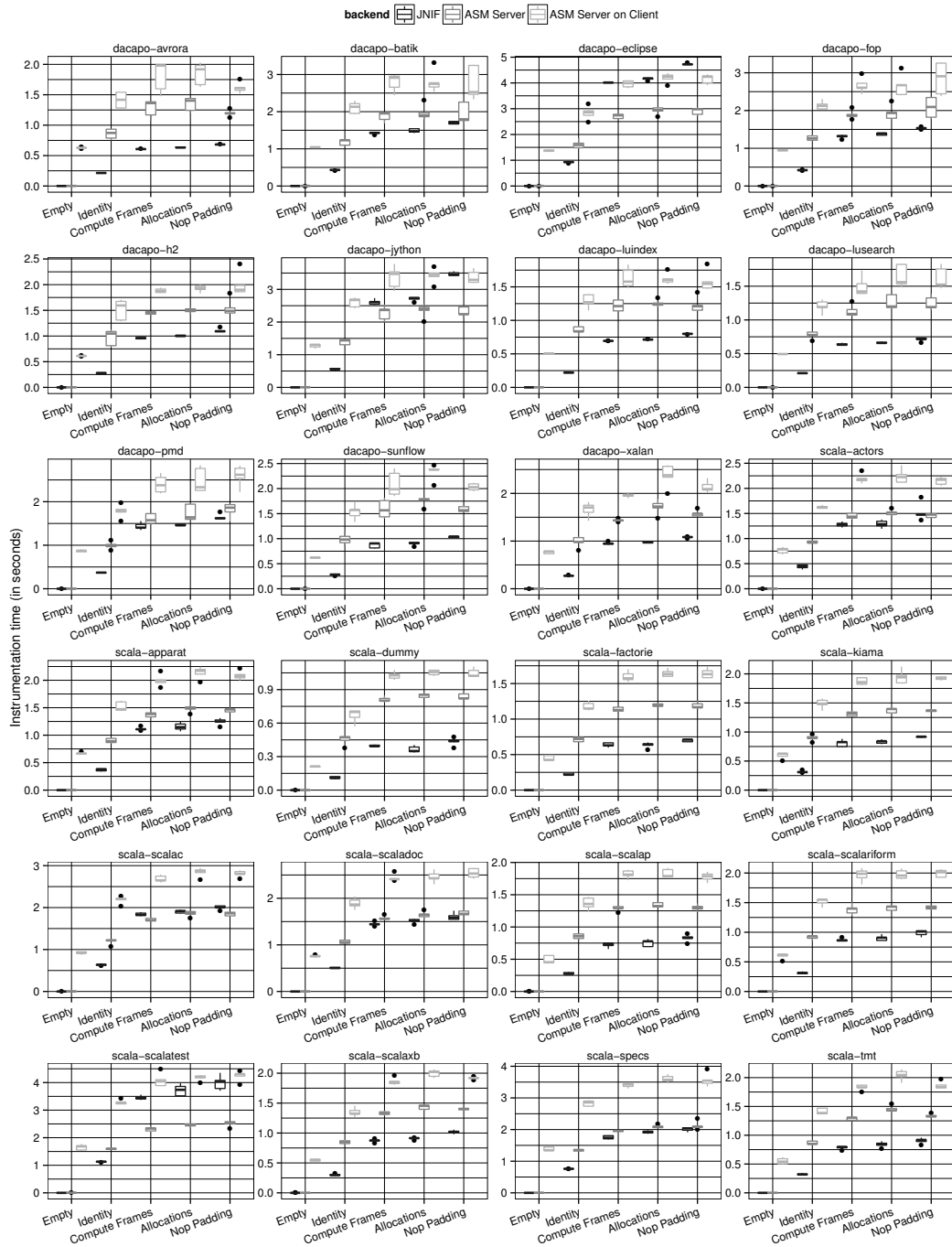


Figure B.1. Instrumentation time on DaCapo and Scala benchmarks

Bibliography

Miltiadis Allamanis and Charles Sutton. Mining source code repositories at massive scale using language modeling. In *2013 10th Working Conference on Mining Software Repositories (MSR)*, pages 207–216, San Francisco, CA, USA, May 2013. IEEE. ISBN 978-1-4673-2936-1 978-1-4799-0345-0. doi: 10.1109/MSR.2013.6624029.

John Altidor, Shan Shan Huang, and Yannis Smaragdakis. Taming the Wildcards: Combining Definition- and Use-site Variance. In *Proceedings of the 32Nd ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '11*, pages 602–613, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0663-8. doi: 10.1145/1993498.1993569.

Nada Amin and Ross Tate. Java and Scala’s Type Systems Are Unsound: The Existential Crisis of Null Pointers. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2016*, pages 838–848, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4444-9. doi: 10.1145/2983990.2984004.

Denis N. Antonioli and Markus Pilz. Analysis of the Java Class File Format. Technical report, University of Zurich, 1998.

Francesca Arcelli, Fabrizio Perin, Claudia Raibulet, and Stefano Ravani. Design Pattern Detection in Java Systems: A Dynamic Analysis Based Approach. In *Evaluation of Novel Approaches to Software Engineering, Communications in Computer and Information Science*, pages 163–179. Springer, Berlin, Heidelberg, May 2008. ISBN 978-3-642-14818-7 978-3-642-14819-4. doi: 10.1007/978-3-642-14819-4_12.

Muhammad Asaduzzaman, Muhammad Ahasanuzzaman, Chanchal K. Roy, and Kevin A. Schneider. How Developers Use Exception Handling

- in Java? In *Proceedings of the 13th International Conference on Mining Software Repositories*, MSR '16, pages 516–519, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4186-8. doi: 10.1145/2901739.2903500.
- Pavel Avgustinov, Oege de Moor, Michael Peyton Jones, and Max Schäfer. QL: Object-oriented Queries on Relational Data. In Shriram Krishnamurthi and Benjamin S. Lerner, editors, *30th European Conference on Object-Oriented Programming (ECOOP 2016)*, volume 56 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 2:1–2:25, Dagstuhl, Germany, 2016. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. ISBN 978-3-95977-014-9. doi: 10.4230/LIPIcs.ECOOP.2016.2.
- David F. Bacon, Perry Cheng, and V. T. Rajan. A Real-time Garbage Collector with Low Overhead and Consistent Utilization. In *Proceedings of the 30th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '03, pages 285–298, New York, NY, USA, 2003. ACM. ISBN 978-1-58113-628-9. doi: 10.1145/604131.604155.
- S. Bajracharya, J. Ossher, and Cristina Lopes. Sourcerer: An internet-scale software repository. In *Tools and Evaluation 2009 ICSE Workshop on Search-Driven Development-Users, Infrastructure*, pages 1–4, May 2009. doi: 10.1109/SUITE.2009.5070010.
- I.D. Baxter, A. Yahin, L. Moura, M. Sant'Anna, and L. Bier. Clone detection using abstract syntax trees. In *Proceedings. International Conference on Software Maintenance (Cat. No. 98CB36272)*, pages 368–377, Bethesda, MD, USA, 1998. IEEE Comput. Soc. ISBN 978-0-8186-8779-2. doi: 10.1109/ICSM.1998.738528.
- Gavin Bierman. JEP 354: Switch Expressions. 2019. URL <https://openjdk.java.net/jeps/354>.
- W. Binder, J. Hulaas, and P. Moret. Reengineering Standard Java Runtime Systems through Dynamic Bytecode Instrumentation. In *Seventh IEEE International Working Conference on Source Code Analysis and Manipulation (SCAM 2007)*, pages 91–100, September 2007. doi: 10.1109/SCAM.2007.20.
- Stephen M. Blackburn, Robin Garner, Chris Hoffmann, Asjad M. Khang, Kathryn S. McKinley, Rotem Bentzur, Amer Diwan, Daniel Feinberg, Daniel Frampton, Samuel Z. Guyer, Martin Hirzel, Antony Hosking,

- Maria Jump, Han Lee, J. Eliot B. Moss, Aashish Phansalkar, Darko Stefanović, Thomas VanDrunen, Daniel von Dincklage, and Ben Wiederemann. The DaCapo Benchmarks: Java Benchmarking Development and Analysis. In *Proceedings of the 21st Annual ACM SIGPLAN Conference on Object-Oriented Programming Systems, Languages, and Applications, OOPSLA '06*, pages 169–190, New York, NY, USA, 2006. ACM. ISBN 978-1-59593-348-5. doi: 10.1145/1167473.1167488.
- Maximilian Bolingbroke and Simon Peyton Jones. Supercompilation by Evaluation. In *Proceedings of the Third ACM Haskell Symposium on Haskell, Haskell '10*, pages 135–146, New York, NY, USA, 2010. ACM. ISBN 978-1-4503-0252-4. doi: 10.1145/1863523.1863540.
- S. Brandauer and T. Wrigstad. Spencer: Interactive Heap Analysis for the Masses. In *2017 IEEE/ACM 14th International Conference on Mining Software Repositories (MSR)*, pages 113–123, May 2017. doi: 10.1109/MSR.2017.35.
- Martin Bravenboer and Yannis Smaragdakis. Exception analysis and points-to analysis: Better together. page 11, a.
- Martin Bravenboer and Yannis Smaragdakis. Strictly declarative specification of sophisticated points-to analyses. page 19, b.
- Kim B. Bruce. Some Challenging Typing Issues in Object-Oriented Languages. *Electronic Notes in Theoretical Computer Science*, 82(8):1–29, October 2003. ISSN 15710661. doi: 10.1016/S1571-0661(04)80799-0.
- E. Bruneton, T. Coupaye, and J. B. Stefani. Recursive and Dynamic Software Composition with Sharing. 2002a.
- Eric Bruneton, Romain Lenglet, and Thierry Coupaye. ASM: A code manipulation tool to implement adaptable systems. In *In Adaptable and Extensible Component Systems*, 2002b.
- Marian Bubak and Dawid Kurzyniec. Creating java to native code interfaces with Janet extension. In *In Proceedings of the First Worldwide SGI Users's Conference*, pages 283–294, 2000.
- Timothy A. Budd, Richard A. DeMillo, Richard J. Lipton, and Frederick G. Sayward. Theoretical and Empirical Studies on Using Program Mutation to Test the Functional Correctness of Programs. In *Proceedings of the*

- 7th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '80, pages 220–233, New York, NY, USA, 1980. ACM. ISBN 978-0-89791-011-8. doi: 10.1145/567446.567468.
- Oscar Callaú, Romain Robbes, Éric Tanter, and David Röthlisberger. How (and why) developers use the dynamic features of programming languages: The case of smalltalk. *Empirical Software Engineering*, 18(6):1156–1194, December 2013. ISSN 1382-3256, 1573-7616. doi: 10.1007/s10664-012-9203-2.
- Luca Cardelli and Peter Wegner. On understanding types, data abstraction, and polymorphism. *ACM Computing Surveys*, 17(4):471–523, December 1985. ISSN 03600300. doi: 10.1145/6041.6042.
- Manuel M. T. Chakravarty, Gabriele Keller, and Simon Peyton Jones. Associated Type Synonyms. In *Proceedings of the Tenth ACM SIGPLAN International Conference on Functional Programming*, ICFP '05, pages 241–253, New York, NY, USA, 2005. ACM. ISBN 978-1-59593-064-4. doi: 10.1145/1086365.1086397.
- Xiliang Chen, Alice Yuchen Wang, and Ewan Tempero. A Replication and Reproduction of Code Clone Detection Studies. page 10.
- James Cheney and Ralf Hinze. First-class phantom types. Technical Report TR2003-1901, Cornell University, July 2003.
- R. J. Cheavance and T. Heidet. Static Profile and Dynamic Behavior of COBOL Programs. *SIGPLAN Not.*, 13(4):44–57, April 1978. ISSN 0362-1340. doi: 10.1145/953411.953414.
- Shigeru Chiba and Muga Nishizawa. An Easy-to-Use Toolkit for Efficient Java Bytecode Translators. In Frank Pfenning and Yannis Smaragdakis, editors, *Generative Programming and Component Engineering*, Lecture Notes in Computer Science, pages 364–376. Springer Berlin Heidelberg, 2003. ISBN 978-3-540-39815-8.
- Roman Kennke Christine H. Flood. JEP 189: Shenandoah: An Ultra-Low-Pause-Time Garbage Collector. 2014. URL <https://openjdk.java.net/jeps/189>.
- Roberta Coelho, Lucas Almeida, Georgios Gousios, and Arie van Deursen. Unveiling Exception Handling Bug Hazards in Android Based on GitHub

- and Google Code Issues. In *Proceedings of the 12th Working Conference on Mining Software Repositories*, MSR '15, pages 134–145, Piscataway, NJ, USA, 2015. IEEE Press. ISBN 978-0-7695-5594-2.
- Tal Cohen and Itay Maman. JTL – the Java Tools Language. page 20.
- Christian Collberg, Ginger Myles, and Michael Stepp. An empirical study of Java bytecode programs. *Software: Practice and Experience*, 37(6):581–641, May 2007. ISSN 1097-024X. doi: 10.1002/spe.776.
- R. P. Cook. An empirical analysis of the Lilith instruction set. *IEEE Transactions on Computers*, 38(1):156–158, January 1989. ISSN 0018-9340. doi: 10.1109/12.8740.
- Robert P. Cook and Insup Lee. A contextual analysis of Pascal programs. *Software: Practice and Experience*, 12(2):195–203, February 1982. ISSN 1097-024X. doi: 10.1002/spe.4380120209.
- Diego Costa, Artur Andrzejak, Janos Seboek, and David Lo. Empirical Study of Usage and Performance of Java Collections. In *Proceedings of the 8th ACM/SPEC on International Conference on Performance Engineering*, ICPE '17, pages 389–400, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-4404-3. doi: 10.1145/3030207.3030221.
- Johannes Dahse and Thorsten Holz. Experience Report: An Empirical Study of PHP Security Mechanism Usage. In *Proceedings of the 2015 International Symposium on Software Testing and Analysis*, ISSTA 2015, pages 60–70, New York, NY, USA, 2015. ACM. ISBN 978-1-4503-3620-8. doi: 10.1145/2771783.2771787.
- Oege de Moor, M. Verbaere, E. Hajiyeve, P. Avgustinov, T. Ekman, N. Ongkingco, D. Sereni, and J. Tibble. Keynote Address: .QL for Source Code Analysis. In *Seventh IEEE International Working Conference on Source Code Analysis and Manipulation (SCAM 2007)*, pages 3–16, September 2007. doi: 10.1109/SCAM.2007.31.
- Coen De Roover, Carlos Noguera, Andy Kellens, and Vivane Jonckers. The SOUL Tool Suite for Querying Programs in Symbiosis with Eclipse. In *Proceedings of the 9th International Conference on Principles and Practice of Programming in Java*, PPPJ '11, pages 71–80, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0935-6. doi: 10.1145/2093157.2093168.

- Sylvia Dieckmann and Urs Hölzle. A study of the Allocation Behavior of the SPECjvm98 Java Benchmarks. In Rachid Guerraoui, editor, *ECOOP'99 — Object-Oriented Programming*, Lecture Notes in Computer Science, pages 92–115. Springer Berlin Heidelberg, 1999. ISBN 978-3-540-48743-2.
- J. Dietrich, K. Jezek, and P. Brada. Broken promises: An empirical study into evolution problems in Java programs caused by library upgrades. In *2014 Software Evolution Week - IEEE Conference on Software Maintenance, Reengineering, and Reverse Engineering (CSMR-WCRE)*, pages 64–73, February 2014. doi: 10.1109/CSMR-WCRE.2014.6747226.
- Jens Dietrich, David J. Pearce, Kamil Jezek, and Premek Brada. Contracts in the Wild: A Study of Java Programs. In Peter Müller, editor, *31st European Conference on Object-Oriented Programming (ECOOP 2017)*, volume 74 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 9:1–9:29, Dagstuhl, Germany, 2017a. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. ISBN 978-3-95977-035-4. doi: 10.4230/LIPIcs.ECOOP.2017.9.
- Jens Dietrich, Henrik Schole, Li Sui, and Ewan Tempero. XCorpus – An executable Corpus of Java Programs. *The Journal of Object Technology*, 16(4):1:1, 2017b. ISSN 1660-1769. doi: 10.5381/jot.2017.16.4.a1.
- M. Doyle and J. Walden. An Empirical Study of the Evolution of PHP Web Application Security. In *2011 Third International Workshop on Security Measurements and Metrics*, pages 11–20, September 2011. doi: 10.1109/Metrise.2011.18.
- R. Dyer, H. A. Nguyen, H. Rajan, and T. N. Nguyen. Boa: A language and infrastructure for analyzing ultra-large-scale software repositories. In *2013 35th International Conference on Software Engineering (ICSE)*, pages 422–431, May 2013a. doi: 10.1109/ICSE.2013.6606588.
- Robert Dyer, Hoan Anh Nguyen, Hridesh Rajan, and Tien N. Nguyen. Mining Source Code Repositories with Boa. In *Proceedings of the 2013 Companion Publication for Conference on Systems, Programming, & Applications: Software for Humanity, SPLASH '13*, pages 13–14, New York, NY, USA, 2013b. ACM. ISBN 978-1-4503-1995-9. doi: 10.1145/2508075.2514570.
- Robert Dyer, Hridesh Rajan, and Tien N. Nguyen. Declarative Visitors to Ease Fine-grained Source Code Mining with Full History on Billions of AST Nodes. In *Proceedings of the 12th International Conference on Generative*

- Programming: Concepts & Experiences*, GPCE '13, pages 23–32, New York, NY, USA, 2013c. ACM. ISBN 978-1-4503-2373-4. doi: 10.1145 / 2517208.2517226.
- Robert Dyer, Hridayesh Rajan, Hoan Anh Nguyen, and Tien N. Nguyen. Mining Billions of AST Nodes to Study Actual and Potential Usage of Java Language Features. In *Proceedings of the 36th International Conference on Software Engineering*, ICSE 2014, pages 779–790, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2756-5. doi: 10.1145/2568225.2568295.
- Robert Dyer, Hoan Nguyen, H Rajan, and T Nguyen. Boa: An enabling language and infrastructure for ultra-large-scale msr studies. pages 593–621, 01 2015. doi: 10.1016/B978-0-12-411519-4.00020-3.
- Torbjörn Ekman and Görel Hedin. The jastadd extensible java compiler. In *Proceedings of the 22nd annual ACM SIGPLAN conference on Object-oriented programming systems and applications*, OOPSLA '07, pages 1–18, New York, NY, USA, 2007. ACM. ISBN 978-1-59593-786-5. doi: 10.1145 / 1297027.1297029. URL <http://doi.acm.org/10.1145/1297027.1297029>.
- Erik Ernst. gbeta - a language with virtual attributes, block structure, and propagating, dynamic inheritance. *DAIMI Report Series*, 29(549), May 2000. doi: 10.7146 / dpb.v29i549.7654. URL <https://tidsskrift.dk/daimipb/article/view/7654>.
- Erik Ernst. Family Polymorphism. In *ECOOP 2001 — Object-Oriented Programming*, Lecture Notes in Computer Science, pages 303–326. Springer, Berlin, Heidelberg, June 2001. ISBN 978-3-540-42206-8 978-3-540-45337-6. doi: 10.1007/3-540-45337-7_17.
- Cormac Flanagan and Stephen N. Freund. The RoadRunner Dynamic Analysis Framework for Concurrent Programs. In *Proceedings of the 9th ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering*, PASTE '10, pages 1–8, New York, NY, USA, 2010. ACM. ISBN 978-1-4503-0082-7. doi: 10.1145/1806672.1806674.
- Justin E Forrester and Barton P Miller. An Empirical Study of the Robustness of Windows NT Applications Using Random Testing. page 10, 2000.
- George Fourtounis, George Kastrinis, and Yannis Smaragdakis. Static Analysis of Java Dynamic Proxies. In *Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis*, ISSTA 2018, pages

- 209–220, New York, NY, USA, 2018. ACM. ISBN 978-1-4503-5699-2. doi: 10.1145/3213846.3213864.
- Milos Gligoric, Darko Marinov, and Sam Kamin. CoDeSe: Fast Deserialization via Code Generation. In *Proceedings of the 2011 International Symposium on Software Testing and Analysis, ISSTA '11*, pages 298–308, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0562-4. doi: 10.1145/2001420.2001456.
- Brian Goetz. JEP 218: Generics over Primitive Types. 2014. URL <https://openjdk.java.net/jeps/218>.
- Brian Goetz. JEP 305: Pattern Matching for instanceof (Preview). 2017a. URL <https://openjdk.java.net/jeps/305>.
- Brian Goetz. JEP 325: Switch Expressions (Preview). 2017b. URL <https://openjdk.java.net/jeps/325>.
- Alessandra Gorla, Ilaria Tavecchia, Florian Gross, and Andreas Zeller. Checking App Behavior Against App Descriptions. In *Proceedings of the 36th International Conference on Software Engineering, ICSE 2014*, pages 1025–1035, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2756-5. doi: 10.1145/2568225.2568276.
- James Gosling, Bill Joy, Guy L. Steele, Jr., Gilad Bracha, and Alex Buckley. *The Java Language Specification, Java SE 7 Edition*. Addison-Wesley Professional, 2013. ISBN 0-13-326022-4 978-0-13-326022-9.
- Georgios Gousios. The GHTorrent Dataset and Tool Suite. In *Proceedings of the 10th Working Conference on Mining Software Repositories, MSR '13*, pages 233–236, Piscataway, NJ, USA, 2013. IEEE Press. ISBN 978-1-4673-2936-1.
- Georgios Gousios, Bogdan Vasilescu, Alexander Serebrenik, and Andy Zaidman. Lean GHTorrent: GitHub Data on Demand. In *Proceedings of the 11th Working Conference on Mining Software Repositories, MSR 2014*, pages 384–387, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2863-0. doi: 10.1145/2597073.2597126.
- Mark Grechanik, Collin McMillan, Luca DeFerrari, Marco Comi, Stefano Crespi, Denys Poshyvanyk, Chen Fu, Qing Xie, and Carlo Ghezzi. An Empirical Investigation into a Large-scale Java Open Source Code Repository. In *Proceedings of the 2010 ACM-IEEE International Symposium on*

- Empirical Software Engineering and Measurement*, ESEM '10, pages 11:1–11:10, New York, NY, USA, 2010. ACM. ISBN 978-1-4503-0039-1. doi: 10.1145/1852786.1852801.
- John Hammond. BASIC - an evaluation of processing methods and a study of some programs. *Software: Practice and Experience*, 7(6):697–711, November 1977. ISSN 1097-024X. doi: 10.1002/spe.4380070605.
- I. R. Harlin, H. Washizaki, and Y. Fukazawa. Impact of Using a Static-Type System in Computer Programming. In *2017 IEEE 18th International Symposium on High Assurance Systems Engineering (HASE)*, pages 116–119, January 2017. doi: 10.1109/HASE.2017.17.
- Mark Hills, Paul Klint, and Jurgen Vinju. An Empirical Study of PHP Feature Usage: A Static Analysis Perspective. In *Proceedings of the 2013 International Symposium on Software Testing and Analysis, ISSTA 2013*, pages 325–335, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-2159-4. doi: 10.1145/2483760.2483786.
- Lei Hu and Kamran Sartipi. Dynamic Analysis and Design Pattern Detection in Java Programs. In *20th International Conference on Software Engineering and Knowledge Engineering, SEKE 2008*, pages 842–846, January 2008.
- Atsushi Igarashi, Benjamin C. Pierce, and Philip Wadler. Featherweight Java: A Minimal Core Calculus for Java and GJ. *ACM Trans. Program. Lang. Syst.*, 23(3):396–450, May 2001. ISSN 0164-0925. doi: 10.1145/503502.503505.
- ISO. *ISO/IEC 14882:2011 Information Technology — Programming Languages — C++*. International Organization for Standardization, Geneva, Switzerland, February 2012.
- Dongseok Jang, Ranjit Jhala, Sorin Lerner, and Hovav Shacham. An Empirical Study of Privacy-violating Information Flows in JavaScript Web Applications. In *Proceedings of the 17th ACM Conference on Computer and Communications Security, CCS '10*, pages 270–283, New York, NY, USA, 2010. ACM. ISBN 978-1-4503-0245-6. doi: 10.1145/1866307.1866339.
- Maria Kechagia and Diomidis Spinellis. Undocumented and Unchecked: Exceptions That Spell Trouble. In *Proceedings of the 11th Working Conference*

- on Mining Software Repositories*, MSR 2014, pages 312–315, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2863-0. doi: 10.1145/2597073.2597089.
- Mary Beth Kery, Claire Le Goues, and Brad A. Myers. Examining Programmer Practices for Locally Handling Exceptions. In *Proceedings of the 13th International Conference on Mining Software Repositories*, MSR '16, pages 484–487, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4186-8. doi: 10.1145/2901739.2903497.
- A. Kinneer, M. B. Dwyer, and G. Rothermel. Sofya: Supporting Rapid Development of Dynamic Program Analyses for Java. In *29th International Conference on Software Engineering - Companion*, 2007. ICSE 2007 Companion, pages 51–52, May 2007. doi: 10.1109/ICSECOMPANION.2007.68.
- Oleg Kiselyov, Ralf Lämmel, and Kean Schupke. Strongly typed heterogeneous collections. In *Proceedings of the ACM SIGPLAN Workshop on Haskell - Haskell '04*, page 96, Snowbird, Utah, USA, 2004. ACM Press. ISBN 978-1-58113-850-4. doi: 10.1145/1017472.1017488.
- Oleg Kiselyov, Simon Peyton Jones, and Chung-chieh Shan. Fun with type functions, April 2009. URL <https://www.microsoft.com/en-us/research/publication/fun-type-functions/>.
- P. Klint, T. v d Storm, and J. Vinju. RASCAL: A Domain Specific Language for Source Code Analysis and Manipulation. In *2009 Ninth IEEE International Working Conference on Source Code Analysis and Manipulation*, pages 168–177, September 2009. doi: 10.1109/SCAM.2009.28.
- Donald E. Knuth. An empirical study of FORTRAN programs. *Software: Practice and Experience*, 1(2):105–133, April 1971. ISSN 1097-024X. doi: 10.1002/spe.4380010203.
- Goh Kondoh and Tamiya Onodera. Finding Bugs in Java Native Interface Programs. In *Proceedings of the 2008 International Symposium on Software Testing and Analysis*, ISSTA '08, pages 109–118, New York, NY, USA, 2008. ACM. ISBN 978-1-60558-050-0. doi: 10.1145/1390630.1390645.
- Guy Korland, Nir Shavit, and Pascal Felber. Noninvasive concurrency with Java STM. January 2010.
- Eugene Kuleshov. *Using the ASM Framework to Implement Common Java Byte-code Transformation Patterns*. 2007.

- D. Landman, A. Serebrenik, and J. J. Vinju. Challenges for Static Analysis of Java Reflection - Literature Review and Empirical Study. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*, pages 507–518, May 2017. doi: 10.1109/ICSE.2017.53.
- Thomas D. LaToza and Brad A. Myers. Developers Ask Reachability Questions. In *Proceedings of the 32Nd ACM/IEEE International Conference on Software Engineering - Volume 1, ICSE '10*, pages 185–194, New York, NY, USA, 2010. ACM. ISBN 978-1-60558-719-6. doi: 10.1145/1806799.1806829.
- Doug Lea. JEP 193: Variable Handles. 2014. URL <https://openjdk.java.net/jeps/193>.
- Daan Leijen and Meijer Erik. Domain-specific embedded compilers. In *In Proceedings of the 2nd Conference on Domain-Specific Languages*, pages 109–122. USENIX Association, October 1999.
- Siliang Li and Gang Tan. Finding Bugs in Exceptional Situations of JNI Programs. In *Proceedings of the 16th ACM Conference on Computer and Communications Security, CCS '09*, pages 442–452, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-894-0. doi: 10.1145/1653662.1653716.
- Tim Lindholm, Frank Yellin, Gilad Bracha, and Alex Buckley. The Java® Virtual Machine Specification. page 626.
- Benjamin Livshits. *Improving Software Security with Precise Static and Runtime Analysis*. PhD thesis, Stanford University, Stanford, California, 2006.
- Benjamin Livshits, John Whaley, and Monica S. Lam. Reflection Analysis for Java. In *Programming Languages and Systems, Lecture Notes in Computer Science*, pages 139–160. Springer, Berlin, Heidelberg, November 2005. ISBN 978-3-540-29735-2 978-3-540-32247-4. doi: 10.1007/11575467_11.
- Cristina V. Lopes, Petr Maj, Pedro Martins, Vaibhav Saini, Di Yang, Jakub Zitny, Hitesh Sajnani, and Jan Vitek. DéjàVu: A Map of Code Duplicates on GitHub. *Proc. ACM Program. Lang.*, 1(OOPSLA):84:1–84:28, October 2017. ISSN 2475-1421. doi: 10.1145/3133908.
- Magnus Madsen and Esben Andreasen. String Analysis for Dynamic Field Access. In David Hutchison, Takeo Kanade, Josef Kittler, Jon M. Kleinberg, Friedemann Mattern, John C. Mitchell, Moni Naor, Oscar Nierstrasz, C. Pandu Rangan, Bernhard Steffen, Madhu Sudan, Demetri Ter-

- zopoulos, Doug Tygar, Moshe Y. Vardi, Gerhard Weikum, and Albert Cohen, editors, *Compiler Construction*, volume 8409, pages 197–217. Springer Berlin Heidelberg, Berlin, Heidelberg, 2014. ISBN 978-3-642-54806-2 978-3-642-54807-9. doi: 10.1007/978-3-642-54807-9_12.
- Lukáš Marek, Alex Villazón, Yudi Zheng, Danilo Ansaloni, Walter Binder, and Zhengwei Qi. DiSL: A domain-specific language for bytecode instrumentation. In *Proceedings of the 11th Annual International Conference on Aspect-Oriented Software Development - AOSD '12*, page 239, Potsdam, Germany, 2012a. ACM Press. ISBN 978-1-4503-1092-5. doi: 10.1145/2162049.2162077.
- Lukáš Marek, Yudi Zheng, Danilo Ansaloni, Walter Binder, Zhengwei Qi, and Petr Tuma. DiSL: An Extensible Language for Efficient and Comprehensive Dynamic Program Analysis. In *Proceedings of the Seventh Workshop on Domain-Specific Aspect Languages, DSAL '12*, pages 27–28, New York, NY, USA, 2012b. ACM. ISBN 978-1-4503-1128-1. doi: 10.1145/2162037.2162046.
- Luis Mastrangelo and Matthias Hauswirth. JNIF: Java Native Instrumentation Framework. In *Proceedings of the 2014 International Conference on Principles and Practices of Programming on the Java Platform: Virtual Machines, Languages, and Tools, PPPJ '14*, pages 194–199, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2926-2. doi: 10.1145/2647508.2647516.
- Luis Mastrangelo, Luca Ponzanelli, Andrea Mocci, Michele Lanza, Matthias Hauswirth, and Nathaniel Nystrom. Use at Your Own Risk: The Java Unsafe API in the Wild. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2015*, pages 695–710, New York, NY, USA, 2015. ACM. ISBN 978-1-4503-3689-5. doi: 10.1145/2814270.2814313.
- Clemens Mayer, Stefan Hanenberg, Romain Robbes, Éric Tanter, and Andreas Stefik. An Empirical Study of the Influence of Static Type Systems on the Usability of Undocumented Software. In *Proceedings of the ACM International Conference on Object Oriented Programming Systems Languages and Applications, OOPSLA '12*, pages 683–702, New York, NY, USA, 2012. ACM. ISBN 978-1-4503-1561-6. doi: 10.1145/2384616.2384666.
- Davood Mazinanian, Ameya Ketkar, Nikolaos Tsantalis, and Danny Dig. Understanding the Use of Lambda Expressions in Java. *Proc. ACM Pro-*

- gram. Lang.*, 1(OOPSLA):85:1–85:31, October 2017. ISSN 2475-1421. doi: 10.1145/3133909.
- Barton P. Miller, Louis Fredriksen, and Bryan So. An empirical study of the reliability of UNIX utilities. *Communications of the ACM*, 33(12):32–44, December 1990. ISSN 00010782. doi: 10.1145/96267.96279.
- Barton P. Miller, David Koski, Cjin Pheow, Lee Vivekananda Maganty, Ravi Murthy, Ajitkumar Natarajan, and Jeff Steidl. Fuzz revisited: A re-examination of the reliability of UNIX utilities and services. Technical report, 1995.
- Robin Milner. A Proposal for Standard ML. In *Proceedings of the 1984 ACM Symposium on LISP and Functional Programming*, LFP '84, pages 184–197, New York, NY, USA, 1984. ACM. ISBN 978-0-89791-142-9. doi: 10.1145/800055.802035.
- Neil Mitchell. Rethinking Supercompilation. In *Proceedings of the 15th ACM SIGPLAN International Conference on Functional Programming*, ICFP '10, pages 309–320, New York, NY, USA, 2010. ACM. ISBN 978-1-60558-794-3. doi: 10.1145/1863543.1863588.
- Neil Mitchell and Colin Runciman. A Supercompiler for Core Haskell. In *Implementation and Application of Functional Languages*, Lecture Notes in Computer Science, pages 147–164. Springer, Berlin, Heidelberg, September 2007. ISBN 978-3-540-85372-5 978-3-540-85373-2. doi: 10 . 1007 / 978-3-540-85373-2_9.
- Andrew C. Myers. Software composition with multiple nested inheritance. 2006.
- Meiyappan Nagappan, Thomas Zimmermann, and Christian Bird. Diversity in Software Engineering Research. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, ESEC/FSE 2013, pages 466–476, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-2237-9. doi: 10.1145/2491411.2491415.
- Meiyappan Nagappan, Romain Robbes, Yasutaka Kamei, Éric Tanter, Shane McIntosh, Audris Mockus, and Ahmed E. Hassan. An Empirical Study of Goto in C Code from GitHub Repositories. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*, ESEC/FSE 2015,

- pages 404–414, New York, NY, USA, 2015. ACM. ISBN 978-1-4503-3675-8. doi: 10.1145/2786805.2786834.
- Mayur Naik. Chord: A Versatile Platform for Program Analysis. 2011.
- Suman Nakshatri, Maithri Hegde, and Sahithi Thandra. Analysis of Exception Handling Patterns in Java Projects: An Empirical Study. In *Proceedings of the 13th International Conference on Mining Software Repositories, MSR '16*, pages 500–503, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4186-8. doi: 10.1145/2901739.2903499.
- Charles Oliver Nutter. JEP 191: Foreign Function Interface. 2014. URL <https://openjdk.java.net/jeps/191>.
- Nathaniel Nystrom, Michael R. Clarkson, and Andrew C. Myers. Polyglot: An Extensible Compiler Framework for Java. In *Compiler Construction, Lecture Notes in Computer Science*, pages 138–152. Springer, Berlin, Heidelberg, April 2003. ISBN 978-3-540-00904-7 978-3-540-36579-2. doi: 10.1007/3-540-36579-6_11.
- Martin Odersky and Matthias Zenger. Independently extensible solutions to the expression problem. In *Proc. FOOL 12*, January 2005. URL <http://homepages.inf.ed.ac.uk/wadler/fool>.
- Diarmuid O'Donoghue, Aine Leddy, James Power, and John Waldron. Bigram Analysis of Java Bytecode Sequences. In *Proceedings of the Inaugural Conference on the Principles and Practice of Programming, 2002 and Proceedings of the Second Workshop on Intermediate Representation Engineering for Virtual Machines, 2002*, PPPJ '02/IRE '02, pages 187–192, Maynooth, County Kildare, Ireland, Ireland, 2002. National University of Ireland. ISBN 978-0-901519-87-0.
- Bruno C. d. S. Oliveira and William R. Cook. Extensibility for the Masses. In David Hutchison, Takeo Kanade, Josef Kittler, Jon M. Kleinberg, Friedemann Mattern, John C. Mitchell, Moni Naor, Oscar Nierstrasz, C. Pandu Rangan, Bernhard Steffen, Madhu Sudan, Demetri Terzopoulos, Doug Tygar, Moshe Y. Vardi, Gerhard Weikum, and James Noble, editors, *ECOOP 2012 – Object-Oriented Programming*, volume 7313, pages 2–27. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012. ISBN 978-3-642-31056-0 978-3-642-31057-7. doi: 10.1007/978-3-642-31057-7_2.

- Bruno C. d. S. Oliveira, Zhiyuan Shi, and João Alpuim. Disjoint intersection types. In *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming*, ICFP 2016, pages 364–377, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4219-3. doi: 10.1145/2951913.2951945. URL <http://doi.acm.org/10.1145/2951913.2951945>.
- Nathalie Oostvogels, Joeri De Koster, and Wolfgang De Meuter. Static Typing of Complex Presence Constraints in Interfaces. In Todd Millstein, editor, *32nd European Conference on Object-Oriented Programming (ECOOP 2018)*, volume 109 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 14:1–14:27, Dagstuhl, Germany, 2018. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. ISBN 978-3-95977-079-8. doi: 10.4230/LIPIcs.ECOOP.2018.14.
- OpenJDK. Project Sumatra. 2013. URL <https://openjdk.java.net/projects/sumatra/>.
- F. Palomba, D. Di Nucci, M. Tufano, G. Bavota, R. Oliveto, D. Poshyvanyk, and A. De Lucia. Landfill: An Open Dataset of Code Smells with Public Evaluation. In *2015 IEEE/ACM 12th Working Conference on Mining Software Repositories*, pages 482–485, May 2015. doi: 10.1109/MSR.2015.69.
- Fabio Palomba, Gabriele Bavota, Massimiliano Di Penta, Rocco Oliveto, Andrea De Lucia, and Denys Poshyvanyk. Detecting bad smells in source code using change history information. In *2013 28th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 268–278, Silicon Valley, CA, USA, November 2013. IEEE. ISBN 978-1-4799-0215-6. doi: 10.1109/ASE.2013.6693086.
- Chris Parnin, Christian Bird, and Emerson Murphy-Hill. Java Generics Adoption: How New Features Are Introduced, Championed, or Ignored. In *Proceedings of the 8th Working Conference on Mining Software Repositories*, MSR ’11, pages 3–12, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0574-7. doi: 10.1145/1985441.1985446.
- Chris Parnin, Christian Bird, and Emerson Murphy-Hill. Adoption and use of Java generics. *Empirical Software Engineering*, 18(6):1047–1089, December 2013. ISSN 1382-3256, 1573-7616. doi: 10.1007/s10664-012-9236-6.
- Renaud Pawlak, Martin Monperrus, Nicolas Petitprez, Carlos Noguera, and Lionel Seinturier. Spoon: A Library for Implementing Analyses

- and Transformations of Java Source Code. *Software: Practice and Experience*, 46:1155–1179, 2015. doi: 10 . 1002 / spe . 2346. URL <https://hal.archives-ouvertes.fr/hal-01078532/document>.
- Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Geoffrey Washburn. Simple unification-based type inference for gadts. In *Proceedings of the ACM SIGPLAN International Conference on Functional Programming*, ICFP’06, pages 50–61. ACM Press, April 2006. URL <https://www.microsoft.com/en-us/research/publication/simple-unification-based-type-inference-for-gadts/>.
- Benjamin C. Pierce. *Types and Programming Languages*. The MIT Press, 1st edition, 2002. ISBN 978-0-262-16209-8.
- Gustavo Pinto, Wesley Torres, Benito Fernandes, Fernando Castor, and Roberto S.M. Barros. A large-scale study on the usage of Java’s concurrent programming constructs. *Journal of Systems and Software*, 106:59–81, 2015. ISSN 0164-1212. doi: <https://doi.org/10.1016/j.jss.2015.04.064>.
- D. Posnett, C. Bird, and P. Devanbu. THEX: Mining metapatterns from java. In *2010 7th IEEE Working Conference on Mining Software Repositories (MSR 2010)*, pages 122–125, May 2010. doi: 10.1109/MSR.2010.5463349.
- L. Prechelt. An empirical comparison of seven programming languages. *Computer*, 33(10):23–29, October 2000. ISSN 0018-9162. doi: 10.1109/2.876288.
- M Pukall, C Kaestner, W Cazzola, S Goetz, A Grebhahn, and R Schroeter. Flexible Dynamic Software Updates of Java Applications: Tool Support and Case Study. page 39.
- Xin Qi and Andrew C. Myers. Masked Types for Sound Object Initialization. In *Proceedings of the 36th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL ’09, pages 53–65, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-379-2. doi: 10.1145/1480881.1480890.
- Baishakhi Ray, Daryl Posnett, Premkumar Devanbu, and Vladimir Filkov. A Large-scale Study of Programming Languages and Code Quality in GitHub. *Commun. ACM*, 60(10):91–100, September 2017. ISSN 0001-0782. doi: 10.1145/3126905.

- M. Rebouças, G. Pinto, F. Ebert, W. Torres, A. Serebrenik, and F. Castor. An Empirical Study on the Usage of the Swift Programming Language. In *2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER)*, volume 1, pages 634–638, March 2016. doi: 10.1109/SANER.2016.66.
- Gregor Richards, Sylvain Lebesne, Brian Burg, and Jan Vitek. An Analysis of the Dynamic Behavior of JavaScript Programs. In *Proceedings of the 31st ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '10*, pages 1–12, New York, NY, USA, 2010. ACM. ISBN 978-1-4503-0019-3. doi: 10.1145/1806596.1806598.
- Gregor Richards, Christian Hammer, Brian Burg, and Jan Vitek. The Eval That Men Do: A Large-scale Study of the Use of Eval in Javascript Applications. In *Proceedings of the 25th European Conference on Object-Oriented Programming, ECOOP'11*, pages 52–78, Berlin, Heidelberg, 2011. Springer-Verlag. ISBN 978-3-642-22654-0.
- Matthias Rieger and Stephane Ducasse. Visual Detection of Duplicated Code. page 6.
- John Rose. JEP 169: Value Objects. 2012a. URL <https://openjdk.java.net/jeps/169>.
- John Rose, Brian Goetz, and Guy Steele. State of the Values. 2014. URL <https://cr.openjdk.java.net/~jrose/values/values-0.html>.
- John R. Rose. Arrays 2.0. 2012b.
- John R. Rose. The isthmus in the VM. 2014. URL <https://cr.openjdk.java.net/~jrose/panama/isthmus-in-the-vm-2014.html>.
- Harry J. Saal and Zvi Weiss. Some Properties of APL Programs. In *Proceedings of Seventh International Conference on APL, APL '75*, pages 292–297, New York, NY, USA, 1975. ACM. doi: 10.1145/800117.803819.
- Harry J. Saal and Zvi Weiss. An empirical study of APL programs. *Computer Languages*, 2(3):47–59, January 1977. ISSN 0096-0551. doi: 10.1016/0096-0551(77)90007-8.
- A Salvadori, J. Gordon, and C. Capstick. Static Profile of COBOL Programs. *SIGPLAN Not.*, 10(8):20–33, August 1975. ISSN 0362-1340. doi: 10.1145/956028.956031.

- Paul Sandoz. Safety Not Guaranteed: Sun.misc.Unsafe and the quest for safe alternatives. 2014. Oracle Inc. [Online; accessed 29-January-2015].
- Paul Sandoz. Personal communication. 2015.
- Z. Shen, Z. Li, and P. C. Yew. An empirical study of Fortran programs for parallelizing compilers. *IEEE Transactions on Parallel and Distributed Systems*, 1(3):356–364, July 1990. ISSN 1045-9219. doi: 10.1109/71.80162.
- Fridtjof Siebert. Eliminating External Fragmentation in a Non-moving Garbage Collector for Java. In *Proceedings of the 2000 International Conference on Compilers, Architecture, and Synthesis for Embedded Systems, CASES '00*, pages 9–17, New York, NY, USA, 2000. ACM. ISBN 978-1-58113-338-7. doi: 10.1145/354880.354883.
- Jeremy G. Siek and Walid Taha. Gradual typing for functional languages. In *Scheme and functional programming workshop*, pages 81–92, 2006.
- Dan Smith. JEP 300: Augment Use-Site Variance with Declaration-Site Defaults. 2014. URL <https://openjdk.java.net/jeps/300>.
- E. Soloway and K. Ehrlich. Empirical Studies of Programming Knowledge. *IEEE Transactions on Software Engineering*, SE-10(5):595–609, September 1984. ISSN 0098-5589. doi: 10.1109/TSE.1984.5010283.
- Christopher Strachey. Fundamental Concepts in Programming Languages. *Higher-Order and Symbolic Computation*, 13(1-2):11–49, April 2000. ISSN 1388-3690, 1573-0557. doi: 10.1023/A:1010000313106.
- Andreas Stuchlik and Stefan Hanenberg. Static vs. Dynamic Type Systems: An Empirical Study About the Relationship Between Type Casts and Development Time. In *Proceedings of the 7th Symposium on Dynamic Languages, DLS '11*, pages 97–106, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0939-4. doi: 10.1145/2047849.2047861.
- Mengtao Sun and Gang Tan. NativeGuard: Protecting Android Applications from Third-party Native Libraries. In *Proceedings of the 2014 ACM Conference on Security and Privacy in Wireless & Mobile Networks, WiSec '14*, pages 165–176, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2972-9. doi: 10.1145/2627393.2627396.

- Gang Tan and Jason Croft. An Empirical Security Study of the Native Code in the JDK. /paper/An-Empirical-Security-Study-of-the-Native-Code-in-Tan-Croft/4c3a84729bd09db6a90a862846bb29e937ec2ced, 2008.
- Gang Tan, Srimat Chakradhar, Raghunathan Srivaths, and Ravi Daniel Wang. Safe Java native interface. In *In Proceedings of the 2006 IEEE International Symposium on Secure Software Engineering*, pages 97–106, 2006.
- E. Tempero, C. Anslow, J. Dietrich, T. Han, J. Li, M. Lumpe, H. Melton, and J. Noble. The Qualitas Corpus: A Curated Collection of Java Code for Empirical Studies. In *2010 Asia Pacific Software Engineering Conference*, pages 336–345, November 2010. doi: 10.1109/APSEC.2010.46.
- Ewan Tempero, James Noble, and Hayden Melton. How Do Java Programs Use Inheritance? An Empirical Study of Inheritance in Java Software. In Jan Vitek, editor, *ECOOP 2008 – Object-Oriented Programming*, Lecture Notes in Computer Science, pages 667–691. Springer Berlin Heidelberg, 2008. ISBN 978-3-540-70592-5.
- N. M. Tiwari, G. Upadhyaya, H. A. Nguyen, and H. Rajan. Candoia: A Platform for Building and Sharing Mining Software Repositories Tools as Apps. In *2017 IEEE/ACM 14th International Conference on Mining Software Repositories (MSR)*, pages 53–63, May 2017. doi: 10.1109/MSR.2017.56.
- N. Tsantalis, T. Chaikalis, and A. Chatzigeorgiou. JDeodorant: Identification and Removal of Type-Checking Bad Smells. In *2008 12th European Conference on Software Maintenance and Reengineering*, pages 329–331, April 2008. doi: 10.1109/CSMR.2008.4493342.
- Michele Tufano, Fabio Palomba, Gabriele Bavota, Rocco Oliveto, Massimiliano Di Penta, Andrea De Lucia, and Denys Poshyvanyk. When and Why Your Code Starts to Smell Bad. In *Proceedings of the 37th International Conference on Software Engineering - Volume 1, ICSE '15*, pages 403–414, Piscataway, NJ, USA, 2015. IEEE Press. ISBN 978-1-4799-1934-5.
- Michele Tufano, Fabio Palomba, Gabriele Bavota, Rocco Oliveto, Massimiliano Di Penta, Andrea De Lucia, and Denys Poshyvanyk. When and Why Your Code Starts to Smell Bad (and Whether the Smells Go Away). *IEEE Transactions on Software Engineering*, 43(11):1063–1088, November 2017. ISSN 0098-5589, 1939-3520. doi: 10.1109/TSE.2017.2653105.

- Valentin F. Turchin. The concept of a supercompiler. *ACM Transactions on Programming Languages and Systems*, 8(3):292–325, June 1986. ISSN 01640925. doi: 10.1145/5956.5957.
- Phillip Merlin Uesbeck, Andreas Stefik, Stefan Hanenberg, Jan Pedersen, and Patrick Daleiden. An Empirical Study on the Impact of C++ Lambdas and Programmer Experience. In *Proceedings of the 38th International Conference on Software Engineering, ICSE '16*, pages 760–771, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-3900-1. doi: 10.1145/2884781.2884849.
- Raoul-Gabriel Urma and Alan Mycroft. Programming language evolution via source code query languages. In *Proceedings of the ACM 4th Annual Workshop on Evaluation and Usability of Programming Languages and Tools - PLATEAU '12*, page 35, Tucson, Arizona, USA, 2012. ACM Press. ISBN 978-1-4503-1631-6. doi: 10.1145/2414721.2414728.
- Raja Vallée-Rai, Phong Co, Etienne Gagnon, Laurie Hendren, Patrick Lam, and Vijay Sundaresan. Soot - a Java Bytecode Optimization Framework. In *Proceedings of the 1999 Conference of the Centre for Advanced Studies on Collaborative Research, CASCON '99*, pages 13–, Mississauga, Ontario, Canada, 1999. IBM Press.
- Mandana Vaziri, Frank Tip, Stephen Fink, and Julian Dolby. Declarative Object Identity Using Relation Types. In *ECOOP 2007 – Object-Oriented Programming*, Lecture Notes in Computer Science, pages 54–78. Springer, Berlin, Heidelberg, July 2007. ISBN 978-3-540-73588-5 978-3-540-73589-2. doi: 10.1007/978-3-540-73589-2_4.
- Jurgen Vinju and James R. Cordy. How to make a bridge between transformation and analysis technologies? In James R. Cordy, Ralf Lämmel, and Andreas Winter, editors, *Transformation Techniques in Software Engineering*, Dagstuhl Seminar Proceedings, Dagstuhl, Germany, 2006. Internationales Begegnungs- und Forschungszentrum für Informatik (IBFI), Schloss Dagstuhl, Germany.
- Kris De Volder. JQuery: A generic code browser with a declarative configuration language. In *Practical Aspects of Declarative Languages, 8th International Symposium, PADL 2006*, pages 88–102. Springer, 2006.
- Philip Wadler. The expression problem, November 1998. URL <https://homepages.inf.ed.ac.uk/wadler/papers/expression/expression.txt>.

- Philip Wadler and Robert Bruce Findler. Well-Typed Programs Can't Be Blamed. In Giuseppe Castagna, editor, *Programming Languages and Systems*, Lecture Notes in Computer Science, pages 1–16. Springer Berlin Heidelberg, 2009. ISBN 978-3-642-00590-9.
- Yanlin Wang and Bruno C. d. S. Oliveira. The expression problem, trivially! In *Proceedings of the 15th International Conference on Modularity*, MODULARITY 2016, pages 37–41, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-3995-7. doi: 10.1145/2889443.2889448. URL <http://doi.acm.org/10.1145/2889443.2889448>.
- Shiyi Wei, Francesca Xhakaj, and Barbara G. Ryder. Empirical study of the dynamic behavior of JavaScript objects. *Software: Practice and Experience*, 46(7):867–889, July 2016. ISSN 1097-024X. doi: 10.1002/spe.2334.
- Johnni Winther. Guarded Type Promotion: Eliminating Redundant Casts in Java. In *Proceedings of the 13th Workshop on Formal Techniques for Java-Like Programs*, FTfJP '11, pages 6:1–6:8, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0893-9. doi: 10.1145/2076674.2076680.
- Baijun Wu and Sheng Chen. How Type Errors Were Fixed and What Students Did? *Proc. ACM Program. Lang.*, 1(OOPSLA):105:1–105:27, October 2017. ISSN 2475-1421. doi: 10.1145/3133929.
- Baijun Wu, John Peter Campora III, and Sheng Chen. Learning User Friendly Type-error Messages. *Proc. ACM Program. Lang.*, 1(OOPSLA):106:1–106:29, October 2017. ISSN 2475-1421. doi: 10.1145/3133930.
- Yang Yuan and Yao Guo. CMCD: Count matrix based code clone detection. In *Proceedings - Asia-Pacific Software Engineering Conference, APSEC*, pages 250–257, December 2011. doi: 10.1109/APSEC.2011.13.
- Yudi Zheng, Danilo Ansaloni, Lukas Marek, Andreas Sewe, Walter Binder, Alex Villazón, Petr Tuma, Zhengwei Qi, and Mira Mezini. Turbo DiSL: Partial Evaluation for High-Level Bytecode Instrumentation. In David Hutchison, Takeo Kanade, Josef Kittler, Jon M. Kleinberg, Friedemann Mattern, John C. Mitchell, Moni Naor, Oscar Nierstrasz, C. Pandu Rangan, Bernhard Steffen, Madhu Sudan, Demetri Terzopoulos, Doug Tygar, Moshe Y. Vardi, Gerhard Weikum, Carlo A. Furia, and Sebastian Nanz, editors, *Objects, Models, Components, Patterns*, volume 7304, pages 353–368. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012. ISBN 978-3-642-30560-3 978-3-642-30561-0. doi: 10.1007/978-3-642-30561-0_24.

Alex Zhitnitsky. The Top 10 Exception Types in Production Java Applications - Based on 1B Events. <https://blog.takipi.com/the-top-10-exceptions-types-in-production-java-applications-based-on-1b-events/>, June 2016.