

UNIVERSITY OF CALGARY

Characterization of Logging Usage:

An Application of Discovering Infrequent Patterns via anti-unification

by

Narges Zirkchianzadeh

A THESIS

SUBMITTED TO THE FACULTY OF GRADUATE STUDIES  
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE  
DEGREE OF MASTER OF SCIENCE

DEPARTMENT OF COMPUTER SCIENCE

CALGARY, ALBERTA

August, 2016

© Narges Zirkchianzadeh 2016

**UNIVERSITY OF CALGARY**  
**FACULTY OF GRADUATE STUDIES**

The undersigned certify that they have read, and recommend to the Faculty of Graduate Studies for acceptance, a thesis entitled “Characterization of Logging Usage: An Application of Discovering Infrequent Patterns via anti-unification” submitted by Narges Zirakchianzadeh in partial fulfillment of the requirements for the degree of MASTER OF SCIENCE.

---

Dr. Robert J. Walker  
Supervisor  
Department of Computer Science

---

Dr. Jörg Denzinger  
Examiner  
Department of Computer Science

---

Dr. Christian J Jacob  
Examiner  
Department of Computer Science

---

Date

# Abstract

Logging has been a common practice to record the runtime behaviour of a software system, typically performed by inserting log statements in source code. So far, several logging frameworks have been specifically created to help developers perform logging tasks but they do not support locating log statements in source code, thus developers usually rely on their common sense to decide on where to log. If logging is properly done, it can provide valuable information for software development and maintenance. On the other hand, ineffective usage of log statements might impose system performance and maintenance overhead. So far, few studies have been conducted to characterize logging usage in real-world applications. This work tries to address the problem of where to log by proposing an automated approach that characterizes the location of log statements through the approximation of an anti-unification (higher-order anti-unification modulo theories) approach and a hierarchical clustering technique to construct a set of anti-unifiers, each describing the commonalities and differences between source code fragments that embody log statements. This approach has been refined in a prototype tool, called ELUS, that greedily identifies the best structural correspondences with respect to the highest similarity and some constraints. I conducted an empirical study through the application of the tool on the source code of four open source systems and manually examined the generated anti-unifiers. My analysis has resulted in ... different main clusters of anti-unifiers in the logging usage. Two empirical evaluations were conducted in this study: (1) an experiment was conducted to validate the effectiveness of the proposed approach through the application of its supporting tool on a test suite. (2) An empirical experiment has been performed to evaluate the quality of the anti-unifiers in describing the location of log statements in source code.

# Acknowledgements

# Table of Contents

Abstract . . . . .	ii
Acknowledgements . . . . .	iii
Table of Contents . . . . .	iv
List of Tables . . . . .	vi
List of Figures . . . . .	vii
List of Symbols . . . . .	viii
1 Introduction . . . . .	1
1.1 Programmatic support for logging . . . . .	3
1.2 Broad thesis overview . . . . .	4
1.3 Thesis statement . . . . .	5
1.4 Thesis organization . . . . .	5
2 Motivational Scenario . . . . .	8
2.1 Summary . . . . .	11
3 Background . . . . .	15
3.1 Concrete syntax trees and abstract syntax trees . . . . .	15
3.2 Eclipse JDT . . . . .	18
3.3 First-order anti-unification . . . . .	22
3.3.1 Higher-order anti-unification . . . . .	25
3.4 Higher-order anti-unification modulo equational theories . . . . .	26
3.5 The Jigsaw Tool . . . . .	29
3.6 Summary . . . . .	31
4 Characterization Study . . . . .	33
4.1 Experiment . . . . .	33
4.1.1 Results . . . . .	34
4.1.2 Analysis . . . . .	35
4.2 Evaluation . . . . .	39
4.3 Summary . . . . .	42
5 Discussion . . . . .	44
5.1 Threats to validity . . . . .	44
5.2 The pitfalls of my tool . . . . .	45
5.2.1 Inaccurate node ordering . . . . .	45
5.2.2 Conflict resolution . . . . .	46
5.3 Applications of anti-unification . . . . .	46
5.4 Summary . . . . .	47
6 Related Work . . . . .	48
6.1 Usage of logging . . . . .	48
6.2 Correspondence . . . . .	50
6.3 API usages patterns . . . . .	50

6.4	Anti-unification . . . . .	51
6.5	Clustering . . . . .	52
6.6	Summary . . . . .	53
7	Conclusion . . . . .	55
7.1	Future Work . . . . .	56

# List of Tables

4.1	The experimental results. . . . .	35
-----	-----------------------------------	----

# List of Figures and Illustrations

1.1	Log statement examples from the Apache Log4j framework. . . . .	3
1.2	Overview of the approach. . . . .	6
2.1	The EditBus class. . . . .	10
2.2	The developer's initial determination of the usage of logging calls. . . . .	12
2.3	The developer's second determination of the usage of log statements. . . . .	12
2.4	The developer's third determination of the usage of log statements. . . . .	13
2.5	The developer's fourth determination of the usage of log statements. . . . .	13
2.6	The developer's final determination of the usage of log statements. . . . .	14
3.1	A simple example Java program. . . . .	16
3.2	The concrete syntax tree for the program of Figure 3.1. . . . .	17
3.3	The abstract syntax tree derived from the concrete syntax tree of Figure 3.2. . . . .	19
3.4	A more abstract AST derived from the concrete syntax tree of Figure 3.2. . . . .	20
3.5	Example 1: A Java method that uses a log statement. . . . .	21
3.6	Example 2: A Java method that uses a log statement. . . . .	21
3.7	Simple AST structure of the examples in Figures 3.5 and 3.6. . . . .	21
3.8	Unification and anti-unification of the terms $f(X, b)$ and $f(a, Y)$ . . . . .	25
3.9	First-order anti-unification of the terms $f(a, b)$ and $g(a, b)$ . . . . .	25
3.10	Higher-order anti-unification of the terms $f(a, b)$ and $g(a, b)$ . . . . .	26
3.11	Higher-order anti-unification modulo theories of the terms $f(a, b)$ and $b$ . . . . .	27
3.12	Complex anti-unification of two structures demonstrating a NIL-theory. . . . .	28
3.13	Ambiguous higher-order anti-unification modulo theories of two terms. . . . .	28
4.1	Summary of the five software systems used in the characterization study. . . . .	34
4.2	Histograms of the number of LMS per cluster. . . . .	35
4.3	The distribution of the categories of anti-unifiers in the logging usage. . . . .	37
4.4	The precision of ELUS. . . . .	41
4.5	The recall of ELUS. . . . .	41



## List of Symbols, Abbreviations and Nomenclature

AST	abstract syntax tree
AU	anti-unification
AUAST	anti-unification abstract syntax tree
HOAUMT	higher-order anti-unification modulo theories
LM	logged method

# Chapter 1

## Introduction

Understanding the similarities and differences between a set of source code fragments is a potentially complex problem that has many actual or potential applications in various software engineering research areas, such as code clones detection [Bulychev and Minea, 2009], automating source code reuse [Cottrell et al., 2008], recommending application programming interface (API) replacements amongst various versions of a software library [Cossette et al., 2014], collating API usage patterns, and automating the merge operation in a version control system. As a specific application, the focus of this study is on characterizing where log statements are used in source code via the determination of structural correspondences between a set of source code fragments enclosing them.

Logging is a conventional programming practice that has been usually used by developers to diagnose the presence or absence of a particular event in a system, to understand the state of an application, and to follow a program's execution flow to find the root causes of an error. The importance of logging is notable in its various applications during software development, such as problem diagnosis [Lou et al., 2010], system behavioural understanding [Fu et al., 2013], quick debugging [Gupta, 2005], performance diagnosis [Nagaraj et al., 2012], easy software maintenance [Gupta, 2005], and troubleshooting [Fu et al., 2009]. Despite the significance of logging for software development and maintenance, few studies have been conducted on understanding its usage in real-world applications, as it has been considered to be a trivial task [Clarke et al., 1999a,b]. However, the availability of several complex frameworks (e.g., Apache Log4j, SLF4J) that assist developers to log suggests that in practice effective logging is not a straightforward task to perform. In addition, a study by Yuan et al. [2012b] showed that developers expend great effort in modifying their logging practices as an afterthought. This indicates that it is not that simple for developers to

perform logging effectively on their first attempt.

The challenges associated with high quality logging arises from the fact that developers are usually left with the burden of deciding on where and what to log manually, thus log statements can be inserted in various locations of source code. For example, a developer may decide to insert log statements at the start and end of every method to record the occurrence of every event of an application. However, three main problems are associated with excessive logging. First, it can produce a lot of redundant information that makes the system log analysis confusing and misleading to perform. Second, excessive logging is costly. It requires extra time and effort to write, debug, and maintain the logging code. Third, it can generate system resource overhead and thus the application performance will be negatively affected. On the other hand, insufficient usage of log statements may result in the loss of run-time information necessary for software analysis. Therefore, logging should be done in an appropriate manner to be effective.

Research on the problem of understanding logging practices can be divided into two main topics: the context and the location of log statements. The context refers to the log text messages, while the location refers to where logs are used in source code. The context of log statements is important to perform high quality logging, as it provides necessary information needed for system analysis. The location of log statements has also a great impact on the quality of logging, as it helps developers to trace the code execution path to identify the root causes of an error within a system. A few studies have been conducted on characterizing log text message modifications [Yuan et al., 2012b] and developing tools to automatically enhance the context of existing log statements [Yuan et al., 2012c, 2010]. Yuan et al. [2012a] proposed Errlog to automatically insert additional log statements into a software system to log all the generic exceptions in order to enhance failure diagnosis. Zhu et al. [2015] applied machine learning techniques to determine the important factors impacting the location of the log statements in source code. In this study, I address the problem of understanding where to log by developing an automated approach that investigates the feasibility of finding patterns in where log statements occur in source code through the construction of a detailed

view of structural generalizations that describe the commonalities and differences between source code fragments containing log statements.

## 1.1 Programmatic support for logging

A typical log statement takes parameters including a log text message and a verbosity level. A log text message consists of static text that describes the logged event and some optional variables related to the event. The verbosity level is intended to classify the severity of a logged event such as a debugging note, a minor issue, or a fatal error. Figure 1.1 provides examples of log statements from the Apache Log4j framework in descending order of severity. The fatal level designates a very severe error event that will likely lead the application to terminate. The error level indicates that a non-fatal but clearly erroneous situation has occurred. The warn level indicates that the application has encountered a potentially harmful situation. The info level designates important information that might be helpful in detecting root causes of an error or in understanding the application behaviour. The debug level provides useful information for debugging an application, and it is usually used by developers only during the development phase. In general, verbosity level is used for classification, in order to avoid the overhead of creating large log files in high performance code.

```
log.fatal ("Fatal Message %s", variable);  
log.error ("Error Message %s", variable);  
log.warn ("Warn Message %s", variable);  
log.info ("Info Message %s", variable);  
log.debug ("Debug Message %s", variable);
```

Figure 1.1: Log statement examples from the Apache Log4j framework.

## 1.2 Broad thesis overview

I aim to create an approach that provides a description of where logs are used in source code by constructing generalizations that represent the detailed structural similarities and differences between methods that make use of log statements, which I call *logged methods* (LMs). In order to evaluate this idea, I implemented the approach to operate on programs written in the Java programming language. To determine how to construct generalizations using the syntax and semantics of the Java programming language, I looked to previous research conducted by Cottrell et al. [2008] that determined the detailed structural correspondences between two Java source code fragments through the application of approximated anti-unification, such that one fragment can be integrated with the other one for small-scale code reuse. However, my problem context is different, as I need to generalize a set of source code fragments with special attention to log statements. Therefore, my approach must take the logs into account when I perform the generalization task via the determination of structural correspondences.

My approach to characterizing logging usage proceeds in four steps (as shown in Figure 1.2). First, potential structural correspondences are determined between the abstract syntax trees (ASTs) of LMs in a pairwise manner, and stored in a novel structure: the *anti-unifier AST* (AUAST), which allows the application of anti-unification on AST structures. Second, I use an approximated anti-unification algorithm to construct a structural generalization (an anti-unifier) representing the commonalities and differences between AUAST pairs, which employs a greedy selection algorithm to approximate the best anti-unifier for the problem by determining the most similar correspondence for each node. The anti-unification algorithm also applies some constraints prior to determining the best correspondences, in order to prevent the anti-unification of log statements with any other types of nodes in the tree structure. The anti-unifier is constructed through the anti-unification of each AUAST node with its best correspondence and then a measure of structural similarity is developed between the two AUASTs. In the third step, I employ a hierarchical clustering algorithm to group the AUASTs into a number of clusters using the structural similarity measure and I then

create a structural generalization from each cluster. The last step involves creating a detailed view of each structural generalization, which I called *logging usage schema* (LUS), that represent the structural commonalities and differences between the set of LMs within each cluster.

To evaluate the approach, I implemented it in a tool called ELUS, written in the Java programming language. I used the Eclipse JDT framework to extract the AST of LMs from a Java program, and employed the Jigsaw framework developed by Cottrell et al. [2008] to find potential structural correspondences. My anti-unifier building tool (built atop Jigsaw) is applied to construct the structural generalizations (Section ??), and my clustering tool is developed atop of it to perform the clustering algorithm described in Section ??.

### 1.3 Thesis statement

The thesis of this work is to characterize where log statements occur in source code by constructing structural generalizations that describe the commonalities and differences between source code fragments containing log statements, thus providing the developers with some guidelines on where to use them effectively in source code.

### 1.4 Thesis organization

The remainder of the thesis is organized as follows.

Chapter 2 motivates the problem of understanding where to use log statements in source code through a scenario in which a developer attempts to perform a logging task. This scenario outlines the potential problems she may encounter and illustrates that the current logging practice is not sufficiently supported.

Chapter 3 provides background information that I build atop: abstract syntax trees (ASTs), which are the basic structure I will use for describing software source code; the Eclipse JDT, an industrial framework for producing and manipulating ASTs for source code written in the Java

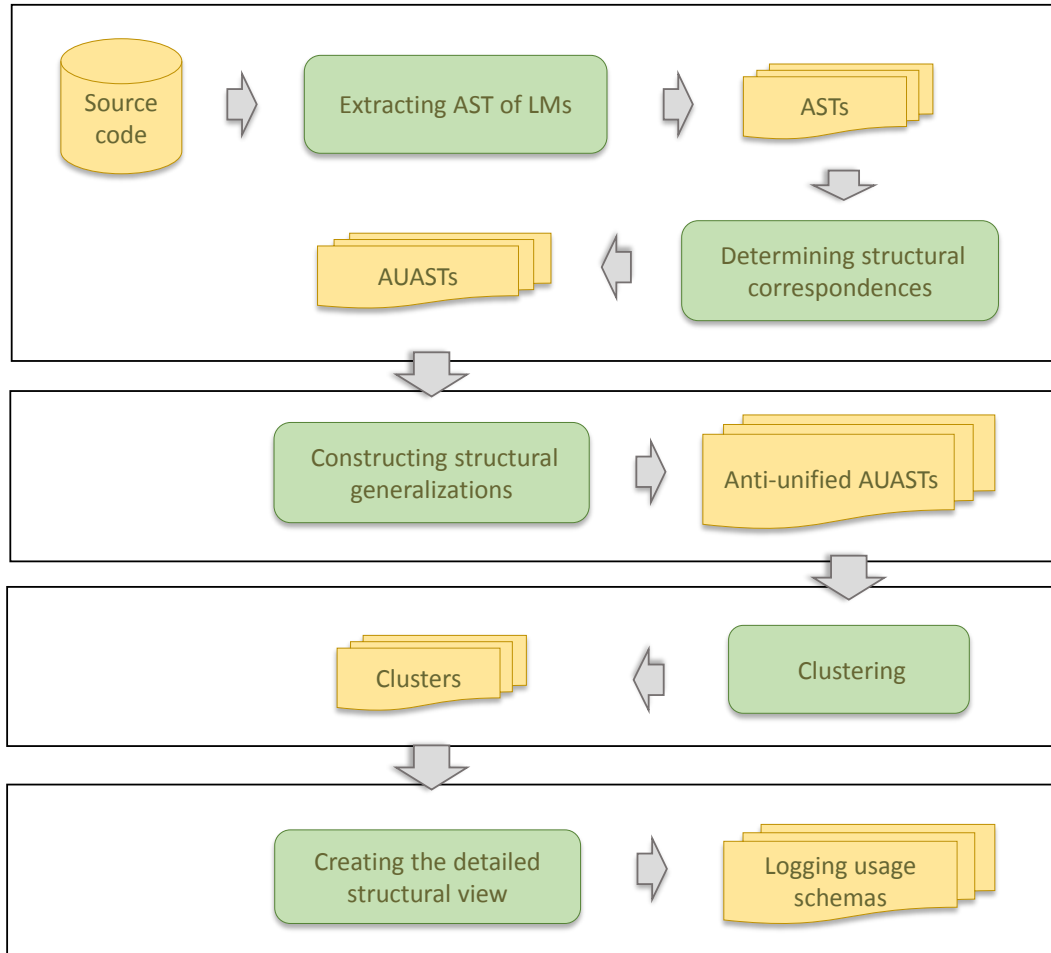


Figure 1.2: Overview of the approach.

programming language; anti-unification, which is a theoretical approach for constructing structural generalizations; and on Jigsaw, a research tool based on the Eclipse JDT for performing anti-unification.

Chapters ??, ??, and ?? present the first three steps of my approach. Determining structural correspondences between AUASTs; constructing structural generalizations from an AUAST pair; and classifying a set of AUASTs into separate clusters, respectively. In each chapter, I discuss the implementation of my approach as an Eclipse plug-in, and conduct an experimental study to assess the effectiveness of my approach through the application of its tool support on a sample test suite

extracted from a real software system.

Chapter 4 presents an empirical study I conducted to characterize the location of log statements of four open-source software systems. Chapter 5 discusses the results and findings of my work, threats to its validity, and remaining issues. Chapter 6 describes work related to my research problem and how it does not adequately address the problem. Chapter 7 concludes the dissertation and presents the contributions of this study and future work.



## Chapter 2

### Motivational Scenario

Printing messages to the console or to a log file is an integral part of software development and can be used to test, debug, and understand what is happening inside an application. In Java programming language, print statements are commonly used to print something on console. However, the availability of tools, frameworks, and APIs for logging that offers more powerful and advanced Java logging features, flexibility, and improvement in the logging quality suggests that using print statements is not sufficient to perform the logging task in real-world applications.

The logging frameworks offer many more facilities that are not provided by the print statements. For example, most logging frameworks (e.g., Log4J, SLF4j, java.util.logging) use different verbosity levels to control the types of information needed to be logged. That is, by logging at a particular verbosity level, logs at that level and higher levels will be recorded whereas the logs with lower levels will be discarded. Each of these verbosity levels can be used for different applications during software development. As an example, the debug log level messages can be used in a test environment, while the error log level messages can be used in a production environment. This feature not only produces fewer log messages for each level, but also improves the performance of an application. Also, most logging frameworks allow the production of formatted log messages, which makes it easier for a developer to monitor the behaviour of a system. In addition, when a developer is working on a server side application, the only way to know what is happening inside the server is by monitoring the log files. Although logging is a valuable practice for software development and maintenance, it imposes extra time and energy on developers to write, test, and run the code, while affecting the application performance. As latency and speed are major concerns for most software systems, it is necessary for a developer to understand and learn logging practices in great detail in order to perform them in an efficient manner.

To illustrate the inherent challenges of effectively performing logging practices in software systems, one may consider a scenario in which a developer is asked to log an event-based mechanism of a text editor tool written in the Java programming language. In this scenario, the developer is trying to log a Java class of the system (Figure 2.1) using the Apache Log4j framework. She knows that components of this application register with the `EditBus` class to receive messages that reflect changes in the application's state, and that the `EditBus` class maintains a list of components that have requested to receive messages. That is, when a message is sent using this class, all registered components receive it in turn. Furthermore, any classes that subscribe to the `EditBus` and implement the `EBComponent` interface define the method `EBComponent.handleMessage(EBMessage)` to handle a message sent on by the `EditBus`. To perform this logging task, the developer might ask herself several fundamental questions, mostly related to where and what to log.

Her first solution might be to simply log at the start and end of every method. However, she believes that logging at the start and end of the `addToBus(EBComponent)`, `removeFromBus(EBComponent)`, and `getComponents()` methods are useless, and will produce redundant information. She assumes that the more she logs, the more she performs file I/O, which slows down the application. Therefore, she decides to log only important information necessary to debug or troubleshoot potential problems. She proceeds to identify the information needed to be logged and then decides on where to use log statements. She thinks that it is important to log the information related to a message sent to a registered component, including the message content and the transmission time, to find the root causes of potential problems in sending messages. She simply wants to begin by using a log statement at the start of the `send()` method (line 2 of Figure 2.2) to log the information. However, she realizes that this log statement does not allow her to log all the information she wants, as the time variable is not initialized at the beginning of this method. Therefore, she proceeds to examine the body of the `send()` method line-by-line and uses another log statement after the time variable is initialized. She aims to log the transmission time in case of potential problems in sending messages. Therefore, she decides to use the logging call inside

```

1 public class EditBus {
2     private static ArrayList components = new ArrayList();
3     private static EBComponent[] copyComponents;
4
5     private EditBus() {
6     }
7
8     public static void addToBus(EBComponent comp) {
9         synchronized(components) {
10             components.add(comp);
11             copyComponents = null;
12         }
13     }
14
15     public static void removeFromBus(EBComponent comp) {
16         synchronized(components) {
17             components.remove(comp);
18             copyComponents = null;
19         }
20     }
21
22     public static EBComponent[] getComponents() {
23         synchronized(components) {
24             if (copyComponents == null) {
25                 EBComponent[] arr = new EBComponent[components.size()];
26                 copyComponents =
27                     (EBComponent[])components.toArray(arr);
28             }
29         }
30         return copyComponents;
31     }
32
33     public static void send(EBMessage message) {
34         EBComponent[] comps = getComponents();
35         for(int i = 0; i < comps.length; i++) {
36             EBComponent comp = comps[i];
37             long start = System.currentTimeMillis();
38             comp.handleMessage(message);
39             long time = (System.currentTimeMillis() - start);
40         }
41     }
42 }

```

Figure 2.1: The EditBus class.

an **if** statement that logs the value of the variable `time`, if it is not within a valid range (shown in lines 9–11 of Figure 2.3).

She also believes that it is important to log an error if any problems occur in sending messages to the components. She decides to use a **try/catch** statement, as it is a common way to handle exceptions in the Java programming language. She creates a **try/catch** block to capture the potential failure in sending messages, and uses a log statement inside the **catch** block to log the exception (shown in lines 2–16 of Figure 2.4). However, she realizes that using this logging call will not allow her to reach the desired functionality, as it does not reveal to which component the problem is related. Thus, she decides to relocate the **try/catch** block inside the **for** statement to log an error in case of a problem in sending messages to any components (shown in lines 5–15 of Figure 2.5).

Figure 2.6 shows the developer’s final determination of the usage of log statements to perform the logging task of the `EditBus` class. By making appropriate decisions about where to use log statements, the developer is in a good position to proceed to write the logging messages by examining the remaining conceptually complex questions. What specific information should I log? How should I choose the log message format? Which information goes to which level of logging? If the developer had reached this point more easily and quickly, she would have had more time and energy to make decisions about the remaining issues and could have completed the logging practice in a timely and appropriate manner.

## 2.1 Summary

This motivational scenario highlights the problems a developer may encounter in performing a logging task. The core problem she faces in this scenario is the difficulty in understanding where to use log statements that enable her to log the desired information. However, having an understanding of how developers usually log in similar situations might assist her to make informed decisions about where to use log statements more effectively, and so she could pay more attention to the remaining, conceptually complex issues to complete the logging task.

```

1 public static void send(EBMessage message){
2     //log statement
3     EBComponent[] comps = getComponents();
4     for (int i = 0; i < comps.length; i++) {
5         EBComponent comp = comps[i];
6         long start = System.currentTimeMillis();
7         comp.handleMessage(message);
8         long time = (System.currentTimeMillis() - start);
9     }
10 }

```

Figure 2.2: The developer's initial determination of the usage of log statements for the send(EBMessage) method.

```

1 public static void send(EBMessage message) {
2     //log statement
3     EBComponent[] comps = getComponents();
4     for(int i = 0; i < comps.length; i++) {
5         EBComponent comp = comps[i];
6         long start = System.currentTimeMillis();
7         comp.handleMessage(message);
8         long time = (System.currentTimeMillis() - start);
9         if (time >= 1000000) {
10            //log statement
11        }
12    }
13 }

```

Figure 2.3: The developer's second determination of the usage of log statements for the send(EBMessage) method.

```

1 public static void send(EBMessage message){
2     try {
3         //log statement
4         EBComponent[] comps = getComponents();
5         for(int i = 0; i < comps.length; i++) {
6             EBComponent comp = comps[i];
7             long start = System.currentTimeMillis();
8             comp.handleMessage(message);
9             long time = (System.currentTimeMillis() - start);
10            if (time >= 1000000) {
11                //log statement
12            }
13        }
14    } catch(Throwable t) {
15        //log statement
16    }
17 }

```

Figure 2.4: The developer's third determination of the usage of log statements for the send(EBMessage) method.

```

1 public static void send(EBMessage message) {
2     //log statement
3     EBComponent[] comps = getComponents();
4     for (int i = 0; i < comps.length; i++) {
5         try {
6             EBComponent comp = comps[i];
7             long start = System.currentTimeMillis();
8             comp.handleMessage(message);
9             long time = (System.currentTimeMillis() - start);
10            if (time >= 1000000) {
11                //log statement
12            }
13        } catch(Throwable t) {
14            //log statement
15        }
16    }
17 }

```

Figure 2.5: The developer's fourth determination of the usage of log statements for the send(EBMessage) method.

```

1 public class EditBus {
2     private static ArrayList components = new ArrayList();
3     private static EBComponent[] copyComponents;
4
5     private EditBus() {
6     }
7
8     public static void addToBus(EBComponent comp) {
9         synchronized(components) {
10             components.add(comp);
11             copyComponents = null;
12         }
13     }
14
15     public static void removeFromBus(EBComponent comp) {
16         synchronized(components) {
17             components.remove(comp);
18             copyComponents = null;
19         }
20     }
21
22     public static EBComponent[] getComponents() {
23         synchronized(components) {
24             if (copyComponents == null) {
25                 EBComponent[] arr = new EBComponent[components.size()];
26                 copyComponents = (EBComponent[])components.toArray(arr);
27             }
28         }
29         return copyComponents;
30     }
31
32     public static void send(EBMessage message) {
33         //log statement
34         EBComponent[] comps = getComponents();
35         for(int i = 0; i < comps.length; i++) {
36             try {
37                 EBComponent comp = comps[i];
38                 long start = System.currentTimeMillis();
39                 comp.handleMessage(message);
40                 long time = (System.currentTimeMillis() - start);
41                 if (time >= 1000000) {
42                     //log statement
43                 }
44             } catch(Throwable t) {
45                 //log statement
46             }
47         }
48     }
49 }

```

Figure 2.6: The developer's final determination of the usage of log statements for the EditBus class.

# Chapter 3

## Background

A programming language is described by the combination of its syntax and semantics. The syntax concerns the legal structures of programs written in the programming language, while the semantics is about the meaning of every construct in that language. Furthermore, the abstract syntactic structure of source code written in a programming language can be represented as an *abstract syntax tree* (AST), in which nodes are occurrences of syntactic structures and edges represent nesting relationships. Since ASTs will be the form in which I represent and analyze source code, I need a means to generalize sets of ASTs in order to understand their commonalities while abstracting away their differences. The theoretical framework of anti-unification is presented as that means.

In this chapter, ASTs are described in Section 3.1, along with their more concrete counterparts, concrete syntax trees. A specific, industrial framework for creating and manipulating ASTs for source code written in the Java programming language—the Eclipse JDT—is described in Section 3.2. Anti-unification is summarized in Section 3.3, starting with its most basic form, first-order anti-unification, and progressing to the form that I will make use of, higher-order anti-unification modulo equational theories, in Section 3.4. A research approach, built atop the Eclipse JDT, for performing anti-unification on Java ASTs—the Jigsaw framework—is described in Section 3.5.

### 3.1 Concrete syntax trees and abstract syntax trees

A concrete syntax tree is a tree (i.e., a kind of graph)  $T = (V, E)$  whose vertices  $V$  (equivalently, nodes) represent the syntactic structures (equivalently, syntactic elements) of a specific program written in a specific programming language and whose directed edges  $E$  represent the nesting relationships amongst those syntactic structures. Non-leaf nodes in a concrete syntax tree (also called a parse tree) represent the grammar productions that were satisfied in parsing the program it



```

1 public class HelloWorld {
2     public static void main(String[] args) {
3         System.out.println("Hello world!");
4     }
5 }

```

Figure 3.1: A simple example Java program.

represents; leaf nodes represent the concrete lexemes, such as literals and keywords.

I focus on the Java programming language and I make use of the grammar in the language specification [Gosling et al., 2012, Chapter 18] to determine the form of the concrete syntax trees. Non-leaf node names are represented by names in “camel-case” written in *italics*. Consider the trivial program in Figure 3.1; its concrete syntax tree is represented in Figure 3.2.

Beyond the fact that the concrete syntax tree is rather verbose and thus occupies a lot of space even for a trivial example, I can see two key problems with it: (1) there are a multitude of redundant nodes such as *expression1*, *expression2*, and *expression3* that are present solely for purposes of creating an unambiguous grammar; and (2) there are no nodes that express key concepts, such as “method declaration” and “method invocation”, that should be obviously present in the example program.

To address these problems, concrete syntax trees are converted to abstract syntax trees (ASTs). An AST is similar in concept to a concrete syntax tree but it does not generally represent the parsing steps followed to differentiate different kinds of syntactic structure. The node types are chosen to represent syntactical concepts; I use the grammar presented for exposition by Gosling et al. [2012], which differs markedly from the grammar they propose in their Chapter 18 for efficient parsing. Note that a given node type constrains the kinds and numbers of child nodes that it possesses. The AST derived from the concrete syntax tree of Figure 3.2 is shown in Figure 3.3. Note that, although I know that (for any normal program) `System` refers to the class `java.lang.System` and `out` is a static field on that class, non-normal programs can occur and a pure syntactic analysis cannot rule out that `System` is a package and that `out` is a class therein declaring a static method



Figure 3.2: The concrete syntax tree for the program of Figure 3.1.

`println (String).`

This is still verbose, so in practice we elide details that are implied or otherwise trivial, to arrive at a more abstract AST as shown in Figure 3.4.

## 3.2 Eclipse JDT

The Eclipse Java Development Tools (JDT) framework provides APIs to access and manipulate Java source code via ASTs. An AST represents Java source code in a tree form, where the typed nodes represent instances of certain syntactic structures from the Java programming language. Each node type (in general) takes a set of child nodes, also typed and with certain constraints on their properties. Groups of children are named on the basis of the conceptual purpose of those groups; optional groups can be empty, which we can represent with the `NIL` element. For example, the simple AST structure of two sample LMs in Figures 3.5 and 3.6 is shown in Figure 3.7, with the `log` statements highlighted in yellow.

In the JDT framework, structural properties of each AST node can be used to obtain specific information about the Java element that it represents. These properties are stored in a map data structure that associates each property to its value; this data is divided into three types:

- *Simple structural properties:* These contain a simple value which has a primitive or simple type or a basic AST constant (e.g., identifier property of a name node whose value is a `String`). For example, all the *identifier* nodes in Figure 3.3 fall in this case; each references an instance of `String` representing the string that constitutes the identifier.
- *Child structural properties:* These involve situations where the value is a single AST node (e.g., name property of a method declaration node). For example, the *classDeclaration* node in Figure 3.3 has a single child that represents its name as an *identifier* node.
- *Child list structural properties:* These involve situations where the value is a list of

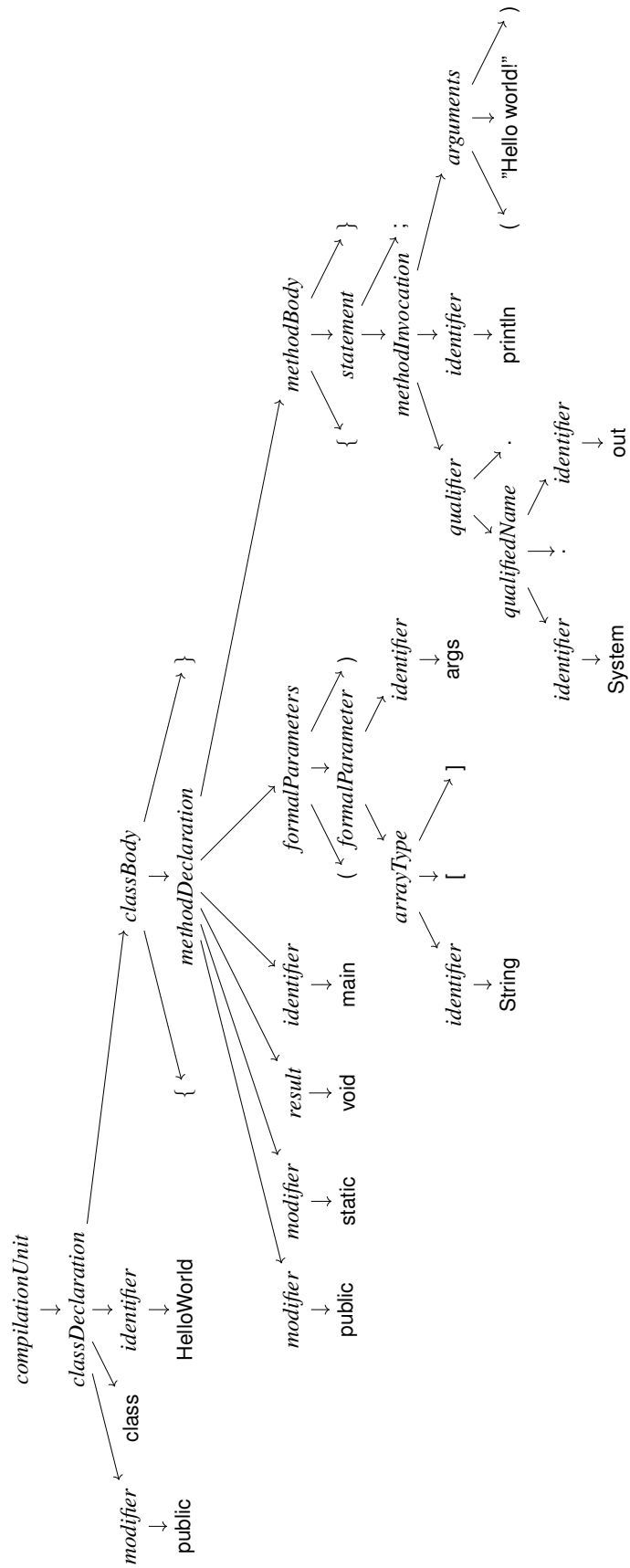


Figure 3.3: The abstract syntax tree derived from the concrete syntax tree of Figure 3.2.

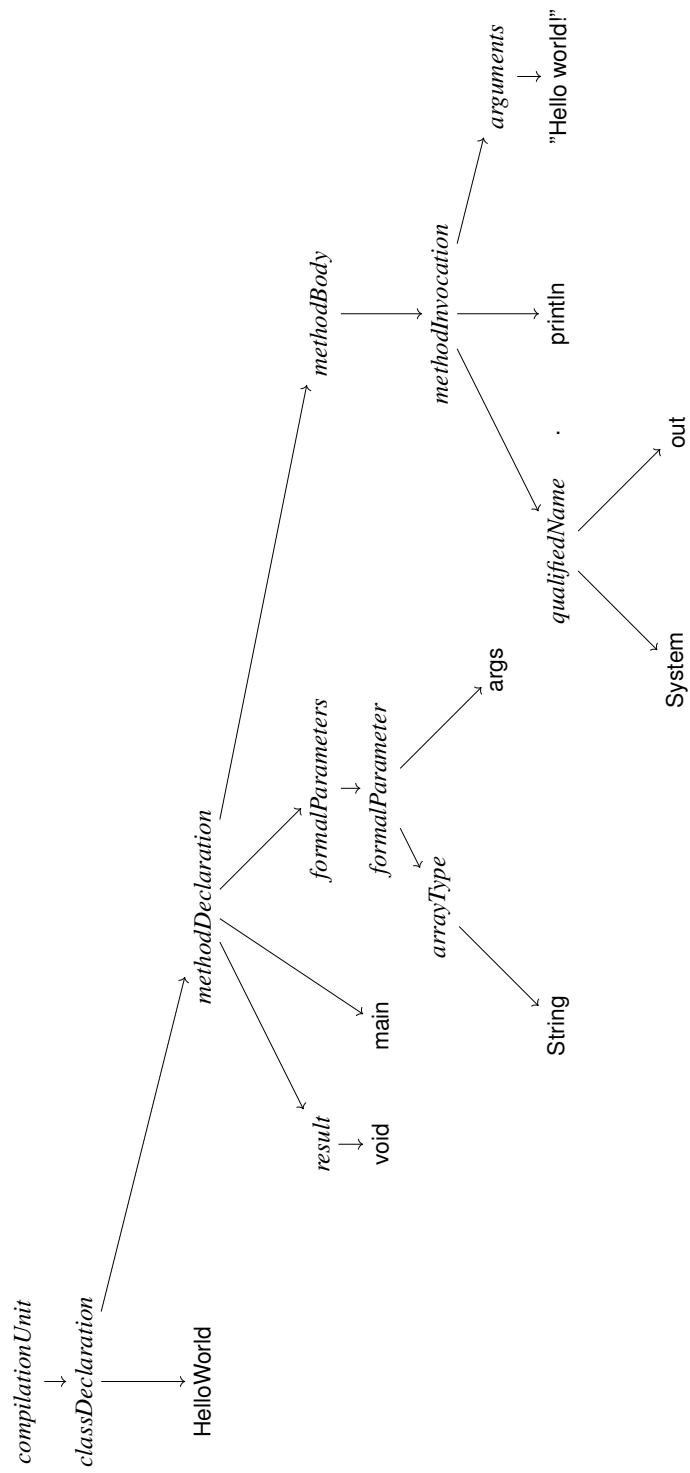


Figure 3.4: A more abstract AST derived from the concrete syntax tree of Figure 3.2.

```

1 public void handleMessage(EBMessage message) {
2     if (seenWarning)
3         return;
4     seenWarning = true;
5     Log.log(Log.WARNING, this, getClassName() + " should extend EditPlugin not EBPlugin
        since it has an empty " + handleMessage());
6 }

```

Figure 3.5: A Java method that uses a log statement. This will be referred to as Example 1.

```

1 public void actionPerformed(ActionEvent evt) {
2     EditAction action = context.getAction(actionName);
3     if (action == null) {
4         Log.log(Log.WARNING, this, "Unknown action: " + actionName);
5     }
6     else{
7         context.invokeAction(evt, action);
8     }
9 }

```

Figure 3.6: A Java method that uses a log statement. This will be referred to as Example 2.

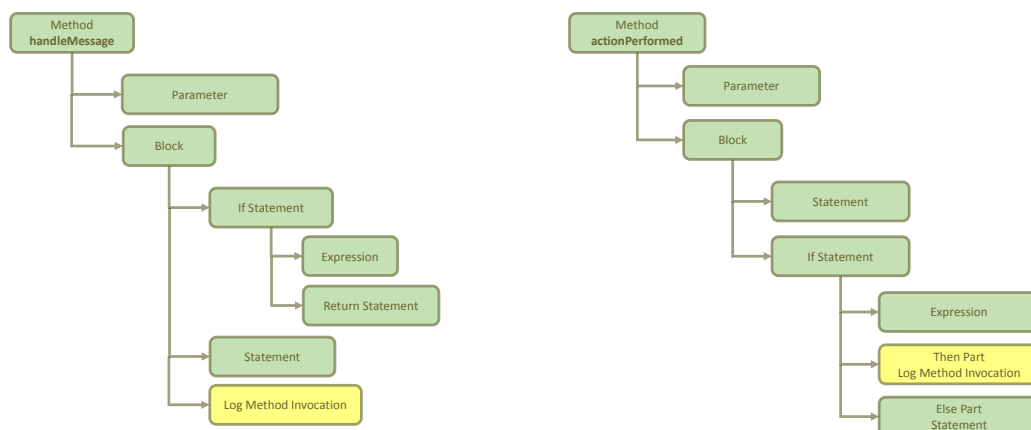


Figure 3.7: Simple AST structure of the examples in Figures 3.5 and 3.6.

child nodes. For example, the *classDeclaration* node in Figure 3.3 can possess multiple *modifiers*.

As an example, the ASTs of the log statements at line 4 of Figure 3.5 and Figure 3.6 can be represented respectively as:

- *methodInvocation*(  
*qualifiedName*(Log, *identifier*(log)),  
*arguments*(  
*qualifiedName*(Log, *identifier*(WARNING)),  
*thisExpression*(),  
*additionExpression*(  
*methodInvocation*(*identifier*(getClassName), *arguments*()),  
*stringLiteral*(" should extend EditPlugin not EBPlugin since it has an empty "),  
*methodInvocation*(*identifier*(handleMessage), *arguments*()))))
- *methodInvocation*(  
*qualifiedName*(Log, *identifier*(log)),  
*arguments*(  
*qualifiedName*(Log, *identifier*(WARNING)),  
*thisExpression*(),  
*additionExpression*(  
*stringLiteral*("Unknown action: "),  
*identifier*(actionName))))

### 3.3 First-order anti-unification

This section defines terms, substitutions, applying a substitution to a term, and instances and anti-instances of a term, as the requirements needed to describe anti-unification theory (and its dual,

unification theory).

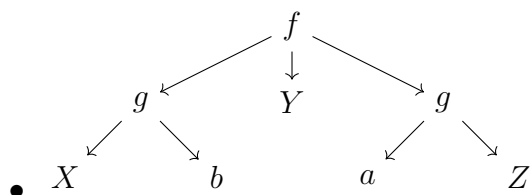
**Definition 3.3.1** (Term). A (first-order) term is defined to be a variable, a constant, or a function symbol followed by a list of terms as the arguments of the function. [Note that function symbols without the subsequent list of terms do not constitute first-order terms.]

Function symbols taking  $n$  arguments are called  $n$ -ary function symbols; 0-ary function symbols are called constants. The identifiers starting with a lowercase letter are used to represent function symbols (e.g.,  $f(a, b)$ ,  $g(a, b)$ ) and constants (e.g.,  $a$ ,  $b$ ), while variables are represented by identifiers starting with an uppercase letter (e.g.,  $X$ ,  $Y$ ). The following are examples of a term:

- $Y$
- $a$
- $f(X, c)$
- $f(g(X, b), Y, g(a, Z))$

Note that for any term there is a unique, equivalent tree and vice versa: constants and (first-order) variables are leaf nodes, while function symbols are non-leaf nodes; a function with given arguments is represented by a non-leaf node (representing the function symbol) with directed edges pointing to leaf nodes representing each argument. For example:

- $Y$
- $a$
- $X \leftarrow f \rightarrow c$





**Definition 3.3.2** (Substitution). A substitution is a set of mappings, each from a variable to a term.

**Definition 3.3.3** (Applying a substitution). Applying a substitution to a term results in the replacement of all occurrences of each variable in the term, by its corresponding term as defined in the substitution.

As an example, applying the substitution  $\Theta = \{X \rightarrow a, Y \rightarrow b\}$  to the term  $f(X, Y)$  results in the replacement of all occurrences of the variable  $X$  by the term  $a$  and all occurrences of the variable  $Y$  by the term  $b$ , and thus  $f(X, Y) \xrightarrow{\Theta} f(a, b)$ .

**Definition 3.3.4** (Instance & anti-instance).  $a$  is an instance of a term  $X$  and  $X$  is an anti-instance of  $a$ , if there is a substitution  $\Theta$  such that applying  $\Theta$  to  $X$  results in  $a$  (i.e.,  $X \xrightarrow{\Theta} a$ ).

**Definition 3.3.5** (Unifier). A unifier is a common instance of two given terms.

Unification usually aims to create the *most general unifier* (MGU); that is,  $U$  is the MGU of two terms such that for all unifiers  $U'$  there exists a substitution  $\Theta$  such that  $U \xrightarrow{\Theta} U'$ . Unification aims to make a more concrete structure in essence, whereas what we need is a more generalized structure, which leads to the use of the dual of unification, called *anti-unification*.

**Definition 3.3.6** (Anti-unifier).  $X$  is an anti-unifier (or generalization) for  $a$  and  $b$ , if  $X$  is an anti-instance for  $a$  and an anti-instance for  $b$  under substitutions  $\Theta_1$  and  $\Theta_2$ , respectively (i.e.,  $X \xrightarrow{\Theta_1} a$  and  $X \xrightarrow{\Theta_2} b$ ).

An anti-unifier contains common pieces of the original terms, while the differences are abstracted away using variables. An anti-unifier for a pair of terms always exists since we can anti-unify any two terms by the anti-instance  $X$ , i.e., a single variable. However, anti-unification usually aims to find the *most specific anti-unifier* (MSA), that is,  $A$  is the MSA of two structures where there exists no anti-unifier  $A'$  such that  $A \xrightarrow{\Theta} A'$ .

As an example, the anti-unifier of two given terms  $f(X, b)$  and  $f(a, Y)$  is the new term  $f(X, Y)$ , containing common pieces of the two original terms. The variable  $Y$  in the anti-unifier  $f(X, Y)$

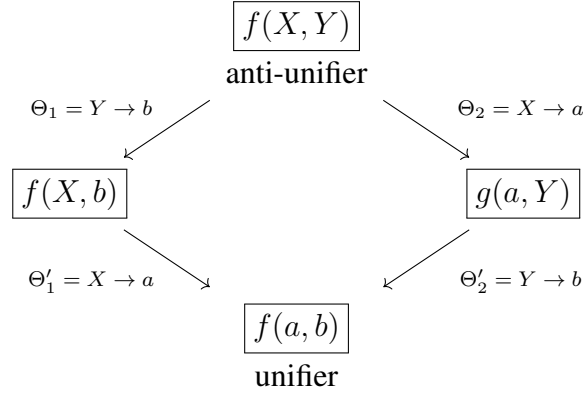


Figure 3.8: Unification and anti-unification of the terms  $f(X, b)$  and  $f(a, Y)$ .

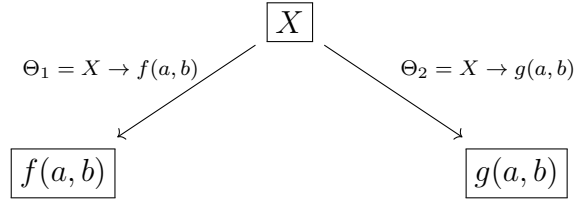


Figure 3.9: First-order anti-unification of the terms  $f(a, b)$  and  $g(a, b)$ .

can be substituted by the term  $b$  to re-create  $f(X, b)$  (with  $\Theta_1 = Y \rightarrow b$ ) and the variable  $X$  in the anti-unifier can be substituted by the term  $a$  to re-create  $f(a, Y)$  (with  $\Theta_2 = X \xrightarrow{\Theta} a$ ), as depicted in Figure 3.8. In addition, the unifier  $f(a, b)$  of the two terms can be instantiated by applying the substitutions  $\Theta'_1 = X \xrightarrow{\Theta} a$  and  $\Theta'_2 = Y \xrightarrow{\Theta} b$  on the terms  $f(X, b)$  and  $f(a, Y)$ , respectively.

The MSA should preserve as much of common pieces of both original terms as possible; however, first-order anti-unification fails to capture complex commonalities as it restricts substitutions to only replace first-order variables by terms. That is, when two terms differ in function symbols, first-order anti-unification fails to capture common details of them. For example, the first-order anti-unifier of the terms  $f(a, b)$  and  $g(a, b)$  is  $X$  as depicted in Figure 3.9.

### 3.3.1 Higher-order anti-unification

Higher-order anti-unification would allow us to create the MSA by extending the set of possible substitutions such that variables can be replaced not only by terms but also by function symbols in

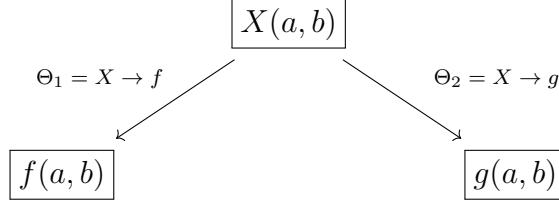


Figure 3.10: Higher-order anti-unification of the terms  $f(a, b)$  and  $g(a, b)$ .

order to retain the detailed commonalities. For example, the higher-order anti-unifier of the terms  $f(a, b)$  and  $g(a, b)$  is  $X(a, b)$  as depicted in Figure 3.10.

Applying higher-order anti-unification could help to construct a structural generalization by maintaining the common pieces and abstracting the differences away using variables. However, it is not comprehensive enough to solve the problem as it does not consider background knowledge about AST structures, such as syntactically different but semantically equivalent structures, missing structures, and different ordering of arguments.

### 3.4 Higher-order anti-unification modulo equational theories

In higher-order anti-unification modulo (equational) theories, a set of equational theories, which treat different structures as equivalent, is defined to incorporate background knowledge. Each equational theory  $=_E$  determines which terms are considered equal and a set of these equations can be applied on higher-order extended structures to determine structural equivalences. For example, we have introduced an equivalence equation  $=_E$ , such that  $f(X, Y) =_E f(Y, X)$  to indicate that the ordering of arguments does not matter in our context.

We have also introduced a theory, called NIL-theory, that adds the concept of a NIL structure, which permits a structure to be equated with nothing, and defines an equivalence equation  $=_E$  for it. The NIL structure can be used to anti-unify two structures when a substructure exists in one but is missing from the other. However, some requirements should be taken to avoid the overuse of NIL structures such that the original structures must have common substructures but vary in the size for dissimilar substructures. For example, we can anti-unify the two structures  $b$  and  $f(a, b)$

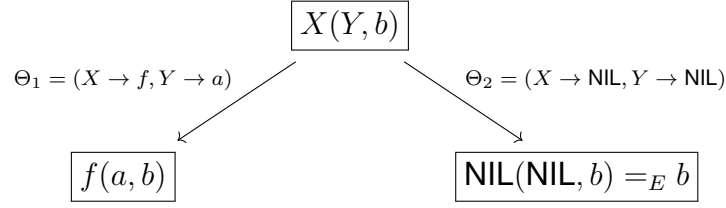


Figure 3.11: Higher-order anti-unification modulo theories of the terms  $f(a, b)$  and  $b$ .

through the application of NIL-theory by creating the term  $\text{NIL}(\text{NIL}, b)$  which is  $=_E$  to  $f(b)$  and anti-unifying  $\text{NIL}(\text{NIL}, b)$  with  $f(a, b)$  as depicted in Figure 3.11.

We have also defined a set of equivalence equations to incorporate semantic knowledge of structural equivalences supported by the Java language specification, as it provides various ways to define the same language specifications. These theories should be applied on higher-order extended structures to anti-unify AST structures that are not identical but are semantically equivalent. For example, consider **for**- and **while**-statements that are two types of looping structure in Java programming language: they have different syntax but semantically cover the same concept. Let us look at the code snippets **for**( $i=0; i<10; i++$ ) and **while**( $i<10$ ), whose ASTs can be represented as **for**(*initializer*( $i, 0$ ); *lessThanExpression*( $i, 10$ ); *updaters*(*postIncrementExpression*( $i$ ))) and **while**(*lessThanExpression*( $i, 10$ )), respectively. We could define an equivalence equation  $=_E$  that allows the anti-unification of **for**- and **while**-statements. We also need to utilize the NIL-theory to handle the varying number of arguments as the **for**-loop has three arguments whereas the **while**-loop only has one. Using the NIL-theory we can create the structure **while**( $\text{NIL}(\text{NIL}, \text{NIL}), \text{lessThanExpression}(i, 10), \text{NIL}(\text{NIL}, \text{NIL}))$  that is  $=_E$  to **while**(*lessThanExpression*( $i, 10$ )) and construct the anti-unifier,  $V_0(V_1(V_2, V_3), \text{lessThanExpression}(i, 10), V_4(V_5(V_2)))$ , as depicted in Figure 3.12.

However, defining complex substitutions in higher-order anti-unification modulo theories results in losing the uniqueness of the MSA. For example, consider the terms  $f(g(a, e))$  and  $f(g(a, b), g(d, e))$ . As described in Figure 3.13, two MSAs exist for these terms: we can anti-unify  $g(a, e)$  and  $g(a, b)$  to create the anti-unifier  $g(a, X_0)$  and anti-unify  $g(d, e)$  with the NIL structure to create the anti-unifier  $Y(Z, X_1)$ ; or we can anti-unify  $g(a, e)$  and  $g(d, e)$  to create the anti-unifier  $g(X_0, e)$

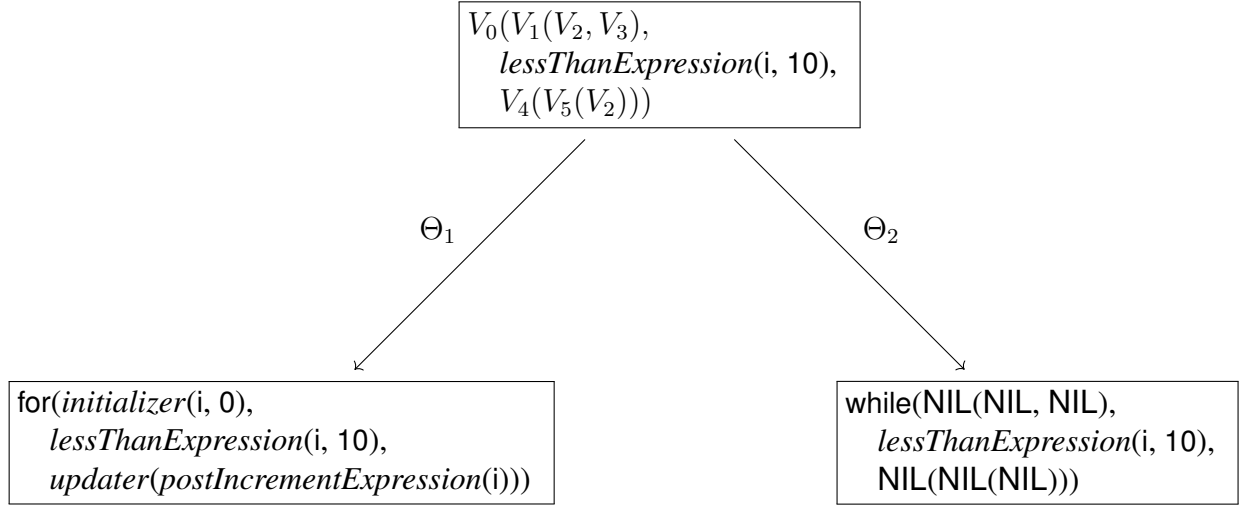


Figure 3.12: Anti-unification of the structures **for**(*initializer*(*i*, 0), *lessThanExpression*(*i*, 10), *updater*(*postIncrementExpression*(*i*))) and **while**(*NIL*(*NIL*, *NIL*), *lessThanExpression*(*i*, 10), *NIL*(*NIL*, *NIL*)). The substitutions are defined as follows:  $\Theta_1 = (V_0 \rightarrow \text{for}, V_1 \rightarrow \text{initializer}, V_2 \rightarrow i, V_3 \rightarrow 0, V_4 \rightarrow \text{updater}, V_5 \rightarrow \text{postIncrementExpression})$ ; and  $\Theta_2 = (V_0 \rightarrow \text{while}, V_1 \rightarrow \text{NIL}, V_2 \rightarrow \text{NIL}, V_3 \rightarrow \text{NIL}, V_4 \rightarrow \text{NIL}, V_5 \rightarrow \text{NIL})$

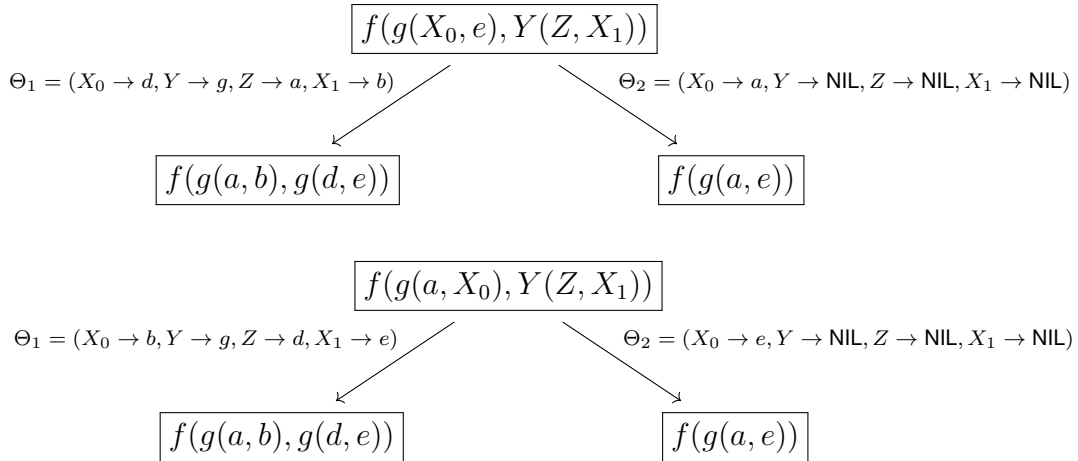


Figure 3.13: Ambiguous higher-order anti-unification modulo theories of the terms  $f(g(a, b), g(d, e))$  and  $f(g(a, e))$ , creating multiple MSAs.

and anti-unify  $g(a, b)$  with the NIL structure to create the anti-unifier  $Y(Z, X_1)$ .

Despite having multiple potential MSAs, I need to determine one single MSA that is the most appropriate in my context. However, the complexity of finding an optimal MSA is undecidable in general [Cottrell et al., 2008] since an infinite number of possible substitutions can be applied to variables in a term. Therefore, I need to use an approximation technique to construct one of the best MSAs that can sufficiently solve the problem.

### 3.5 The Jigsaw Tool

Jigsaw is a plug-in to the Eclipse integrated development environment (IDE), which was developed by Cottrell et al. [2008] to support small-scale source code reuse via structural correspondence. A small-scale reuse task can be divided into two phases. The first phase involves the developer identifying a source code snippet that implements functionality that is missing within a target system. The second phase involves integrating the source code snippet within the target system. Jigsaw supports the small scale reuse task by identifying structural correspondences between the code snippet and the context into which the code should be pasted, in order to suggest to developers what parts already exist within the target system, what parts are missing, and what parts need to be modified to fit into the context. In summary, the Jigsaw tool determines the structural correspondences between two Java source code fragments through the application of higher-order anti-unification modulo equational theories such that one fragment can be integrated to the other one for small-scale code reuse.

In general, the proposed approach by Cottrell et al. proceeds in three steps. First, it generates an augmented form of AST, called a *correspondence AST* (CAST), where each node holds a list of candidate correspondence connections, each implicitly representing an anti-unifier. To find candidate correspondences amongst the CASTs of the original system and the target system, it uses a similarity measure that relies on syntactic similarity along with simple knowledge of semantic equivalences supported the Java language specifications. Although the CAST structure may repre-

sent many anti-unifiers, they used a greedy selection algorithm to select the best fit for each node via thresholding in order to approximate the optimal generalization. That is, the correspondence connections with a similarity value below a threshold are removed. Second, when there is more than one candidate correspondence connection for a node, the developer is prompted to resolve the conflict by selecting the best fit for his functionality. Third, the best correspondences are used to semi-automatically perform the integration task by replacing the references to variables in the original system by the references to variables in the target system. The Jigsaw tool is a proof-of-concept implementation of this approach.

Underlying the Jigsaw tool is the Jigsaw framework for determining likely structural correspondences between two ASTs; I simply refer to “Jigsaw” henceforth to intend the Jigsaw framework. The Jigsaw similarity function returns a value in  $[0, 1]$  where zero indicates complete lack of similarity and one indicates perfect similarity. In general, this function returns a value above zero if the compared nodes are of identical type, and thus it returns a similarity of 0 for the nodes of different types. However, it uses several heuristics to improve the utility of the similarity measurement by defining an arbitrary value for the nodes that are syntactically different but are semantically relevant. For example, the similarity between names of AST nodes is measured using a normalized computation based on the length of the longest common substring. The comparison of the **int** and **long** nodes is another example, where an arbitrary value of 0.5 is defined as the similarity, since they are not of syntactically identical types but have a semantic equivalence. This function also detects the correspondence between **for**-, enhanced-**for**-, **while**-, and **do**-loop statements; and **if** and **switch** conditional statements.

As I intend to construct a structural generalization from ASTs of two logged methods via structural correspondence, it could be helpful to use the first phase of the proposed approach to find candidate correspondences using the similarity measure. However, the second phase does not help determine the best correspondences needed in my context, as the CAST generated via thresholding neither resolves the conflicts that occur in constructing one single anti-unifier automatically, nor

prevents the anti-unification of log statement nodes with any other nodes. There, the Jigsaw similarity function does not enable us to measure how similar are the usages of log statements inside methods. In addition, as the problem of this study is different, the integration phase of the approach is not related to my work. Instead, I should develop an algorithm to construct a detailed view of the generalization describing the structural commonalities and differences between logged methods. However, the CAST structure does not suffice to construct an anti-unifier: it does not allow the insertion of *structural variables* in place of nodes in the tree structure, and thus an extended form is required. In the following chapters, I will discuss my approach to create a structural generalization and its implementation by means of the higher-order anti-unification modulo theories.

### 3.6 Summary

I described abstract syntax trees (ASTs) as a standard syntactic representation of source code. Every AST can also be represented in a functional format (and vice versa) which constitute the standard theoretical concept of terms. I presented Eclipse JDT as a concrete framework that can be used to manipulate ASTs of a source code written in the Java programming language.

I demonstrated how the theoretical framework of anti-unification can be used as a technique to construct a common generalization of two given terms, and hence of two ASTs. First-order anti-unification permits terms to be replaced with variables and vice versa, but it is limited in that low-level commonality can be discarded due to high-level differences. Higher-order anti-unification overcomes this by permitting substitution relative to function symbols as well as terms. A further extension allows for insertion and deletion by declaring equivalence with the NIL structure, as well as other arbitrary equational theories to embed knowledge of semantic equivalence. Unfortunately, this approach of higher-order anti-unification modulo theories leads to ambiguity and the potential for an infinite number of possible substitutions for every structural variable. To make use of that technique despite its weaknesses, we must apply an approximation technique to select amongst the best MSAs in order to reach a solution that is reasonable in practice. I also introduced Jigsaw, an



existing framework for determining structural correspondences between ASTs and why it does not adequately address my problem. To address my problem, I describe in subsequent chapters how I extended the Jigsaw framework.

# Chapter 4

## Characterization Study

To characterize the location of log statements in source code, I conducted an experimental study that addresses the following research questions:

- RQ1: *“Is it possible to find patterns of where log statements occur in source code?”* I aim to investigate whether there are clusters containing a large number of LMs. This suggests that there might be common ways of locating log statements in source code.
- RQ2: *“What common structural characteristics do logged methods have?”* I conducted a manual analysis on the detailed view of structural generalizations produced by ELUS to identify the common structural characteristics of LMs in each cluster.

### 4.1 Experiment

In this experiment, I will analyze logging usage of five popular open-source software systems: Apache Tomcat, Hibernate ORM, Apache Camel, Apache Solr, and .... Each system is written in the Java programming language and they all utilize the same logging framework, Apache Log4j. I decided to study the usage of log4j statements in these systems, as Apache Log4j is ranked as the most commonly used logging package for Java<sup>1</sup>. The studied systems are from different application domains: Apache Tomcat is a Java Servlet; Hibernate ORM is a object relational-mapping framework; Apache Camel is a rule-based routing and mediation engine; and Apache Solr is an enterprise search platform. I chose these systems as my study subject due to their popularity in their area of application (7000+ commits to the GitHub repository) and long history of development (9 to 13 years). Table 4.1 represents the details about these software systems. I also decided

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Java\\_logging\\_framework](https://en.wikipedia.org/wiki/Java_logging_framework)

to exclude the log4j statements at the trace- and debug- verbosity level, as they are usually used by developers only during the software development phase. I believe that studying these systems could give us an insight about logging usage in real-world applications.

Software system	Description	Version	Start time	LOC	Log statements
Tomcat	Server	9.0.11	2003	306,704	3,117
Hibernate ORM	Framework	4.2.23	2004	509,734	1,939
Camel	Middleware	2.18.0	2007	120,528	2,177
Solr	Platform	6.2.1	2007	128,824	2,319

Figure 4.1: Summary of the five software systems used in the characterization study.

My proof-of-concept implementation takes the source code of these systems as input, extracts the ASTs of their LMs, applies the proposed algorithm to construct AUASTs, classifies the AUASTs into clusters, and outputs the detailed view of structural generalization (LUS) for each cluster.

#### 4.1.1 Results

The experimental results for each software system are presented in Table 4.1. This table describes the total number of detected log4j statements (debug- and trace-level log statements are excluded), the number of logged methods (LMs); the number of generated clusters; the number of generalized clusters containing more than one LM; the number of singleton clusters that only contain one LM; and the reduction percentage calculated by the Equation 4.1. In addition, Figure 4.2 shows the histograms of the number of LMs per cluster for each system.

$$reduction = \frac{|Primitive\ clusters| - |Clusters|}{|Primitive\ clusters|} \quad (4.1)$$

	Tomcat	Hibernate	Camel	Solr
Log4j statements	1098	128	632	1484
LMs	658	81	490	818
Primitive clusters	1098	128	632	1484
Clusters	43	12	49	50
Generalized clusters	20	9	18	27
Singleton clusters	23	3	31	23
Reduction	96%	91%	92%	96%

Table 4.1: The experimental results.

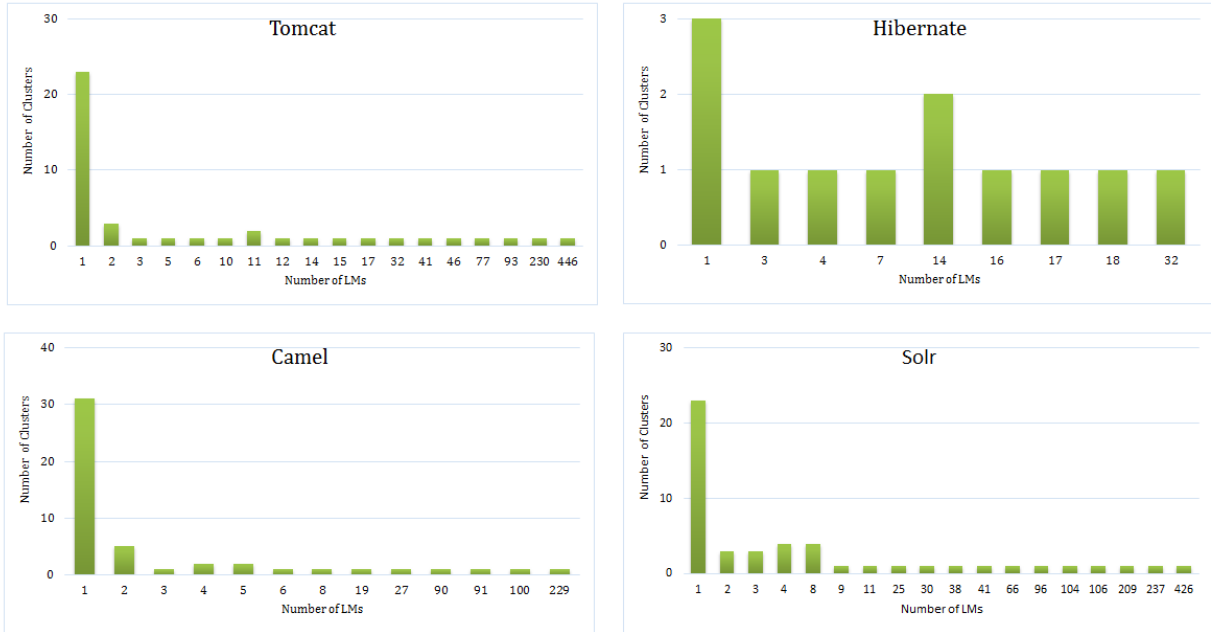


Figure 4.2: Histograms of the number of LMS per cluster.

#### 4.1.2 Analysis

The first research question is : *"Is it possible to find patterns of where log statements occur in source code?"*. As it is shown in Table 4.1, the number of clusters has been reduced by more than 90% in all the studied systems, indicating that developers follow some patterns for locating the log statements in source code. Furthermore, histograms depicted in Figure 4.2 show that in all the

studied systems, a few clusters contain a large number of LMs; however, the other clusters contain a very few number of LMs (only one or two). This indicates that in these cases developers follow a more complex or rare way of locating log statements.

The second research question is : “*What common structural characteristics do logged methods have?*” To address this question, I have manually went through the LUSs of anti-unifiers to identify the common structural characteristics of locating the log statements in source code.

### *Categorizing logging usage*

In this section, I will describe the anti-unifiers of logging usage by examining the LUSs produced by ELUS. In general, there are six categories of anti-unifiers in the logging usage. In the following sections, I will describe the common structural characteristics of each category represented by the anti-unifiers. In addition, Figure 4.3 presents the number of LMs in each category and its percentage to the total number of LMs for each of the software systems.

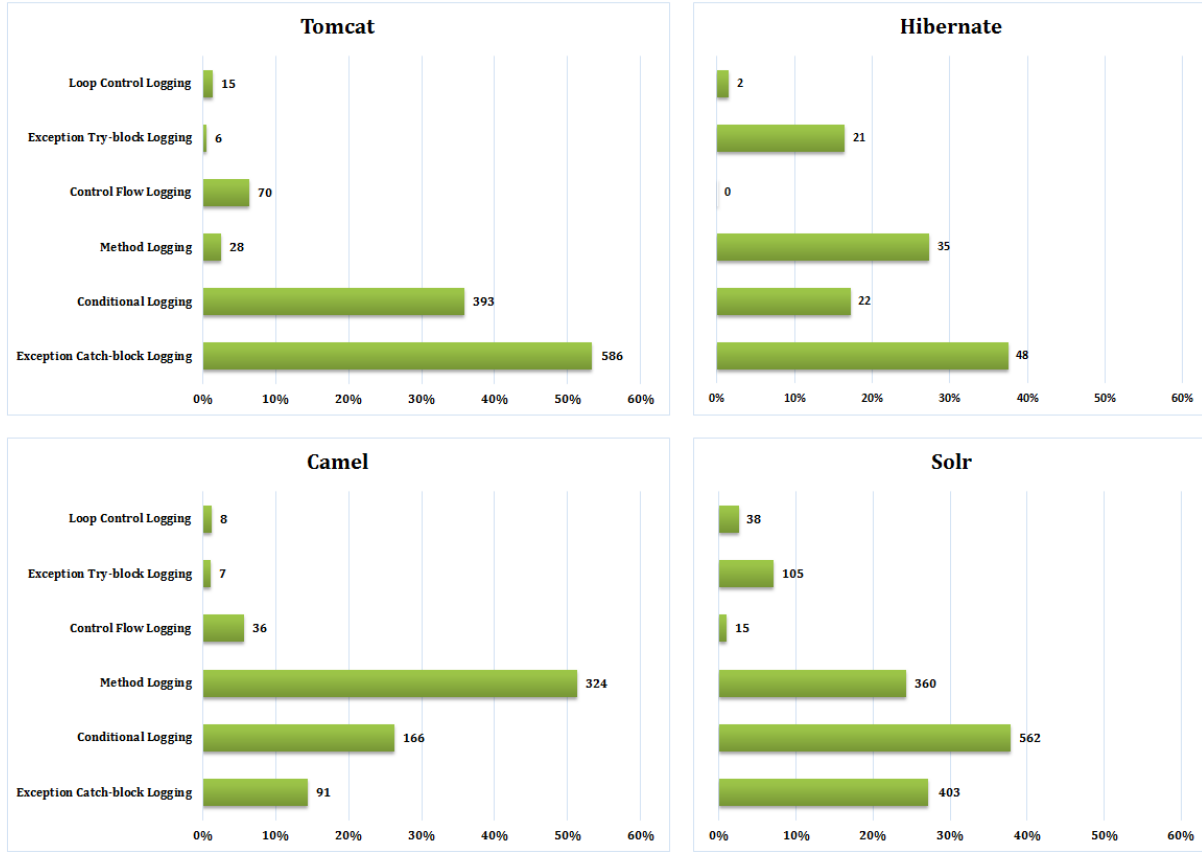


Figure 4.3: The distribution of the categories of anti-unifiers in the logging usage.

#### A. *Exception Catch-block Logging*

The main common structural characteristics of the anti-unifiers of this category are the **try** and **catch** statements, where the log statements are located inside the body of a Catch Clause. As shown in Figure 4.3, 14% to 53% of the total LMs are described by the anti-unifiers of this category, and it is the most commonly used logging usage category in the Tomcat and Hibernate software systems. The popularity of this category among all the studied systems is due to the fact that exception handling using the **try/catch** blocks is a common error handling technique in the Java programming language.

### *B. Conditional Logging*

In this category, log statements are enclosed by **if** statements that their expressions are mostly among Infix Expression, Method Invocation, or Binary Expression nodes. The Infix Expressions mostly check if the value of a variable either equals to the Null Literal or is outside of a valid range; the **if** statements with Method Invocations mostly check if the return value of an invoked method is an indicator of a potential problem within the system; and the **if** statements with Binary Expressions mostly checks if a Boolean Literal is incorrect. As shown in Figure 4.3, 17% to 36% of the total LMs are described by the anti-unifiers of this category, and it is the most commonly used logging category in the Solr system.

### *C. Method Logging*

In this category, the log statements are located inside the body of Method Declaration nodes. A common structural characteristic of the anti-unifiers in this category is that they mostly use the Throw Statements to throw an exception if an error occurs. The percentage of LMs that are described by the anti-unifiers of this category ranges from 3% to 51%, and it is the most common logging usage category in the Apache Camel. This suggests that developers use logging to record important method granularity information about the state of a software system. This information might be used later to detect the root causes of an application problem.

### *D. Exception Try-block Logging*

In this category, the log statement is located inside the body of the **try** statement of a **try/catch** block. These log statements can be used to record important information about the code that may throw an exception. According to the Figure 4.3, 1% to 16% of the total LMs of the studied systems are described by the anti-unifiers of this category.

### *E. Control Flow Logging*

In this category, the log statements are located inside the body of either **switch**– or **if**–then–**else**– statements. These log statements can be used to reveal necessary information to track the

location of root causes of a potential problem in a software system. According to the Figure 4.3, 0% to 6% of the total LMs are described by the anti-unifiers of this category.

#### *F. Loop Control Logging*

The log statements of this category are located inside the **for**–, enhanced–**for**–, **while**–, or **do**–**while**– statements. These log statements are mostly used to log important information about every object of a Java Collection. This information might be helpful to diagnose the root causes of a failure within a system. According to the Figure 4.3, a low percentage of LMs (1% to 3%) are described by the anti-unifiers of this category. This suggest that in practice, it is not a common way of using log statements in source code.

## 4.2 Evaluation

An empirical study is conducted to evaluate the quality of the anti-unifiers generated by ELUS in describing the location of log statements in source code. Section 4.2 describe the process of evaluating the precision and recall of ELUS.

#### *Calculating the precision and recall*

To find the locations in source code that are described by an anti-unifier using ELUS, I applied the DETERMINE-LOCATIONS algorithm, which takes the anti-unifier and a list of all methods in source code and outputs a list of methods that their AUAST matches the anti-unifier AUAST. This algorithm anti-unifies each method in the list with the anti-unifier using the ANTIUNIFY algorithm described in Section (Lines 2–3). If the result equals the anti-unifier, that method will be added to the list of locations matching the anti-unifier (Lines 4–5). EQUALS is a procedure that takes two AUAST nodes and checks whether they are equal or not. To evaluate the generalizability of the anti-unifiers, I have implemented this procedure in two ways: (1) when variables are considered to be *constrained*, it tests that the non-variable nodes are identical in the two AUASTs and checks if the constraints of variable are identical or not. (2) When variables are considered to be



*unconstrained*, it tests that the non-variable nodes are identical in the two AUASTs, but permits unconstrained variables to differ. I ran my tool on the source code of the five studied systems and applied this algorithm to find the locations in the code that matches the structure of the generated anti-unifiers. Then, the precision and recall metrics are calculated using the equations 4.2 and 4.3, respectively.

---

**Algorithm 4.1** DETERMINE-LOCATIONS(*antiUnifier*, *methods*) finds the locations in source code that matches an anti-unifier.

---

**DETERMINE-LOCATIONS**(*antiUnifier*, *methods*)

```

1: locations  $\leftarrow$  ()
2: for method  $\in$  methods do
3:   result  $\leftarrow$  ANTIUNIFY(antiUnifier, method)
4:   if EQUALS(result, antiUnifier) then
5:     APPEND(method, locations)
6:   end if
7: end for
8: return locations

```

---

$$precision = \frac{TP}{TP + FP} \quad (4.2)$$

$$recall = \frac{TP}{TP + FN} \quad (4.3)$$

Where TP is the number of correct locations obtained, FP is the number of incorrect locations retrieved, and FN is the number of correct locations that are not retrieved. Figures 4.4 and 4.5 show the precision and recall results for each software system where the experiment was run once with constrained variables and once with unconstrained variables.

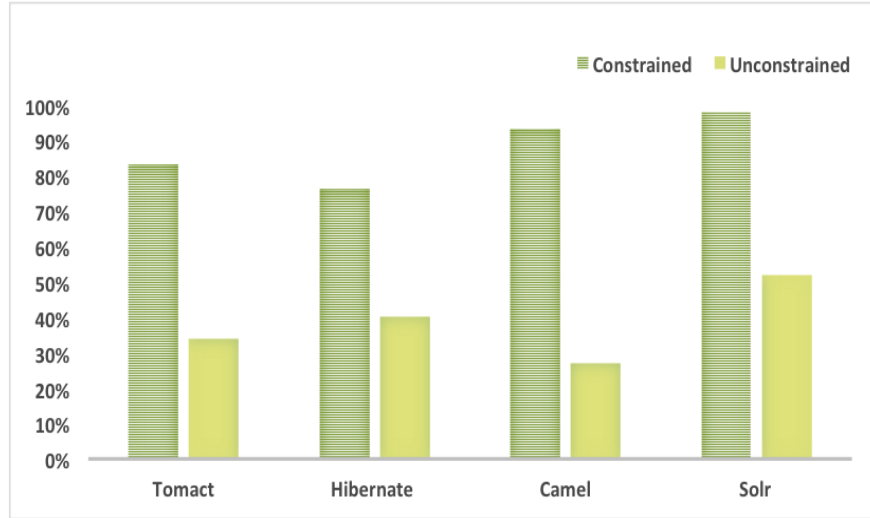


Figure 4.4: The precision of ELUS.

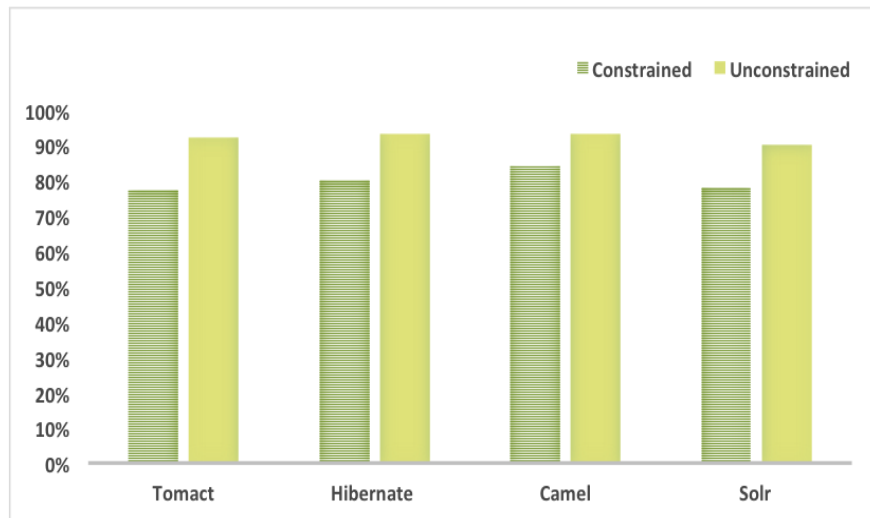


Figure 4.5: The recall of ELUS.

### *Precision Results*

The green and yellow bars in Figure 4.4 show the precision results when the experiment was run with constrained and unconstrained variables, respectively. I have also calculated the overall average precision of ELUS, by averaging the precision values between the five software systems. The average precision for ELUS is 88% and 38% for constrained and unconstrained variable experiments, respectively. In general, the precision for constrained variables are fairly high. The main

reason behind the high precision is that in these cases the variables can only be substituted with some particular nodes, which makes the anti-unifier very specific. Furthermore, according to the Figure 4.4, the precision is fairly low for unconstrained variables. The main reason for the low precision for these cases is the fact that the unconstrained variables can be substituted by any nodes, which makes the anti-unifiers too general. As a result, the tool finds many incorrect locations that matches the anti-unifiers.

### *Recall Results*

The green and yellow bars in Figure 4.5 show the recall results when the experiment was run with constrained and unconstrained variables, respectively. I have also calculated the overall average recall of ELUS, by averaging the recall values between the studied systems. The average recall for ELUS is 80% and 92% for the constrained and unconstrained variable experiments, respectively. In general, when variables are constrained, ELUS can detect many correct locations, as the recalls for all the studied systems are fairly high. Also, ELUS can detect most of the correct locations in source code when no constraints are taken on variable nodes.

The main reason behind ELUS's failure in detecting the correct locations is the potential complexities in constructing anti-unifiers from a large set of source code fragments. As in some cases, the anti-unifier might not maintain the correct locations of nodes in the AST hierarchy, and thus ELUS would not be able to detect correct locations of log statements in the source code.

## 4.3 Summary

I conducted an experimental study to characterize the location of log statements by applying my tool on the source code of five full software systems that make use of the Apache Log4j logging framework. My tool inputs the source code of these systems, extracts ASTs of LMs, applies the proposed anti-unification and clustering algorithms, and outputs the anti-unifier for each cluster. I also conducted an experimental study to evaluate the precision and the recall of ELUS in constructing the anti-unifiers that describe the location of log statements in source code. This experiment

shows that ELUS has achieved promising results in terms of precision and recall. Furthermore, The results taken from the characterization experiment shows that there are common ways of locating log statements I manually examined the detailed view of structural generalizations to categorize the anti-unifiers of logging usage.

# Chapter 5

## Discussion

In this chapter, I discuss the validity of my evaluation and the characterization study (Section 5.1), and a number of remaining issues, including the limitations and pitfalls of my approach and the tool support (Section 5.2), and the usage of anti-unification theory for other applications (Section 5.3).

### 5.1 Threats to validity

Prior to applying our tool for characterizing logging usage in real-world software systems, I have conducted three experiments to investigate the effectiveness of the proposed approach. However, there are several potential threats regarding the validity of these experiments. First, the results of my manual examination might be biased, as I determined the correct correspondences between AUASTs and the correct way of classifying the set of AUASTs in our test suite based on a similarity measurement. To limit the bias, other people can be involved to double check the accuracy of my manual inspection in a future work. Secondly, the experiments have examined one test suite containing a set of LMs from a real-world software system, though different test suites may generate different results. Although I cannot claim that the LMs in my test suite are a good representative of all LMs in real-world software systems, the results are still promising, as logging calls are used in various ways in Java methods of my test suite, and have sufficed to indicate the effectiveness of my approach in constructing structural generalizations. Another potential thread is that the successful rate of detecting correspondences by our tool might happened accidentally only for our test suite. To resolve this doubt, I examined the cases where our tool fails to detect correct correspondences, and I found that the failures are due to the fundamental limitations and complexities in the construction of structural generalization through the use of structural correspondence. That is, our tool creates structural generalizations successfully with regard to what

our algorithm should generate. A potential thread to the validity of our characterization study is the degree to which our sample set of software systems is a good representation of all real-world logging usage. To address this issue, we selected various open-source software projects in terms of application, including a programming text editor, a web server, and an application server. These software systems are among the most popular applications in their own product category, and they all have at least 10 years of history in software development. However, our findings might not be able to reflect the characteristics of logging usage in other types of systems such as commercial software systems, or software written in other programming languages.

## 5.2 The pitfalls of my tool

There are some issues that the approximation approach and my tool support is not able to handle perfectly, including inaccurate node ordering, and the resolution of conflicts happened in constructing the anti-unifiers.

### 5.2.1 Inaccurate node ordering

Our anti-unification algorithm does not guarantee to maintain the correct sequence of statements in the body of methods when anti-unifying two method declaration nodes, since the order of statement nodes is not considered in determining the best correspondences. For example, consider we have two corresponding methods  $method_1$  and  $method_2$  embodying  $a_1, a_2, a_3$  and  $b_1, b_2$  sequences of statements, respectively. If our tool finds that the  $b_1$  and  $b_2$  nodes are the best correspondences for the  $a_3$  and  $a_1$  nodes respectively, the output generalization view for the set of statement nodes would be  $a_1$ -or- $b_2$ ,  $a_2$ -or-NIL,  $a_3$ -or- $b_1$ . Therefore, the generalization view does not preserve the correct ordering of nodes in the original structures.

### 5.2.2 Conflict resolution

The decisions I have made to resolve the conflicts occurred in constructing structural generalizations might affect the accuracy of our results. For example, in situations where I have two correspondences with the same similarity value in the ordered list of correspondence connections, my approach picks the one which involves two subtrees with higher number of leaves, though it might be not the best choice for all cases. In addition, I consider AST hierarchies to perform anti-unification. That is, my algorithm does not anti-unify two nodes if their parent nodes are not found to be corresponded. As a result, situations can occur where in fact two nodes should be anti-unified with each other, while they are not anti-unified by the tool. Though these decisions led me to get approximate results, they helped to limit the complexity of my approach allowing the implementation of it as a practical solution.

## 5.3 Applications of anti-unification

Our study demonstrates the application of an extended form of anti-unification (HOAUMT) to infer usage patterns of log statements in source code via the creation of structural generalizations. Anti-unification and its extensions have been already applied to solve several theoretical and practical problems, such as analogy making [Schmidt, 2010], determining lemma generation in equational inductive proofs [Burghardt, 2005], and detecting the construction laws for a sequence of structures [Burghardt, 2005].

Higher-order anti-unification modulo theories can be used to create generalizations in different contexts, and therefore the set of equational theories should be developed particularly for the higher-order structure used in each problem context. That is, the utility of these theories are highly dependent on how well they allow the incorporation of semantic knowledge of structures. In addition, these theories should ensure that only a finite number of anti-instances exist for each structure. Taking all these considerations into account enables HOAUMT to anti-unify sets of structures in a particular context. The practical tests I have conducted through the application of my tool on a test

suite demonstrate that our approximation of HOAUMT was successful in constructing structural generalizations required to solve our problem.

## 5.4 Summary

I discussed the potential threads to validity of my evaluation and characterization study. To limit the bias of the experiments I conducted to evaluate the effectiveness of my approach and the tool support, I selected the test cases from a real system with various levels of similarity in the usage of logging calls. Furthermore, I examined the failed test cases to assure that our tool works when it should work with regard to the proposed algorithm. I will also make our test suite available for public examination to further check the accuracy of our manual inspection. For the characterization study, I selected various software systems in terms of functionality that are widely used by many developers for a long period of time. I also discussed the remaining issues with the tool support, including inaccurate node ordering and handling the conflicts happened in the construction of anti-unifiers.

This work aims to provide a detailed view of structural generalizations constructed from a set of source code fragments that use log statements via the application of an approximated anti-unification and clustering. However, I argued how higher-order anti-unification modulo theories can be effectively approximated for various applications by means of developing an appropriate set of equational theories particularly for the higher-order structure used in each problem context.



# Chapter 6

## Related Work

In this chapter, we review related work to the topics of our study including: the application of logging in real-world software systems (Section 6.1), determining correspondences in source code (Section 6.2), data mining approaches to extract API usage patterns (Section 6.3), anti-unification and its application to detect structural correspondences and construct generalizations (Section 6.4), and clustering (Section 6.5).

### 6.1 Usage of logging

Logging is a conventional programming practice to record a software system’s runtime information that can be used in post-modern analysis to trace the root causes of systems’ activities. Log analysis is most often performed for failure diagnosis, system behavioral understanding, system security monitoring, and performance diagnostics purposes as described below:

- **Log analysis for failure diagnosis:** Xu et al. [2009] use statistical techniques to learn a decision tree based signature from console logs and then utilize the signature to diagnose anomalies. SherLog [Yuan et al., 2010] uses failure log messages to infer the source code paths that might have been executed during a failure.
- **Log analysis for system behaviour understanding:** Fu et al. [2013] present an approach for understanding system behaviour through contextual analysis of logs. They first extracted execution patterns reflected by a sequence of system logs and then utilized the patterns to find contextual factors from logs that causes a specific system behavior. The Linux Trace Toolkit [Yaghmour and Dagenais, 2000] was created to record and analyze system behavior by providing an efficient kernel-level event logging infrastructure. A more flexi-

ble approach is taken by DTrace [Cantrill et al., 2004] which allows dynamic modification of kernel code.

- **Log analysis for system security monitoring:** Bishop [1989] proposes a formal model of system's security monitoring using logging and auditing. Peisert et al. [2007] have developed a model that demonstrates a mechanism for extracting logging information to detect how an intrusion occurs in software systems.
- **Log analysis for performance diagnosis:** Nagaraj et al. [2012] developed an automated tool to assist developers in diagnosis and correction of performance issues in distributed systems by analyzing system behaviours extracted from the log data.

Jiang et al. [2009] study the effectiveness of logging in problem diagnosis. Their study shows that customer problems in software systems with logging resolve faster than those without logging by investigating the correlations between failure root causes and diagnosis time. Despite the importance of logging for software development and maintenance, few studies have been conducted in pursuit of understanding logging usage in real-world software. Yuan et al. [2012b] provides a quantitative characteristic study to investigate log message modifications on four open-source software systems by mining their revision history. Their study shows that developers spend great effort modifying log statements as after-thoughts, which indicates that they are not satisfied with the log quality in their first attempt. They also characterize where developers spend most of their time in modifying the log messages.

Yuan et al. [2012a] study the problem of lack of log messages for error diagnosis and suggests to log when generic error conditions happens. LogEnhancer [Yuan et al., 2012c] automatically enhances existing log messages by detecting variables containing important values and inserting them into the log messages. However, these studies only consider source code fragments containing bugs that are needed to be logged and do not consider the other code fragments with no bugs but still needed to be logged. Moreover, these studies mainly research log message modifications

and potential enhancements of them, however, the focus of this study is on understanding where logging calls are used in source code.

## 6.2 Correspondence

Several studies have been conducted to find the similarities and differences between source code fragments. Baxter et al. [1998] develop an algorithm to detect code clones in source code that uses hash functions to partition subtrees of ASTs of a program and then find common subtrees in the same partition through a tree comparison algorithm. Apiwattanapong et al. [2004] present a top-down approach to detect differences and correspondences between two versions of a Java program, through comparison of the control flow graphs created from source code. Holmes et al. [2005] recommend relevant code snippet examples from a source code repository for the sake of helping developers to find examples of how to use an API by heuristically matching the structure of the code under development with the code in the repository. Coogle [Sager et al., 2006] was developed to detect similar Java classes by converting ASTs to a normalized format and then comparing them through tree similarity algorithms. However, none of these approaches construct a detailed view of structural generalizations needed in our context.

Cossette et al. [2014] present a new approach, called matching via structural generalization (MSG), to recommend replacements for API migration. They used Jigsaw to find structural correspondences, however, the proposed algorithm does not suffice to construct structural generalizations that represent the detailed commonalities and differences of a set of source code fragments with special attention to log statements, which is required to solve our problem.

## 6.3 API usages patterns

Various data mining approaches have been used to extract API usages patterns out of source code such as unordered pattern mining and sequential pattern mining [Robillard et al., 2013]. Unordered pattern mining, such as association rule mining and itemset mining, extracts a set of API usage

rules without considering their order [Agrawal et al., 1994]. CodeWeb [Michail, 2000] uses data mining association rules to identify reuse patterns between a source code under development and a specific library. PR-Miner [Li and Zhou, 2005] uses frequent itemset mining to extract implicit programming rules from source code and detect violations. The sequential pattern mining technique is different from the unordered one in the way that it considers the order of API usage. As an example, MAPO [Xie and Pei, 2006] combines frequent subsequence mining with clustering to extract API usage patterns from source code.

Another technique for extracting API usage patterns is through statistical source code analysis. For example, PopCon [Holmes and Walker, 2008] is a tool developed to help developers understanding how to use APIs in their source code through calculating popularity statistics for each API of a library. Acharya et al. [2007] present a framework to extract API usage scenarios as partial orders, as specifications were extracted from frequent partial orders. They adapted a compile time model checker to generate control-flow-sensitive static traces of APIs, from which API usage scenarios were extracted. However, none of these approaches suffice to construct detailed structural generalizations needed in our context.

## 6.4 Anti-unification

Anti-unification is the problem of finding the most specific generalization of two terms. First-order syntactical anti-unification was introduced by Plotkin [1970] and Reynolds [1970], independently. Burghardt and Heinz [1996] extend the notion of anti-unification to E-anti-unification to incorporate background knowledge to syntactical anti-unification, which is required for some applications. Anti-unification and its extensions have been applied in various studies for program analysis. Bulychiev and Minea [2009] suggest an anti-unification algorithm to detect clones in ASTs. Their approach consists of three stages: first, identifying similar statements through anti-unification and grouping them into clusters; second, determining similar sequences of statements with the same cluster identifier; third, refining candidate statement sequences using an anti-unification based sim-

ilarity measurement to generate final clones. However, their approach does not construct structural generalizations.

Cottrell et al. [2007] propose Breakaway to automatically determine structural correspondences between a pair of ASTs to create a generalized correspondence view. However, their approach does not allow the determination of the best structural correspondence for each AST node required to our context. Cottrell et al. [2008] developed Jigsaw to help developers integrate small-scale reused code into their own source code by determining structural correspondences through the application of higher-order anti-unification modulo theories. Although I used the Jigsaw framework to find potential correspondences between AST nodes, their approach does not suffice to construct structural generalizations of a set of source code fragments by considering the limitations of this study in determining correspondences.

## 6.5 Clustering

Clustering is an unsupervised machine mining technique that aims to organize a collection of data into clusters, such that intra-cluster similarity is maximized and the inter-cluster similarity is minimized [Karypis et al., 1999, Grira et al., 2004]. We divided existing clustering approaches into two major categories: partitional clustering and hierarchical clustering. Partitional clustering try to classify a data set into  $k$  clusters such that the partition optimizes a pre-determined criterion [Karypis et al., 1999]. The most popular partitional clustering algorithm is  $k$ -means, which repeatedly assigns each data point to a cluster with the nearest centroid and computes the new cluster centroids accordingly until a pre-determined number of clusters is obtained [Bouguettaya et al., 2015]. However, the  $k$ -means clustering algorithm is not a good fit to our problem, as it requires to predefine the number of clusters we need to come up with, which is not reasonable in our context.

Hierarchical clustering algorithms produce a nested grouping of clusters, with single point clusters at the bottom and an all-inclusive cluster at the top [Karypis et al., 1999]. Agglomerative hierarchical clustering is one of the main stream clustering methods [Day and Edelsbrunner, 1984]

and has applications in document retrieval [Voorhees, 1986] and information retrieval from a search engine query log [Beeferman and Berger, 2000]. It starts with singleton clusters, where each contains one data point. Then it repeatedly merges the two most similar clusters to form a bigger one until a pre-determined number of clusters is obtained or the similarity between the closest clusters becomes below a pre-determined threshold value. Hierarchical clustering algorithms work implicitly or explicitly with the  $n \times n$  similarity matrix such that an element in row  $i$  and column  $j$  represents the similarity between the  $i$ th and the  $j$ th clusters [Karypis et al., 1999].

There are various versions of agglomerative hierarchical algorithms that mainly differ in how they update the similarity between clusters. There are various methods to measure the similarity between clusters, such as single linkage, complete linkage, average linkage, and centroids [Rasmussen, 1992] **[RW: Put into bibtex]**. In the single linkage method, the similarity is measured by the similarity of the closest pair of data points of two clusters. In the complete linkage method, the similarity is computed by the similarity of the farthest pair of data points of two clusters. In the average linkage method, the similarity is measured by the average of all pairwise similarities between data points of two clusters. In the centroids methods, each cluster is represented by a centroid of all data points in the cluster, and the similarity between two clusters is measured by the similarity of the clusters' centroids. However, in our application, each cluster is composed of one AUAST, and the similarity between two clusters is measured by the similarity between the clusters' AUASTs, which is the ratio of the number of common pieces over the total number of pieces of their anti-unifier.

## 6.6 Summary

Despite the great importance of logging and its various applications in software development and maintenance, few studies have focused on understanding logging usage in source code. Some work has been done on characterizing log messages modifications made by developers and to help them enhance the content of log messages. However, to my knowledge, no study has been

conducted on characterizing where logging is used in source code via structural generalization and clustering. Several data mining and statistical source code analysis techniques have been used to extract API usage patterns, however, none of them enable us to construct the detailed structural generalizations of a set of source code fragments. On the other hand, using higher-order anti-unification modulo theories and an agglomerative hierarchical clustering algorithm allow us to construct generalizations representing the commonalities and differences between ASTs of logged Java methods and grouping them into clusters based on structural correspondence.

# Chapter 7

## Conclusion

Determining the detailed structural similarities and differences between a set of source code fragments is a complex task that can be applied to solve several source code analysis problems. As a specific application, the focus of this study is on detecting usage patterns of logging calls in source code via structural generalization and clustering.

Logging is a pervasive practice and has various applications in software development and maintenance. However, it is a challenging task for developers to understand how to use logging calls in source code. I have presented an approach to automatically characterize where logging calls happen in source code. I have developed a prototype tool implementing my proposed approach that proceeds in three steps. First, it extracts the ASTs of logged Java methods using the Eclipse JDT framework, extends the AST structures to AUAST, and determines potential structural correspondences between AUAST nodes via the Jigsaw framework. Second, it constructs an anti-unifier from AUASTs of two given LMs with a focus on logging calls through the implementation of higher-order anti-unification modulo theories. Due to the problem of undecidability of HOAUMT, it employs an approximation technique which greedily determines the best correspondence for each node with the highest similarity. It applies several constraints prior to determining the best correspondences to prevent the anti-unification of logging calls with anything else. It also develops a measure of structural similarity that determines how similar is the usage of logging calls in these Java methods. Third, it classifies a set of logged Java methods via a hierarchical clustering algorithm suited to our application.

The application of HOAUMT to construct generalizations via determining structural correspondences is novel to the problem of extracting usage patterns of log statements in source code. To evaluate the effectiveness of this approach in constructing generalizations and clustering logged



Java methods, three experiments were conducted on a sample test suite. I found that my tool was successful in determining correct correspondences for my application in % of test cases. It was also successful in clustering logged Java methods of our test suite. This work also shows how the Jigsaw framework could be effectively used to construct structural generalizations for a particular problem context by determining structural correspondences.

Furthermore, an study was conducted to infer logging usage patterns of the source code of three software systems on a method-granularity basis via my tool that generates a detailed view of structural generalizations describing the commonalities and differences between the usage of logging calls in Java methods. Our characterization study shows

In summary, our study makes the following contributions:

- An approach to constructing a structural generalization from ASTs of two logged Java methods with special attention to log statements by determining structural correspondences and developing an approximation of higher-order anti-unification modulo theories for our context.
- An approach to developing a similarity measure that indicates the level of similarity between the usage of logging calls in two Java methods.
- A hierarchical clustering algorithm to clustering a set of Java methods showing different usages of logging calls into different clusters.
- An approach to detecting usage patterns of log statements in source code via structural generalization and clustering.

## 7.1 Future Work

Future work could be directed to address the remaining issues of this study as described below:

- **Improving logging practices:** characterizing logging usage could be a huge step towards improving logging practices through the provision of some guidelines that might help de-

developers in making decisions about where to log. Further studies could be conducted to investigate the feasibility of predicting the location of logging calls based on the detected usage patterns. Future work can also be done to develop recommendation tool supports that not only save developers time and effort for making decisions about where to log, but also improve the quality of logging practices.

- **Further validation of this study:** The characterization study can be conducted on more software systems to further validate the findings of my study. In addition, a survey can be conducted to ask developers on the factors they consider to decide on where to log. It might also be helpful to recognize important structural and semantic information that should be taken into account for characterizing logging usage.
- **Further extensions to my approach and the tool support:** data flow analysis can be performed to detect the problems related to node ordering in the construction of anti-unifiers. This approach can also be extended to examine more advanced semantical and contextual information of source code fragments enclosing logging calls in addition to structural information. In addition, further analyses can be done to detect and resolve all the conflicts happen in deciding the best correspondences to construct an approximation of the best anti-unifier for our problem. However, the complexity of applying all these extensions must be kept restricted to maintain the approach as a practical one.
- **Other applications:** any applications that are involved in the inference of structural patterns in source code even infrequently-used patterns might benefit from my tools underlying framework. Furthermore, understanding the commonalities and differences among source code fragments has application in several areas of software engineering, such as API usage pattern collation, code clone detection, recommending replacements for API migration, and merging different branches of version control systems. Our tool's functionality to construct a detailed view of structural generalizations of a set of source code fragments via structural correspondence and clustering could be used to improve the results of these studies as well.

## Bibliography

- Mithun Acharya, Tao Xie, Jian Pei, and Jun Xu. Mining api patterns as partial orders from source code: from usage scenarios to specifications. In *Proceedings of the the 6th joint meeting of the European software engineering conference and the ACM SIGSOFT symposium on The foundations of software engineering*, pages 25–34. ACM, 2007.
- Rakesh Agrawal, Ramakrishnan Srikant, et al. Fast algorithms for mining association rules. In *Proc. 20th int. conf. very large data bases, VLDB*, volume 1215, pages 487–499, 1994.
- Taweessup Apiwattanapong, Alessandro Orso, and Mary Jean Harrold. A differencing algorithm for object-oriented programs. In *Proceedings of the 19th IEEE international conference on Automated software engineering*, pages 2–13. IEEE Computer Society, 2004.
- Ira D Baxter, Andrew Yahin, Leonardo Moura, Marcelo Sant’Anna, and Lorraine Bier. Clone detection using abstract syntax trees. In *Software Maintenance, 1998. Proceedings., International Conference on*, pages 368–377. IEEE, 1998.
- Doug Beeferman and Adam Berger. Agglomerative clustering of a search engine query log. In *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 407–416. ACM, 2000.
- Matt Bishop. A model of security monitoring. In *Computer Security Applications Conference, 1989., Fifth Annual*, pages 46–52. IEEE, 1989.
- Athman Bouguettaya, Qi Yu, Xumin Liu, Xiangmin Zhou, and Andy Song. Efficient agglomerative hierarchical clustering. *Expert Systems with Applications*, 42(5):2785–2797, 2015.
- Peter Bulychiev and Marius Minea. An evaluation of duplicate code detection using anti-unification. In *Proc. 3rd International Workshop on Software Clones*. Citeseer, 2009.

- J. Burghardt. E-generalization using grammars. *Artificial Intelligence Journal*, 165(1):1–35, June 2005. doi: 10.1016/j.artint.2005.01.008.
- Jochen Burghardt and Birgit Heinz. Implementing anti-unification modulo equational theory. arbeitspapier 1006, 1996.
- Bryan Cantrill, Michael W Shapiro, Adam H Leventhal, et al. Dynamic instrumentation of production systems. In *USENIX Annual Technical Conference, General Track*, pages 15–28, 2004.
- Siobhán Clarke, William Harrison, Harold Ossher, and Peri Tarr. The dimension of separating requirements concerns for the duration of the development lifecycle. In *First Workshop on Multi-Dimensional Separation of Concerns in Object-oriented Systems (at OOPSLA)*, 1999a.
- Siobhán Clarke, William Harrison, Harold Ossher, and Peri Tarr. Subject-oriented design: towards improved alignment of requirements, design, and code. *ACM SIGPLAN Notices*, 34(10):325–339, 1999b.
- Bradley Cossette, Robert Walker, and Rylan Cottrell. Using structural generalization to discover replacement functionality for API evolution. Technical Report 2014-745-10, Department of Computer Science, University of Calgary, Calgary, Canada, May 2014.
- Rylan Cottrell, Joseph J. C. Chang, Robert J. Walker, and Jörg Denzinger. Determining detailed structural correspondence for generalization tasks. In *Proceedings of the European Software Engineering Conference/ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 165–174, 2007. doi: 10.1145/1287624.1287649.
- Rylan Cottrell, Robert J. Walker, and Jörg Denzinger. Semi-automating small-scale source code reuse via structural correspondence. In *Proceedings of the ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 214–225, 2008. doi: 10.1145/1453101.1453130.

- William HE Day and Herbert Edelsbrunner. Efficient algorithms for agglomerative hierarchical clustering methods. *Journal of classification*, 1(1):7–24, 1984.
- Qiang Fu, Jian-Guang Lou, Yi Wang, and Jiang Li. Execution anomaly detection in distributed systems through unstructured log analysis. In *ICDM*, volume 9, pages 149–158, 2009.
- Qiang Fu, Jian-Guang Lou, Qingwei Lin, Rui Ding, Dongmei Zhang, and Tao Xie. Contextual analysis of program logs for understanding system behaviors. In *Proceedings of the 10th Working Conference on Mining Software Repositories*, pages 397–400. IEEE Press, 2013.
- James Gosling, Bill Joy, Guy Steele, Gilad Bracha, and Alex Buckley. *The Java Language Specification*. Addison-Wesley, Java SE 7 edition, 2012. URL <http://docs.oracle.com/javase/specs/jls/se7/html/index.html>.
- Nizar Grira, Michel Crucianu, and Nozha Boujemaa. Unsupervised and semi-supervised clustering: a brief survey. *A review of machine learning techniques for processing multimedia content*, 1:9–16, 2004.
- Samudra Gupta. Pro apache log4j: Java application logging using the open source apache log4j api. *Apress®*, USA, 2005.
- Reid Holmes and Robert J Walker. A newbie’s guide to eclipse apis. In *Proceedings of the 2008 international working conference on Mining software repositories*, pages 149–152. ACM, 2008.
- Reid Holmes, Robert J Walker, and Gail C Murphy. Strathcona example recommendation tool. In *ACM SIGSOFT Software Engineering Notes*, volume 30, pages 237–240. ACM, 2005.
- Weihang Jiang, Chongfeng Hu, Shankar Pasupathy, Arkady Kanevsky, Zhenmin Li, and Yuanyuan Zhou. Understanding customer problem troubleshooting from storage system logs. In *FAST*, volume 9, pages 43–56, 2009.
- George Karypis, Eui-Hong Han, and Vipin Kumar. Chameleon: Hierarchical clustering using dynamic modeling. *Computer*, 32(8):68–75, 1999.

- Zhenmin Li and Yuanyuan Zhou. PR-Miner: Automatically extracting implicit programming rules and detecting violations in large software code. In *SIGSOFT Software Engineering Notes*, volume 30, pages 306–315, 2005.
- Jian-Guang Lou, Qiang Fu, Shengqi Yang, Ye Xu, and Jiang Li. Mining invariants from console logs for system problem detection. In *USENIX Annual Technical Conference*, 2010.
- Amir Michail. Data mining library reuse patterns using generalized association rules. In *Proceedings of the ACM/IEEE International Conference on Software Engineering*, pages 167–176, 2000. doi: 10.1109/ICSE.2000.870408.
- Karthik Nagaraj, Charles Killian, and Jennifer Neville. Structured comparative analysis of systems logs to diagnose performance problems. In *Presented as part of the 9th USENIX Symposium on Networked Systems Design and Implementation (NSDI 12)*, pages 353–366, 2012.
- Sean Peisert, Matt Bishop, Sidney Karin, and Keith Marzullo. Toward models for forensic analysis. In *Systematic Approaches to Digital Forensic Engineering, 2007. SADFE 2007. Second International Workshop on*, pages 3–15. IEEE, 2007.
- Gordon D Plotkin. A note on inductive generalization. *Machine intelligence*, 5(1):153–163, 1970.
- John C Reynolds. Transformational systems and the algebraic structure of atomic formulas. *Machine intelligence*, 5(1):135–151, 1970.
- Martin P Robillard, Eric Bodden, David Kawrykow, Mira Mezini, and Tristan Ratchford. Automated api property inference techniques. *Software Engineering, IEEE Transactions on*, 39(5): 613–637, 2013.
- Tobias Sager, Abraham Bernstein, Martin Pinzger, and Christoph Kiefer. Detecting similar java classes using tree algorithms. In *Proceedings of the 2006 international workshop on Mining software repositories*, pages 65–71. ACM, 2006.

- Martin Schmidt. Restricted higher-order anti-unification for heuristic-driven theory projection. Bachelor's thesis, University of Osnabrück, Osnabrück, Austria, 2010. URL <http://ikw.uni-osnabrueck.de/en/system/files/31-2010.pdf>.
- Ellen M Voorhees. Implementing agglomerative hierarchic clustering algorithms for use in document retrieval. *Information Processing & Management*, 22(6):465–476, 1986.
- Tao Xie and Jian Pei. MAPO: Mining API usages from open source repositories. In *Proceedings of the International Workshop on Mining Software Repositories*, pages 54–57, 2006. doi: 10.1145/1137983.1137997.
- Wei Xu, Ling Huang, Armando Fox, David Patterson, and Michael I Jordan. Detecting large-scale system problems by mining console logs. In *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles*, pages 117–132. ACM, 2009.
- Karim Yaghmour and Michel R Dagenais. Measuring and characterizing system behavior using kernel-level event logging. 2000.
- Ding Yuan, Haohui Mai, Weiwei Xiong, Lin Tan, Yuanyuan Zhou, and Shankar Pasupathy. Sherlock: error diagnosis by connecting clues from run-time logs. In *ACM SIGARCH computer architecture news*, volume 38, pages 143–154. ACM, 2010.
- Ding Yuan, Soyeon Park, Peng Huang, Yang Liu, Michael M Lee, Xiaoming Tang, Yuanyuan Zhou, and Stefan Savage. Be conservative: enhancing failure diagnosis with proactive logging. In *Presented as part of the 10th USENIX Symposium on Operating Systems Design and Implementation (OSDI 12)*, pages 293–306, 2012a.
- Ding Yuan, Soyeon Park, and Yuanyuan Zhou. Characterizing logging practices in open-source software. In *Proceedings of the 34th International Conference on Software Engineering*, pages 102–112. IEEE Press, 2012b.

Ding Yuan, Jing Zheng, Soyeon Park, Yuanyuan Zhou, and Stefan Savage. Improving software diagnosability via log enhancement. *ACM Transactions on Computer Systems (TOCS)*, 30(1):4, 2012c.

Jieming Zhu, Pinjia He, Qiang Fu, Hongyu Zhang, Michael R Lyu, and Dongmei Zhang. Learning to log: Helping developers make informed logging decisions. In *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, volume 1, pages 415–425. IEEE, 2015.