

Automatic Improvement of Software Maintainability through Automated Refactoring of Code Clones

Simon Baars

simon.mailadres@gmail.com

30 June 2019, 55 pages

Research supervisor: Dr. Ana Oprescu, ana.oprescu@uva.nl

Host/Daily supervisor: Xander Schrijen, x.schrijen@sig.eu

Host organisation/Research group: Software Improvement Group (SIG), <http://sig.eu/>



UNIVERSITY OF AMSTERDAM

FACULTY OF PHYSICS, MATHEMATICS AND INFORMATICS

MASTER SOFTWARE ENGINEERING

<http://www.software-engineering-amsterdam.nl>

Abstract

Duplication in source code can have a major negative impact on the maintainability of source code. There are several techniques that can be used in order to merge clones, reduce duplication, improve the design of the code and potentially also reduce the total volume of a software system. In this study, we look into the opportunities to aid in the process of refactoring these duplication problems for object-oriented programming languages. We focus primarily on the Java programming language, as refactoring in general is very language-specific.

We first look into redefinitions for different types of clones that have been used in code duplication research for many years. These redefinitions are aimed towards flagging only clones that are useful for refactoring purposes. Our definition defines additional rules for type 1 clones to make sure two cloned fragments are actually equal. We also redefined type 2 clones to reduce the number of false positives resulting from it.

We have conducted measurements that have indicated that more than half of the duplication in code is related to each other through inheritance, making it easier to refactor these clones in a clean way. Approximately a fifth of the duplication can be refactored through method extraction, the other clones require other techniques to be applied.

Contents

1	Introduction	4
1.1	Problem statement	4
1.1.1	The problem	4
1.1.2	Research questions	5
1.1.3	Research method	5
1.2	Contributions	6
1.2.1	Identification	6
1.2.2	Refactoring	6
1.3	Scope	6
1.4	Outline	6
2	Background	8
2.1	Code clone terminology	8
2.1.1	Example	8
2.2	Clone Types	9
2.2.1	Type 4 clones	10
2.3	Refactoring methods	10
2.3.1	Extract Method	10
2.3.2	Move method	10
2.3.3	Clone refactoring in relationship to its context	11
2.4	Internal and external classes	11
2.4.1	Maven	11
3	Defining refactoring-oriented clone types	12
3.1	Shortcomings of clone types	12
3.1.1	Type 1 clones	12
3.1.2	Type 2 clones	13
3.1.3	Type 3 clones	14
3.2	Refactoring-oriented clone types	14
3.2.1	Type 1R clones	14
3.2.2	Type 2R clones	15
3.2.3	Type 3R clones	17
3.2.4	Clone types summarized	19
3.2.5	The challenge of detecting these clones	19
3.3	Suitability of existing Clone Detection Tools for detecting these clones	20
4	CloneRefactor	21
4.1	JavaParser	21
4.2	Thresholds	22
4.3	Clone Detection	22
4.3.1	Generating the clone graph	23
4.3.2	Detecting Clones	25
4.3.3	Validating the type 2R variability threshold	27
4.3.4	Checking for type 3 opportunities	31
4.4	Context Analysis of Clones	33
4.4.1	Relation	33
4.4.2	Location	35

4.4.3	Contents	36
4.5	Method extraction opportunities	36
5	Results	37
5.1	The corpus	37
5.2	Clone types	37
5.3	Thresholds	38
5.4	Context Analysis of Clones	38
5.4.1	Relations Between Clone Instances	39
5.4.2	Clone instance location	39
5.5	Clone instance contents	40
5.6	Merging duplicate code through refactoring	40
5.7	Refactoring clones	40
6	Discussion	41
6.1	Clone Type Definitions	41
6.2	Clone Context Analysis	42
6.2.1	Refactorability	42
7	Related work	45
8	Conclusion	46
8.1	Threats to validity	46
8.2	Future work	46
8.2.1	Automated Refactoring for more metrics	47
8.2.2	Type 4 clones	48
8.2.3	Naming of refactored methods/classes/etc.	48
8.2.4	Looking into other languages/paradigms	49
8.2.5	Running our experiments with a different set of software projects	49
8.2.6	Automatically refactoring code that is duplicated with a library	49
	Bibliography	51
	Appendix A Non-crucial information	55

Chapter 1

Introduction

Refactoring is used to improve quality related attributes of a codebase (maintainability, performance, etc.) without changing the functionality. There are many methods that have been introduced to help with the process of refactoring [1, 2]. However, most of these methods still require manual assessment of where and when to apply them. Because of this, refactoring takes up a significant portion of the development process [3, 4], or does not happen at all [5]. For a large part, refactoring requires domain knowledge to do it right. However, there are also refactoring opportunities that are rather trivial and repetitive to execute. In this thesis, we take a look at the challenges and opportunities in automatically refactoring duplicated code, also known as “code clones”. The main goal is to improve maintainability of the refactored code.

Duplication in source code is often seen as one of the most harmful types of technical debt. *“Clones are problematic for the maintainability of a program, because if the clone is altered at one location to correct an erroneous behaviour, you cannot be sure that this correction is applied to all the cloned code as well. Additionally, the code base size increases unnecessarily and so increases the amount of code to be handled when conducting maintenance work”* [6]. Bruntink et al. [7] show that code clones can contribute up to 25% of the code size.

In this study, we use refactoring techniques to automatically reduce duplication in software systems. This allows us to obtain before- and after-refactoring snapshots of software systems. We use software maintainability metrics to measure the impact of refactoring clones. This way we can determine better definitions of clones. We also look into what variability we can allow between code fragments to still consider them clones, while still improving maintainability when refactoring these clones. Furthermore, we look into the thresholds that should be used while detecting clones to find clones that should be refactored.

We perform several quantitative experiments on a large corpus of open source software to collect statistical data. With these experiments we map the context of clones: where they reside in the codebase and what the relation is between duplicate fragments. We use the results to find appropriate refactoring opportunities for specific clones. We then automate the process of applying such refactorings, to be able to measure the impact on maintainability when refactoring clones found by certain definitions and thresholds.

1.1 Problem statement

In this section we describe the problem we address with this study and the research questions that we answer in order to contribute to solving the problem.

1.1.1 The problem

The maintainability of a codebase has a large impact on the time and effort spent on building the desired software system [8, 9]. The maintainability of software is one of the factors to be kept under control in order to avoid major delays and unexpected costs as a result of a software project [1]. One factor that has a major impact on the maintainability of a software system is the amount of duplicate code present in a codebase [10, 11].

The process of improving maintainability through the refactoring of duplicate code is time consuming and error-prone. This process mainly consists of these two aspects:

- Find refactoring candidates, either tool-assisted or manually.
- Refactor identified candidates, either tool-assisted or manually.

There are tools that assist in the process of finding duplicate code, but none of these tools identify all duplication problems. Often, these clone detection tools result in false-positives and false-negatives [12]. We define false-negatives as clones that can and should be refactored but are not found due to the configuration by which they are detected. We define false-positives as clones that do not improve the system maintainability when refactored. Many false-negatives are due to the thresholds that are used. Many clone detection tools [13, 14] only consider clones that have a minimum of 6 lines. Using such an approach is bound to result in many false-negatives, but using a lower threshold is shown to increase the number of false-positives.

A study by Batova et al. [15] shows that the process of refactoring often leaves side effects in the code. This study reports that refactorings involving hierarchies (e.g., pull up method), tend to induce faults very frequently. *This suggests more accurate code inspection or testing activities when such specific refactorings are performed.* Such specific refactorings often have to be used when dealing with duplication in software [1, 16]. Because of that, refactoring code clones has been empirically proven to cause bugs or other side effects in code.

1.1.2 Research questions

There is a lot of research on the topic of code clone detection. This research often results in tools that can detect clones by several clone type definitions. However, there is no research yet that looks into how these definitions align with refactoring opportunities. We will align clone type definitions as used in literature [12] with their corresponding refactoring methods [1]. For this, we answer the following research question:

Research Question 1:

How can we define clone types such that they **can** be automatically refactored?

As a result, we expect to formulate clone type definitions that can be refactored. On basis of this results we can perform analyses on the context of clones by these definitions. The context of a clone (location, relation between clones, etc.) has a big impact on how a clone should be refactored. We will create categories by which we map the context of clones and perform a statistical analysis on this. This results in the following research question:

Research Question 2:

How can we prioritize refactoring opportunities based on the **context** of clones?

We expect this research question to result in a prioritization of refactoring opportunities: *with what refactoring method can what percentage of clones be refactored?* As a result of these first two research questions we expect to have clone type definitions that can be refactored together with a prioritization of the refactoring methods that can be used. On basis of this we expect to be able to build a script that can automatically refactor the highest prioritized clones. With this script, we expect to be able to answer our final research question:

Research Question 3:

What are the discriminating factors to decide when a clone **should** be refactored?

Not in all cases will refactoring duplicated code result in a more maintainable codebase. Because of that, we compare the refactored code to the original code and measure the difference in maintainability. To do this, we will use a practical model consisting of metrics to measure maintainability [10]. Based on this, we can look into what *thresholds* result in better maintainable code when refactored. These thresholds consist of the size of clones and the variability we can allow between cloned fragments to still consider them clones.

1.1.3 Research method

We perform an **exploratory** study to look into the opportunities to automatically refactor code clones. To do this, we combine knowledge from literature with our own experience to develop definitions for refactorable clones. We also develop a tool to detect, analyze and refactor such clones. Using this tool, we perform **quantitative** experiments in which we statistically collect information about duplication

in open source software. In these experiments we control several variables to see their impact on the results. During this process we explore concepts and develop understanding, because of which decisions in the study design are based upon the results of the experiments.

1.2 Contributions

Many studies report that code clones negatively affect maintainability [monden2002software, juergens2009code, chatterji2013effects, 10]. However, no studies yet show in what cases code clones can reduce maintainability in source code. Refactoring often includes tradeoffs between design alternatives. With some code clones, the refactored alternative is worse maintainable than keeping the duplication [aversano2007clones, hotta2010duplicate, kim2005empirical, krinke2007study, saha2010evaluating, 17]. In this study, we analyze the maintainability of refactored code clones, in order to improve the suggestion of code clones that should be refactored. This assists with both the identification and refactoring of code clones.

1.2.1 Identification

There are many tools to detect clones. The goal of most of these tools is to assist developers in reducing duplication in their code, i.e. assisting in the refactoring process. The problem is that these tools have no limited insight on the impact of refactoring such clones on the design of the software. In this study, we can analyze a before- and after-refactoring snapshot of the code to determine the impact. A higher maintainability after refactoring increases the support for the clone being a true-positive. This way the results of our study can support the clone identification process.

1.2.2 Refactoring

The tool that results from this research can aid in the process of applying refactorings to clones. The tool will only apply a refactoring if the clone is refactorable and the maintainability of the source code increases as a result of applying the refactoring. The tool applies only refactorings that do not, in any way, influence the functional correctness of the program. Because of this, potential bugs as a result of refactoring can be avoided [15].

1.3 Scope

In this study, we perform research efforts to be able to detect refactorable clones. We will apply refactoring techniques to a subset of these clones and analyze the maintainability of the resulting source code.

There is a lot of study on how what metrics to consider to measure maintainability. In this study, we focus solely on the practical maintainability model by Heitlager et al. [10]. We consider the maintainability scores described in that paper as an sufficiently accurate indication of maintainability. This will be used to quantitatively determine the impact on maintainability when applying refactoring techniques to code clones.

Applying refactoring often entails creating new method declarations, class declarations, etc. Each of these declarations needs to have a name. Because the quality of the name is not included in the maintainability metrics we use, finding appropriate names for automatically refactored code fragments is out of the scope for this study. For our automated refactoring efforts, we will use generated names for these declarations.

It is very disputable whether unit tests apply to the same maintainability metrics that applies to the functional code. Because of that, for this research, unit tests are not taken into scope. The findings of this research may be applicable to those classes, but we will not argue the validity.

1.4 Outline

In Chapter 2 we describe the background of this thesis. In Chapter 3 we list shortcomings with clone definitions from literature and propose a set of clone type definitions which can be automatically refactored. In Chapter 4 we propose a tool to detect, analyze and automatically refactor such clones that

can be refactored. Using this tool we perform a set of experiments, of which the results are shown in Chapter 5 and discussed in Chapter 6. Chapter 7 contains the work related to this thesis. Finally, we present our concluding remarks in Chapter 8 together with future work.

Chapter 2

Background

This chapter will present the necessary background information for this thesis. Here, we define some basic terminology that will be used throughout this thesis.

2.1 Code clone terminology

Many studies present different definitions for code clone concepts. For this study, we mainly use the definitions from Bruntink et al. [7], Roy et al. [12] and Jiang et al. [18]. We created a summary of these concepts, which can be found in table 2.1.

Symbol	Meaning	Definition	Description	Properties
T	Token	-	Tokens are the basic lexical building blocks of source code. For this study, this is the smallest relevant entity of a program.	Range (R): The range this token spans. Category (C): Identifier, Keyword, Literal, Separator, Operator, Comment or Whitespace.
N	Node [18]	Set of tokens.	A statement or declaration in a codebase.	Range (R): The range this node spans.
I	Clone instance [7, 12]	Set of cloned nodes.	A code fragment that appears in multiple locations.	File (F): The file in which this clone instance is found.
C	Clone class [7, 12]	Set of clone instances.	A set of similar code fragments in different locations. Each of these code fragments is called a “clone instance”.	-
S	Clone class collection [7]	Set of clone classes.	All clone classes that have been found for a certain software project.	-

Table 2.1: Clone related terminology and how it maps to the source code.

A **range** denotes the line and column at which a code fragment starts and ends. In this study, we use the symbols given in this table in our formulas.

2.1.1 Example

This section gives an example to make the terminology of table 2.1 more concrete. The code fragment in figure 2.1 shows two clone classes. One of the clone classes consists of two clone instances and has three nodes per instance. The other clone instance has three clone instances and consists of two nodes per instance.

1	<code>// File1.java</code>	1	<code>// File2.java</code>	1	<code>// File3.java</code>
2	<code>doA();</code>	2	<code>doA();</code>	2	<code>doA();</code>
3	<code>doB();</code>	3	<code>doB();</code>	3	<code>doB();</code>
4	<code>doC();</code>	4	<code>doC();</code>	4	<code>doD();</code>

Figure 2.1: Two clone classes: One clone class with three clone instances and one with two clone instances.

This example can be represented as a tree structure using the concepts from table 2.1. This tree structure is displayed in figure 2.2.



Figure 2.2: Representation of figure 2.1 using the terminology introduced in table 2.1

2.2 Clone Types

Duplication in code is found in many different forms. Most often duplicated code is the result of a programmer reusing previously written code [19, 20]. Sometimes this code is then adapted to fit the new context. To reason about these modifications, several clone types have been proposed. These clone types are described in Roy et al [12]:

Type I: Identical code fragments except for variations in whitespace (may be also variations in layout) and comments.

Type II: Structurally/syntactically identical fragments except for variations in identifiers, literals, types, layout and comments.

Type III: Copied fragments with further modifications. Statements can be changed, added or removed in addition to variations in identifiers, literals, types, layout and comments.

Type IV: Two or more code fragments that perform the same computation but implemented through different syntactic variants.

A higher type of clone means that it's harder to detect and refactor. There are many studies that adopt these clone types, analyzing them further and writing detection techniques for them [13, 21, 22].

2.2.1 Type 4 clones

For this study we have chosen to take type 4 clones out of the scope, because they are both hard to detect and hard to refactor. A study by Kodhai et al [23] looks into the distribution of the different types of clones in several open source systems (see table 6 of his study). It becomes apparent that type 4 clones exist way less in source code than all of the other types of clones. For instance, for the J2sdk-swing system he finds 8115 type 1 clones, 8205 type 2 clones, 11209 type 3 clones and only 30 type 4 clones. Because of that, we can conclude that type 4 clones are relatively less relevant to study.

2.3 Refactoring methods

Martin Fowler recently released the second edition of his Refactoring book [1]. In this book he exclaims that *“if you see the same code structure in more than one place, you can be sure that your program will be better if you find a way to unify them”* [1, 11]. He describes several methods to deal with duplication, dependent on the context of the clone. In this section we describe the methods which we use in this study.

2.3.1 Extract Method

The main method for dealing with duplication in source code, as described by Martin Fowler, is to extract the duplication to a common place. This refactoring method is called “Extract Method”. Several studies have already concluded that most duplication in source code is found in the body of methods [24–26]. Method extraction can move matching functionality in method bodies to a common place, namely a new method. An example of this procedure is displayed in figure [fig:extractmethod].

<pre> 1 // Original 2 public void doStuff() { 3 doA(); 4 doB(); 5 doC(); 6 doA(); 7 doB(); 8 }</pre>	<pre> 1 // Refactored 2 public void doStuff() { 3 doAandB(); 4 doC(); 5 doAandB(); 6 } 7 8 public void doAandB() { 9 doA(); 10 doB(); 11 }</pre>
--	---

Figure 2.3: Refactoring duplication through method extraction.

2.3.2 Move method

In the example of the previous section, both duplicated parts are in the same method. However, a study by Fontana et al. [16] shows that this is most often not the case. Based on the relation between clone instances, the extracted method must be moved to be accessible by all locations of the clone instances. This might require different techniques.

2.3.2.1 Pull up method

A refactoring technique to move a method up in its inheritance structure is called “Pull up method” [1]. This way, if cloned methods are related in any way through inheritance, they can be called by both classes by placing the method in a class they both have in common. This way its possible to refactor both fully cloned methods (by just pulling up the method) and partially cloned methods (by first performing method extraction and then pulling up the refactored method).

2.3.2.2 Create class abstraction based on implicit relations

Duplication in source code is an implicit relation between fragments of source code. If two classes have many of these implicit relations, then the implementation should be refactored to make this relation explicit. If the classes do not yet have a parent/super class, a parent class can be created and the common functionality can be placed in this newly created class. This makes the relation between these classes explicit and reduces duplication.

2.3.2.3 Providing default implementations for common functionality

Looking specifically at the Java programming language, recent versions of Java introduce a new method of reducing duplication through refactoring. Java has the concept of *interfaces* as an abstract type that is used to specify a behavior that classes must implement. Since Java version 8, Java also supports default implementations to be provided in interfaces [27]. This gives an opportunity to make the relation between classes explicit and reduce duplication. It can be used in instances where creating a new parent class for duplicated classes is undesirable.

2.3.3 Clone refactoring in relationship to its context

How to refactor clones is highly dependent on their context. Method-level clones can be extracted to a method [23] if all occurrences of the clone reside in the same class. If a method level clone is duplicated among classes in the same inheritance structure, we might need to pull-up a method in the inheritance structure. If instances of a method level clone are not in the same inheritance structure, we might need to either make a static method or create an inheritance structure ourselves. So not only a single instance of a clone has a context, but also the relationship between individual instances in a clone class. This is highly relevant to the way in which the clone has to be refactored.

2.4 Internal and external classes

In this study we differentiate between *internal* and *external* classes. Classes are the components of a software system that contain its functionality. When analyzing a software system, *Internal classes* are classes that belong to this software system specifically, and its source code is included in the project.

External classes are classes that a software system uses, but do not belong to this specific software system. Typically, the source code of external classes are not included in the source code of the project. Most often these external dependencies are referenced in a file that its build automation system uses to fetch a projects’ dependencies.

Regarding refactoring, a software system can often not change the source code of its dependencies. This can sometimes be a burden, when an external dependency does not use proper abstractions that are required by a software system. Because of that, in this study, we differentiate between internal and external classes to see what impact they have on the refactoring process.

2.4.1 Maven

For this study we perform a shallow analysis on external dependencies, in order to derive more context for internally used concepts. As these external dependencies are most often not included in a software systems’ source code, we must first use the projects’ build automation system to gather the dependencies. To limit the scope of this process, we decided to focus on only the *Maven* build automation system.

Maven is a build automation tool, mainly used with the Java programming language. Maven has a simple ecosystem to configure and fetch the binaries and source code of all software projects a given codebase is dependent on. Maven can also be used to run tests for the systems and to package the project as an executable file.

Chapter 3

Defining refactoring-oriented clone types

In section 2.2 we introduced the four clone types as defined in literature. These simple definitions are suitable for analysis of a codebase. Their detection results in simple to understand numbers to argue about a codebase. However, these clone types have a few flaws which makes it hard to argue to what extend two fragments of code can and should be refactored. For each of type 1-3 clones [12] we list our solutions to their shortcomings to increase the chance that we can refactor the clone while improving the design.

We also look into clone detection tools for their suitability to support the proposed clone type definitions. We selected a few criteria. Most clone detection tools support these definitions of clone types. However, many of these tools use a vastly different approach. A study by Saini et al [28] outlines different clone detection tools and compares their results for each of type 1-3 clones. Even though they operate on the same type definitions, the tools used in this study yield different results.

3.1 Shortcomings of clone types

Clone types 1-3 (further explained in section 2.2) allow reasoning about the duplication in a software system. Clones by these definitions can relatively easily and efficiently be detected. This has allowed for large scale analyses of duplication [29]. However, these clone type definitions have shortcomings which makes the clones detected in correspondence with these definitions less valuable for (automated) refactoring purposes.

In this section, we discuss the shortcomings of the different clone type definitions. Because of these shortcomings, clones found by these definitions are often found to require additional judgment whether they should and can be refactored.

3.1.1 Type 1 clones

Type 1 clones are *identical clone fragments except for variations in whitespace and comments* [12]. This allows for the detection of clones that are the result of copying and pasting existing code, along with other reasons why duplicates might get into a codebase.

Type 1 clones are implemented as textual equality between code fragments (except for whitespace and comments) by most clone detection tools [14, 30–33]. Although textually equal, method calls can still refer to different methods, type declarations can still refer to different types and variables can be of a different type. In such cases, refactoring opportunities could be invalidated. This can make type 1 clones less suitable for refactoring purposes, as they require additional judgment regarding the refactorability of such a clone. When aiming to automatically refactor clones, applying refactorings to type 1 clones is bound to be error prone and can result in an uncompileable project or a difference in functionality.

1	<code>package com.sb.game;</code>	1	<code>package com.sb.fruitgame;</code>
2		2	
3	<code>import java.util.List;</code>	3	<code>import java.awt.List;</code>
4		4	
5	<code>public class GameScene</code>	5	<code>public class LoseScene</code>
6	<code>{</code>	6	<code>{</code>
7	<code>public void addToList(List l) {</code>	7	<code>public void addToList(List l) {</code>
8	<code>l.add(getClass().getName());</code>	8	<code>l.add(getClass().getName());</code>
9	<code>}</code>	9	<code>}</code>
10		10	
11	<code>public void addTen(int x) {</code>	11	<code>public void concatTen(String x) {</code>
12	<code>x = x + 10;</code>	12	<code>x = x + 10;</code>
13	<code>Notifier.notifyChanged(x);</code>	13	<code>Notifier.notifyChanged(x);</code>
14	<code>return x;</code>	14	<code>return x;</code>
15	<code>}</code>	15	<code>}</code>
16	<code>}</code>	16	<code>}</code>

Figure 3.1: Example of a type 1 clone that is functionally different.

In the example in figure 3.1, we see a type 1 clone displaying two clone classes. Defining an automatic way to refactor these clone classes is nearly impossible, as both cloned fragments describe different functional concepts. The first cloned fragment is a method that adds something to a `List`. However, the `List` objects to which something is added are different. Looking at the `import` statement above the class, one fragment uses the `java.util.List` and the other uses the `java.awt.List`. Both happen to have an `add` method, but apart from that their implementation is completely different.

The second cloned fragment shows how equally named variables can have different types and thus yield different functional concepts. The cloned fragment on the left adds a specific amount to an integer. The cloned fragment on the right concatenates a number to a `String`.

This shows that not all type 1 clones can easily be automatically refactored. In section we describe an alternate approach towards detecting type 1 clones, which results in only clones that can be refactored.

3.1.2 Type 2 clones

Type 2 clones are *structurally/syntactically identical fragments except for variations in identifiers, literals, types, layout and comments* [12]. This definition allows for the reasoning about code fragments that were copied and pasted, and then slightly modified. However, the definition does not adequately differentiate between slight modification and completely different fragments that just happen to have the same structure.

For refactoring purposes, this definition is unsuitable; if we allow any change in identifiers, literals, and types, we cannot distinguish between different variables, different types and different method calls anymore. This could render two methods that have an entirely different functionality as clones. Merging such clones can be harmful instead of helpful.

1	<code>public boolean redCircles</code>	1	<code>public Apple getEdibleApple</code>
2	<code>(List<Circle> circles){</code>	2	<code>(Basket<Apple> basket){</code>
3	<code>return circles.stream()</code>	3	<code>return basket.getFruit()</code>
4	<code>.allMatch(Shape::isRed);</code>	4	<code>.getApple(Fruit::notEaten);</code>
5	<code>}</code>	5	<code>}</code>

Figure 3.2: Example of a type 2 clone with high variability between fragments.

The example in figure 3.2 shows a type 2 clone class. Both methods are, except for their matching structure, completely different in functionality. They operate on different types, call different methods, return different things, etc. Having such a method flagged as a clone does not provide much useful information.

When looking at refactoring, type 2 clones can be difficult to refactor. For instance, if we have variability in types, the code can describe operations on two completely dissimilar types. Type 2 clones do not differentiate between primitives and reference types, which further undermines the usefulness of clones detected by this definition.

3.1.3 Type 3 clones

Type 3 clones are *copied fragments with further modification (with added, removed or changed statements)* [12]. Detection of clones by this definition can be hard, as it may be hard to detect whether a fragment was copied in the first place if it was severely changed. Because of this, most clone detection implementations of type 3 clones work on basis of a similarity threshold [18, 31, 32, 34]. This similarity threshold has been implemented in different ways: textual similarity (for instance using Levenshtein distance) [35], token-level similarity [13] or statement-level similarity [36].

Having a definition that allows for any change in code poses serious challenges on refactoring. A Levenshtein distance of one can already change the meaning of a code fragment significantly, for instance, if the name of a type differs by a character (and thus referring to different types).

3.2 Refactoring-oriented clone types

To resolve the shortcomings of clone types as outlined in the previous section, we propose alternative definitions for clone types directed at detecting clones that can and should be refactored. We have named these clones T1R (type 1R), T2R and T3R clones. These definitions address problems of the corresponding literature definitions. The “R” stands for refactoring-oriented (and may be less useful for other analyses).

3.2.1 Type 1R clones

To solve the issues identified in section 3.1.1, we introduce an alternative definition: cloned fragments have to be both textually *and* functionally equal. Like type 1 clones, type 1R clones do not consider comments and whitespace. Therefore, T1R clones are a subset of type 1 clones.

We check functional equality of two fragments by validating the equality of the fully qualified identifier (FQI) for referenced types, methods and variables. If an identifier is fully qualified, it means we specify the full location of its declaration (e.g. `com.sb.fruit.Apple` for an `Apple` object).

3.2.1.1 Referenced Types

Many object-oriented programming languages (like Java, Python and C#) require the programmer to import a type (or the class in which it is declared) before it can be used. Based on what is imported, the meaning of the name of a type can differ. For instance, if we import `java.util.List`, we get the interface which is implemented by all list datastructures in Java. However, importing `java.awt.List`,

we get a listbox GUI component for the Java Abstract Window Toolkit (AWT). To be sure we compare between equal types, type 1R clones compare the FQI for all referenced types.

3.2.1.2 Called methods

A codebase can have several methods with the same name. The implementation of these methods might differ. When we call two methods with an identical name, we can in fact call different methods. This is another reason that textually identical code fragments can differ functionally.

Because of this, for type 1R clones, we compare the fully qualified method signature for all method references. A fully qualified method signature consists of the fully qualified name of the method, the fully qualified type of the method plus the fully qualified type of each of its arguments. For instance, an eat method could become `com.sb.AppleCore com.sb.fruitgame.Apple.eat(com.sb.fruitgame.Tool)`.

3.2.1.3 Variables

In typed programming languages, each variable declaration should declare a name and a type. When we reference a variable, we only use its name. If, in different code fragments, we use variables with the same name but different types, the code can be functionally unequal but still textually equal. As an example, see the code in figure 3.3.

<pre> 1 public int addFive(int x){ 2 return x + 5; 3 }</pre>	<pre> 1 public String appendFive(String x){ 2 return x + 5; 3 }</pre>
--	---

Figure 3.3: Variables with different types but the same name.

The body of both methods in figure 3.3 is equal. However, their functionality is not. The first method adds two numbers together and the other concatenates an integer to a String.

For type 1R clones variable references should be compared by both type and name.

3.2.2 Type 2R clones

Type 2R clones are modelled after type 2 clones, which allow any change in identifiers, literals, types, layout, and comments. For refactoring purposes, this definition is unsuitable; if we allow any change in identifiers, literals, and types, we cannot distinguish between different variables, different types and different method calls anymore. This could render two methods that have an entirely different functionality as clones (as shown in figure 3.2 previously).

We tackle these problems with type 2R clones to be able to detect such clones that can and should be refactored. Type 1R clones are a subset of type 2R clones, meaning that each node found as cloned for type 1R will also be found as cloned for type 2R. Similar to type 1R, for type 2R we consider the fully qualified identifiers of type-, method- and variable-references. Additionally, type 2R clones allow variability in a controlled set of expressions and identifiers. The identifiers and expressions in which we allow variability must have one of the following properties:

- We allow a difference in identifiers between cloned fragments if it has no impact on the refactoring process.
- We allow a difference in expressions between cloned fragments if the expression can be extracted to a method argument when applying the “Extract Method” refactoring.

Allowing variability between expressions includes a tradeoff. When many expressions vary between cloned the refactored method might require many additional arguments. Because of that, we constrain this variability by a threshold, which is explained in more detail in section 3.2.2.2.

3.2.2.1 Allow any variability in some identifiers

When refactored, some identifiers have no detrimental effect on the design if they vary between cloned instances. This section explains several patterns of variability that we can allow between cloned fragments

without changing the method by which the fragments can be refactored.

3.2.2.1.1 Declaration names The names of declarations describe the implementation of their body. Examples of declarations are:

- Class declaration
- Method declaration
- Interface declaration (not all languages have a separate declaration for these)
- Enumeration declaration (not all languages have a separate declaration for these)
- Annotation declaration (not all languages support these)
- Etc.

If two declaration bodies and signatures are cloned, but their names differ, one of both names should be redundant. When refactoring such clones, we can choose one instance to keep and one to remove. Such a refactoring doesn't affect maintainability in any other way than refactoring a type 1R clone would. Because of this, type 2R allows any variability in declaration names. However, this will only open up a good refactoring opportunity if the entire body and signature of these declarations are cloned.

3.2.2.1.2 Variable names Cloned code fragments using variables with different names can be refactored without a design tradeoff if the following conditions apply:

- The cloned variables are locally defined or the refactoring method used requires extraction to method argument anyways.
- The cloned variables have the same type.
- The cloned variables are used at the same places in cloned fragments.

Show ex-
amples

3.2.2.2 A threshold for variability in literals, variables and method calls

Type 2 clones allow any variability in literals, variables and method identifiers. However, this information tells a lot about the meaning of the code fragment. Most clone detection tools do not differentiate between a type 2 clone that differs by a single literal/identifier and one that differs by many [37]. However, this does have a big impact on the meaning of the code fragment.

For type 2R clones we define a threshold for variability in literals, variables and method calls. We calculate the variability in literals, variables and method calls using the following formula:

$$\text{T2R Variability} = \frac{\text{Number of different expressions}}{\text{Total number of expressions in clone instance}} * 100 \quad (3.1)$$

In this formula, *number of different expressions* refers to the number of literals, variables and method calls that differ from other clone instances in a clone class. We divide this by the total number of expressions in the clone instance. Based on this threshold, we decide whether a clone should be considered for refactoring. A concrete example of applying this formula to calculate a threshold is given in 4.3.3.2.

3.2.2.2.1 Literal and variable variability We allow only variability in the value of literals and variables, but not in their types. This is because a difference in literal/variable type may have a big impact on the refactorability of the cloned fragment. When we refactor different literals/variables that have both the same type, in case of an “Extract Method” refactoring, we have to create a parameter for this literal and pass the corresponding literal/variable from cloned locations. However, if two literals have different types, this might not be possible (or will have a negative effect on the design of the system). This is because a lot of variability in literals will result in more parameters required in the extracted method, which is detrimental for the design of the system.

Consider the example in figure 3.4. In this example, the two methods have two literals that differ between them. We can perform an “Extract method” refactoring on these to get the result that is displayed on the right. In this process, we create a method parameter for the corresponding literal.

<pre> 1 // Original 2 public void printNice(String s){ 3 System.out.println("====="); 4 System.out.println(s); 5 System.out.println("====="); 6 } 7 8 public void printNice2(String s){ 9 System.out.println("-----"); 10 System.out.println(s); 11 System.out.println("-----"); 12 } 13 14 </pre>	<pre> 1 // Refactored 2 public void printNice(String s){ 3 printNice(s, "====="); 4 } 5 6 public void printNice2(String s){ 7 printNice(s, "-----"); 8 } 9 10 public void printNice(String s, String decoration){ 11 System.out.println(decoration); 12 System.out.println(s); 13 System.out.println(decoration); 14 } </pre>
--	---

Figure 3.4: Literal variability refactored.

3.2.2.2.2 Method call variability Most modern programming languages (like Java, Python and C#) allow to pass method references as a parameter to a method. This helps reducing duplication, as it is possible to refactor two code fragments which differ only by a method call. However, a method call does often not consist of a single token (like variables and literals). For instance, a method call `System.out.println()` consists of several segments: a type reference to the `System` type, a reference to the static `out` field and a call of the `println()` method.

Type 2R clones allow called methods to vary as long as they have the same argument types and return type. As with type 1R clones, these types are compared using their fully qualified identifiers. An example of this is shown in figure 3.5. In this example, we have two methods (`System.out.println` and `myFancyPrint`). We use the “Extract Method” refactoring method to extract a new method and use a parameter to pass the used method.

The method call variability property of type 2R clones imply that type 2R clones are not a subset of type 2 clones. Because methods calls can have a different structure, type 2R clones can be structurally slightly different. The example as shown in figure 3.5 can be a type 2R clone (dependent on the thresholds used), but is not a type 2 clone.

<pre> 1 // Original 2 public void printNice(){ 3 System.out.println("====="); 4 myFancyPrint("The weather is nice"); 5 System.out.println("====="); 6 } 7 8 public void printNice2(){ 9 System.out.println("====="); 10 System.out.println("The weather is nice"); 11 System.out.println("====="); 12 } 13 14 public void myFancyPrint(String s){ 15 Arrays.stream(s.toCharArray()).boxed() 16 .map(e -> "("+e+")") 17 .collect(Collectors.joining()); 18 } 19 20 </pre>	<pre> 1 // Refactored 2 public void printWeather(){ 3 printNice(this::myFancyPrint); 4 } 5 6 public void printNice2(){ 7 printNice(System.out::println); 8 } 9 10 public void printNice(Consumer<String> printFunction){ 11 System.out.println("====="); 12 printFunction.accept("The weather is nice"); 13 System.out.println("====="); 14 } 15 16 public void myFancyPrint(String s){ 17 Arrays.stream(s.toCharArray()).boxed() 18 .map(e -> "("+e+")") 19 .collect(Collectors.joining()); 20 } </pre>
---	---

Figure 3.5: Method variability refactored.

3.2.3 Type 3R clones

Type 3 clones allow any change in statements, often bounded by a similarity threshold. This means that type 3 clones allow the inclusion of a statement that is not detected by type 1 or 2 clone detection. When looking at how we can refactor a statement that is not included by one clone instance but is in another, we find that we require a conditional block to make up for the difference in statements. See figure 3.6 for an example of such a clone.

```

1 // Original
2 public String wordOfNumbers(int max){
3     Random r = new Random();
4     StringBuilder word = new StringBuilder();
5     for(int i = 0; i<=size; i++){
6         word.append(i);
7     }
8     word.append("that's it");
9     return word.toString();
10 }
11
12 public void printWordOfNumbers(int max){
13     Random r = new Random();
14     StringBuilder word = new StringBuilder();
15     for(int i = 0; i<=size; i++){
16         word.append(i);
17         System.out.println("Appended "+i);
18     }
19     word.append("that's it");
20     System.out.println(word);
21 }
    
```

```

1 // Refactored
2 public List<String> wordOfNumbers(int size){
3     return wordOfNumbers(n, size, false);
4 }
5
6 public void printWordOfNumbers(int size){
7     System.out.println(wordOfNumbers(n, size, true));
8 }
9
10 public StringBuilder wordOfNumbers(int size,
11                                     boolean doPrint){
12     Random r = new Random();
13     StringBuilder word = new StringBuilder();
14     for(int i = 0; i<size; i++){
15         word.append(i);
16         if(doPrint) System.out.println("Appended "+i);
17     }
18     word.append("that's it");
19     return word;
20 }
21
    
```

Figure 3.6: Added statement between cloned methods.

In figure 3.6 a single statement is added. This statement is found in between cloned lines. We have named this difference between the two clones, in this examples containing a single statement, the *gap*. This gap has a threshold for type 3R, which is calculated by the following formula:

$$\text{T3R Gap Size} = \frac{\text{Number of nodes in gap}}{\text{Number of nodes in clones}} * 100 \quad (3.2)$$

When applying this formula with the example given in figure 3.7, we get a gap size of 25%:

$$\text{T3R Gap Size} = \frac{|\{\text{doC}()\}|}{|\{\text{doA}(), \text{doB}(), \text{doD}(), \text{doE}()\}|} * 100 = \frac{1}{4} * 100 = 25\% \quad (3.3)$$

```

1 public void doStuff(){
2     doA(); // Clone instance of
3     doB(); // two statements
4     doC(); // Gap of one statement
5     doD(); // Clone instance of
6     doE(); // two statements
7 }
    
```

```

1 public void doStuff(){
2     doA(); // Clone instance of
3     doB(); // two statements
4     doF(); // Gap of one statement
5     doD(); // Clone instance of
6     doE(); // two statements
7 }
    
```

Figure 3.7: Gapped clones.

Apart from this threshold, the gap in between clones may not span over different blocks. Look at figure 3.8 for an example of this. We cannot refactor both statements into a single conditional block. We could however use two conditional blocks, but due to the detrimental effect on the design of the code (as each conditional block adds a certain complexity), we decided not to allow this for type 3R clones.

```
1  public String wordOfNumbers(int max){
2      Random r = new Random();
3      StringBuilder word = new StringBuilder();
4      for(int i = 0; i<=size; i++){
5          word.append(i);
6      }
7      word.append("that's it")
8      return word.toString();
9  }
10
11 public void printWordOfNumbers(int max){
12     Random r = new Random();
13     StringBuilder word = new StringBuilder();
14     for(int i = 0; i<=size; i++){
15         word.append(i);
16         System.out.println("Appended "+i);
17     }
18     System.out.println(word);
19     word.append("that's it");
20 }
```

Figure 3.8: Statements between clones in different blocks.

3.2.4 Clone types summarized

The given clone definitions (types 1R, 2R and 3R) are refactoring-oriented in the sense that they were designed after the literature type definitions but with a concrete refactoring opportunity in mind. Summarized, these types can be explained as follows:

- **Type 1R:** Allows no difference between cloned fragments (both functionally and textually), making it possible to refactor both fragments to a method call that contains the code of both locations.
- **Type 2R:** Allows difference between cloned fragments in a controlled set of expressions. Refactoring opportunities for these controlled features are known, allowing refactoring with a minor tradeoff.
- **Type 3R:** Allows any difference. When refactored, this difference must be wrapped in a conditionally executed block, which entails a major tradeoff.

Regarding the nodes found as cloned by these clone type definitions, there exists a subset relation between all types:

$$\text{Cloned nodes} = T1R \subseteq T2R \subseteq T3R \quad (3.4)$$

Comparing these clone type definitions to the ones used in literature, type 1R clones are a subset of type 1 clones. However, this subset relation does not exist for type 2R and 3R clones. This is because these definitions allow variability in multi-token method calls, which type 2 and type 3 do not allow. Because of that, the following is true for these clone types:

$$\text{Cloned nodes} = T1R \subseteq T1, \text{ but } T2R \not\subseteq T2 \text{ and } T3R \not\subseteq T3 \quad (3.5)$$

3.2.5 The challenge of detecting these clones

To detect each type of clone, we need to parse the fully qualified identifier of all types, method calls and variables. This comes with serious challenges, regarding both performance and implementation. Also, to be able to parse all fully qualified identifiers, and trace the declarations of variables, we might need to follow cross file references. The referenced types/variables/methods might even not be part of the project, but rather of an external library or the standard libraries of the programming language. All these factors need to be considered for the referenced entity to be found, on basis of which a fully qualified identifier can be created.

3.3 Suitability of existing Clone Detection Tools for detecting these clones

We conducted a short survey on (recent) clone detection tools that we could use to analyze refactoring possibilities. The results of our survey are displayed in table 3.1. We chose a set of tools that are open source and can analyze a popular object-oriented programming language. Next, we formulate the following four criteria by which we analyze these tools:

Is this
still up
to date?

1. **Should find clones in any context.** Some tools only find clones in specific contexts, such as only method level clones. We want to perform an analysis of all clones in projects to get a complete overview.
2. **Finds clone classes in control projects.** We assembled a number of control projects to assess the validity of clone detection tools.
3. **Can analyze resolved symbols.** When detecting the types proposed in section 3.2, it is important that we can analyze resolved symbols (for instance a type reference). The rationale for this is further explained in 3.2.5.
4. **Extensive detection configuration.** Detecting our clone definitions, as proposed in section 3.2, require to have some understanding about the meaning of tokens in the source code (whether a certain token is a type, variable, etc.). The tool should recognize such structures, in order for us to configure our clone type definitions in the tool.

Table 3.1: Our survey on clone detection tools.

	(1)	(2)	(3)	(4)
Siamese [34]				✓
NiCAD [32, 38]	✓	✓		
CPD [37]	✓	✓		
CCFinder [30]	✓	✓		
D-CCFinder [29]				
CCFinderSW [31]	✓			✓
SourcererCC [13]	✓			✓
Oreo [39]				
BigCloneEval [14]	✓	✓		
Deckard [18]	✓		✓	
Scorpio [36, 40]	✓		✓	✓

None of the state-of-the-art tools we identified implement all our criteria, so we decided to implement our own clone detection tool, which is further described in the next chapter.

Chapter 4

CloneRefactor

CloneRefactor is the name of our clone detection and refactoring tool. It features the following novel functions:

- Detection of clone classes rather than clone pairs.
- A novel detection method, aimed at extensibility.
- Detection of refactoring-oriented clone types, in addition to the literature clone types.
- Allows for automated refactoring of a subset of the detected duplication issues.

In this section we describe our approach and rationale for the design decisions regarding this tool.

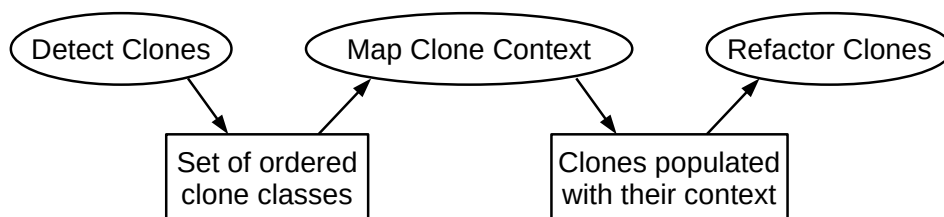


Figure 4.1: CloneRefactor overall process.

Figure 4.12 shows the overall process used by CloneRefactor. First, we detect clones on basis of a Java codebase, given a Java project from disk and a configuration. The clone detection process is further explained in section 4.3. After all clones are found, CloneRefactor maps the context of the clones. On basis of this context, CloneRefactor applies transformations to the source code for clones for which we have configured a refactoring.

4.1 JavaParser

A very important design decision for CloneRefactor is the usage of a library named JavaParser [41]. JavaParser is a Java library which allows to parse Java source files to an abstract syntax tree (AST¹). JavaParser allows to modify this AST and write the result back to Java source code. This allows us to apply refactorings to the detected problems in the source code.

Integrated in JavaParser is a library named SymbolSolver. This library allows for the resolution of symbols using JavaParser. For instance, we can use it to trace references (methods, variables, types, etc) to their declarations (these referenced identifiers are also called “symbols”). This is very useful for the detection of our refactoring-oriented clone types, as they make use of the fully qualified identifiers of symbols.

In order to be able to trace referenced identifiers SymbolSolver requires access to not only the analyzed Java projects, but also all its dependencies. This requires us to include all dependencies with the project. Along with this, SymbolSolver solves symbols in the JRE System Library (the standard libraries coming with every installation of Java) using the active Java Virtual Machine (JVM). This has a big impact on performance efficiency.

¹An AST is a tree-representation of a source code file.

Because of the requirement of symbol resolution, the refactoring-oriented clone types are less suitable for large scale clone analysis.

4.2 Thresholds

CloneRefactor operates on a set of thresholds to determine the validity of clones that are found. In general, these thresholds are:

- **Minimum Amount of Nodes:** The minimum amount of nodes that should be in each instance of cloned fragments for them to be considered a clone class.
- **Minimum Amount of Tokens:** The minimum amount of tokens that should be in each instance of cloned fragments for them to be considered a clone class.
- **Minimum Amount of Lines:** The minimum amount of lines that should be in each instance of cloned fragments for them to be considered a clone class.

When we analyze type 2R or type 3R clones we use the following additional threshold:

- **T2R Variability:** The percentage of variability we allow in variables, method calls and literals. How this threshold is calculated is explained in section 4.3.3.2.

When we analyze type 3 or type 3R clones we use the following additional threshold:

- **T3R Gap Size:** The size of the gap we allow between two valid clones. This is a percentage of the size of the gap against the size of both clones combined. How this threshold is calculated is further explained in section 4.3.4.

4.3 Clone Detection

To detect clones, CloneRefactor parses the AST acquired from JavaParser to an unweighted graph structure. On basis of this graph structure, clones are detected. Dependent on the type of clones being detected, transformations may be applied. The way in which CloneRefactor was designed does not allow for several clone types to be detected simultaneously, in accordance with our clone type philosophy as described in chapter ??.

The overall process regarding clone detection is displayed in figure 4.2. First of all, we use JavaParser to read a project from disk and build an AST, one class file at a time. Each AST is then converted to a directed graph that maps relations between statements, further explained in section 4.3.1. On basis of this graph, we detect clone classes and verify them using three thresholds (in order of importance):

- Amount of tokens.
- Amount of statements.
- Amount of lines.

If the detection was configured to detect either type 2R, 3 or 3R clones we perform some type specific transformations on the resulting set of clones.

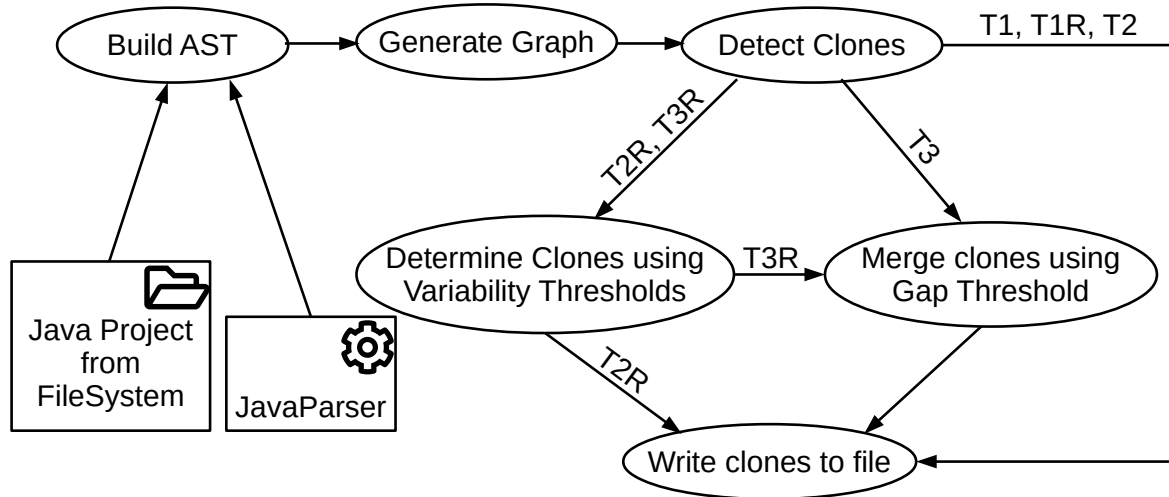


Figure 4.2: CloneRefactor clone detection process.

4.3.1 Generating the clone graph

First of all, we parse the AST obtained from JavaParser into a directed graph structure. We have chosen to base our clone detection around statements as the smallest unit of comparison. This means that a single statement cloned with another single statement is the smallest clone we can find. The rationale for this lies in both simplicity and performance efficiency. This means we won't be able to find when a single expression matches another expression, or even a single token matching another token. This is in most cases not a problem, as expressions are often small and do not span the minimal size to be considered a clone in the first place (more about this in section 5.3).

4.3.1.1 Filtering the AST

As a first step towards building the clone graph, we preprocess the AST to decide which AST nodes should become part of the clone graph. We have decided to consider declarations and statements as the smallest compared entities. The main reasoning for this is because considering smaller AST constructs, like expressions, significantly increases the complexity and CPU usage of our clone detection and refactoring efforts.

Additionally, we exclude package declarations and import statements. These are omitted by most clone detection tools, as clones in import statements hold limited valuable information.

4.3.1.2 Building the clone graph

Building the clone graph consists of walking the AST in-order for each declaration and statement. For each declaration/statement found, we map the following relations:

- The declaration/statement preceding it.
- The declaration/statement following.
- The last **preceding** declaration/statement with which it is cloned.

We do not create a separate graph for each class file, so the statement/declaration preceding or following could be in a different file. While mapping these relations, we maintain a hashed map containing the last occurrence of each unique statement. This map is used to efficiently find out whether a statement is cloned with another. An example of such a graph is displayed in figure 4.3.

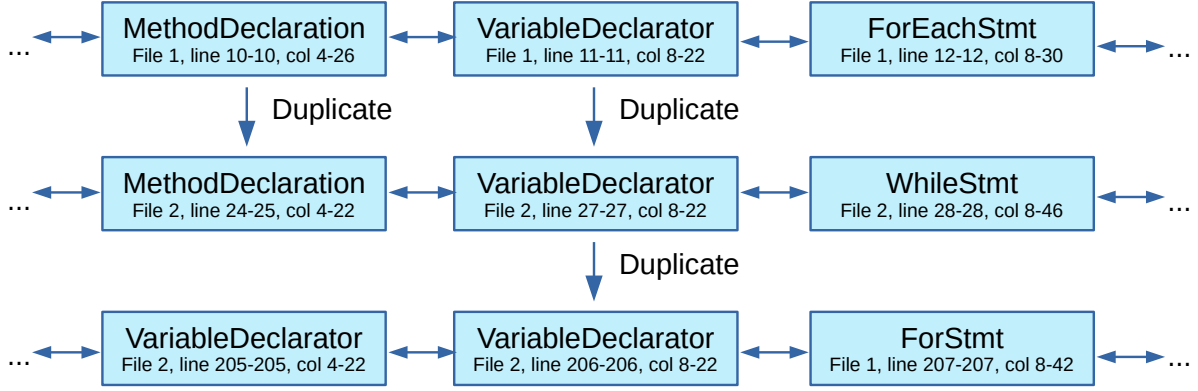


Figure 4.3: Abstract example of a part of a possible clone graph as built by CloneRefactor.

We refer to the declarations and statements in this graph as *nodes*. The relations *next* and *previous* in this graph are represented as an twodirectional arrow. The relations representing duplication are directed. This is a restriction we’ve chosen as it creates an important constraint for the clone detection process. This process is explained in section 4.3.2.

4.3.1.3 Comparing Statements/Declarations

In the previous section we described a “duplicate” relation between nodes in the clone graph built by CloneRefactor. Whether two nodes in this graph are duplicates of each other is dependent on the clone type. In this section, we will describe for each type how we compare statements and declarations to assess whether they are clones of each other.

CloneRefactor detects six different types of clones: T1, T2, T3, T1R, T2R and T3R. These types are further explained in chapter 3. For **type 1** clones, CloneRefactor filters the tokens of a node to exclude its comments, whitespace and end of line (EOL) characters and then compares these tokens. For **type 2** clones, the tokens are further filtered to omit all identifiers and literals. **Type 3** clones do the same duplication comparison as type 2 clones.

For **type 1R** clones, this comparison is a lot more advanced. For *method calls* we trace their declaration and use its fully qualified method signature for comparison with other nodes. For all *referenced types* we trace their declarations and use assemble their fully qualified identifier for comparison with other nodes. For *variables* we trace their declaration and their types. If the variable type is a primitive we can directly use it for comparison. If it is a referenced type, we have to trace this type first in order to collect their fully qualified identifier for comparison.

Type 2R clones allow any variation in literals, variables and method calls at this stage in the clone detection process. However, for *literals* we do resolve their type in order to verify that they are of the same type. For *variables* we also only verify that their types are the same (but not their names). For *method calls*, we trace their declaration but only compare the fully qualified identifier for its return type and each of its arguments’ types. Apart from that, we do not compare the names of method, class, interface and enum declarations.

In this stage, **type 3R** clones have the same compare rules as type 2R clones.

4.3.1.4 Mapping graph nodes to code

The clone graph, as explained in section 4.3.1.2, contains all declarations and statements in a source code. However, declarations and statements may themselves have child declarations and statements. To avoid redundant duplication checks, we exclude child declarations and statements from each node. Look at figure 4.4 for an example of how source code maps to AST nodes.

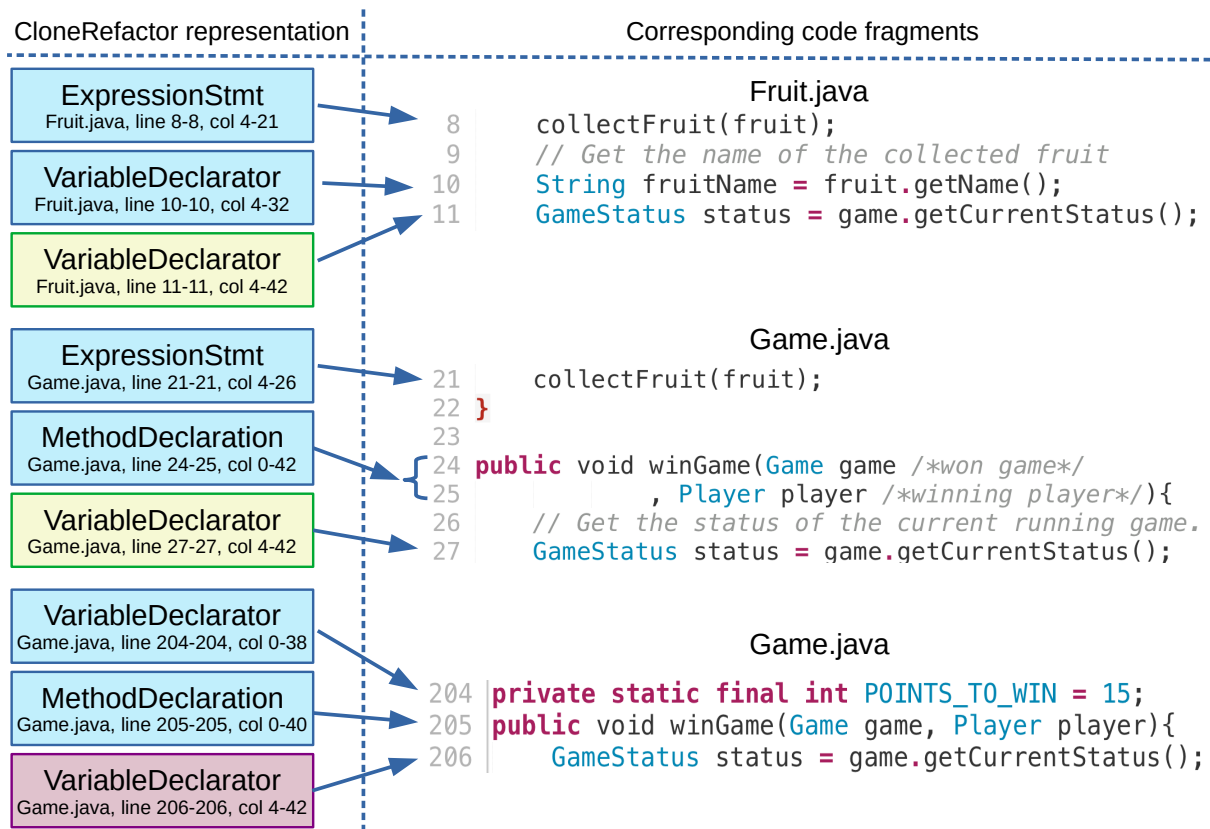


Figure 4.4: CloneRefactor extracts statements and declarations from source code.

In line 24-25 of the code fragment, we see a MethodDeclaration. The node corresponding with this MethodDeclaration denotes all tokens found on these two lines, line 24 and 25. Although the statements following this method declaration (those that are part of its body) officially belong to the method declaration, they are not included in its graph node. Because of that, in this example, the MethodDeclaration on line 24-25 will be considered a clone of the MethodDeclaration on line 205 even though their bodies might differ. Even the range (the line and column that this node spans) does not include its child statements and declarations.

4.3.2 Detecting Clones

After building the clone clone graph, we use it to detect clones. We decided to focus on the detection of clone classes rather than clone pairs because clone pairs do not provide a general overview of all entities containing the clones, with all their related issues and characteristics [42]. Although clone classes are harder to manage, they provide all information needed to plan a suitable refactoring strategy, since this way all instances of a clone are considered. Another issue that results from grouping clones by pairs: clone reference amount increases according to the binomial coefficient formula (two clones form a pair, three clones form three pairs, four clones form six pairs, and so on), which causes a heavy information redundancy [42].

As stated in the previous section, nodes in the graph link to *preceding* cloned statement. This implies that the first node that is cloned does not have any clone relation, as there are no clones preceding it (only following it). Because of this, we start our clone detection process at the final location encountered while building the graph. For this node, we collect all nodes it is cloned with. Even though the final node only links to the preceding node it is cloned with, we can collect all clones. This is because the preceding clone also has a preceding clone (if applicable) and we can follow this trail to collect all clones of a single node. As an example, we convert the code example shown in figure 4.4 to a clone graph as displayed in figure 4.5.

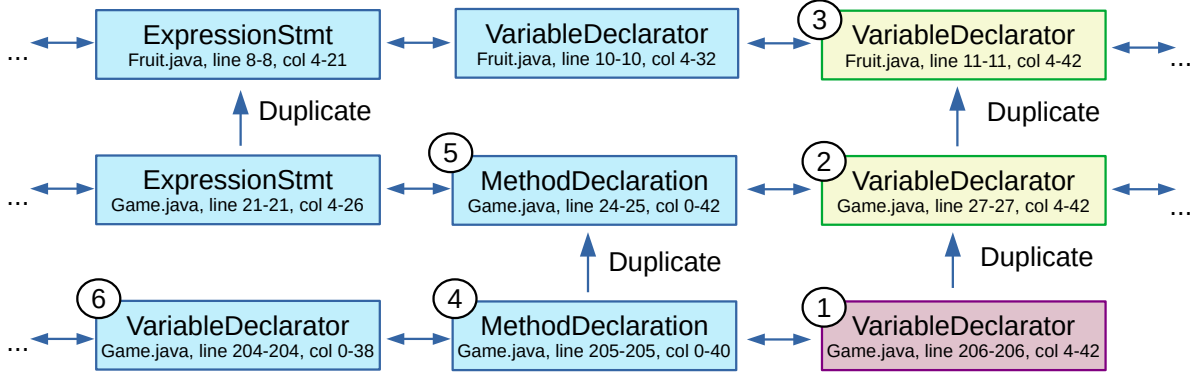


Figure 4.5: Example of a clone graph built by CloneRefactor.

Using the example shown in figure 4.4 and 4.5 we can explain how we detect clones on basis of this graph. Suppose we are finding clones for two files and the final node of the second file is a variable declarator. This final node is represented in the example figure by the purple box (1). We then follow all “duplicate” relations until we have found all clones of this node (2 and 3). We now have a single statement clone class of three clone instances (1, 2 and 3).

Next, we move to the previous line (4). Here again, we collect all duplicates of this node (4 and 5). For each of these duplicates, we check whether the node following it is already in the clone class we collected in the previous iteration. In this case, (2) follows (5) and (1) follows (4). This means that node (3) does not form a ‘chain’ with other cloned statements. Because of this, the clone class of (1, 2 and 3) comes to an end. It will be checked against the thresholds, and if adhering to the thresholds, considered a clone.

We then go further to the previous node (6). In this case, this node does not have any clones. This means we check the (2 and 5, 1 and 4) clone class against the thresholds, and if it adheres, consider it a clone. Dependent on the thresholds, this example can result in a total of two clone classes.

Eventually, following only the “previous node” relations, we can get from (6) to (2). When we are at that point, we will find only one cloned node for (2), namely (3). However, after we check this clone against the thresholds, we check whether it is a subset of any existing clone. If this is the case (which it is for this example), we discard the clone.

4.3.2.1 Removing redundant clone classes

The clone detection method used by CloneRefactor can, for various reasons, result in redundant clone classes. After the insertion of each newly detected clone, we check whether it is redundant and/or any of the existent clones has become redundant by adding this clone. A clone is redundant if it is a subset of another clone. We define the subset relation between clones as follows:

$$C_1 \subseteq C_2 := \forall (i_1 \in C_1) \exists (i_2 \in C_2) F i_1 = F i_2 \wedge R i_1 \subseteq R i_2 \quad (4.1)$$

Where C refers to a clone class (a set of clone instances), i refers to a clone instance, F is the file in which a clone instance is located and R is the range of tokens that a clone instance spans. For each clone found, we remove all existing clones that are a subset of the found clone:

$$S_{after} = S_{before} \setminus \{C_{existing} \subseteq C_{new} \mid C_{existing} \in S_{before}\} \quad (4.2)$$

Where S_{before} is the clone class collection containing all clones that are found up in until this point and C_{new} is the clone class that was just found. S_{after} is the clone class collection after

We should not add the new clone to our list of clones if its a subset of an existing clone. Because of that, we check for each clone added whether there exists a clone class of which the found clone class is a subset:

$$\{C_{existing} \subseteq C_{new} \mid C_{existing} \in S\} = \emptyset \Rightarrow S_{after} = S \cup C_{new} \quad (4.3)$$

If the newly added clone is a subset of an existing clone, we do not add it to the set of clone classes. This way we avoid redundant clone classes being detected by CloneRefactor.

4.3.3 Validating the type 2R variability threshold

In the definition of type 2R clones (see section 3.2.2) we described how type 2R clones work on basis of a variability threshold. This threshold is checked by CloneRefactor to assess the validity of found clone classes. Implementing such a threshold involves some important design decisions and has a lot of underwater complexity. In this section we explain how CloneRefactor detects type 2R clones on basis of this variability threshold. This process is done as a postprocessing step after clone detection. The type 2R clone detection process is described in more detail in section 4.3.1.3. In short, this process detects clones allowing for any variability between expressions. This postprocessing step then determines which (parts of) these clones are valid by the configured variability threshold.

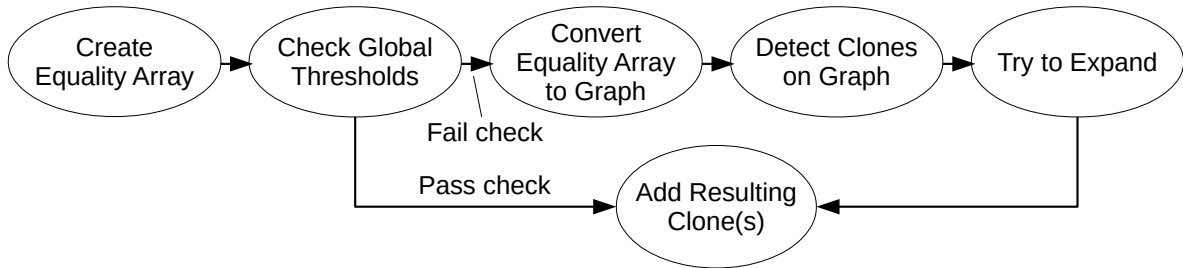


Figure 4.6: Process used to check the variability threshold for T2R clones.

Figure 4.6 shows the steps that CloneRefactor performs to find clones conforming with the type 2R variability threshold. Each of the following paragraphs will explain a step from this figure.

4.3.3.1 Create equality array

To determine the difference in literals, method calls and variables, we convert the code to an equality array. This equality array converts each (group of) token(s) to a number unique to that (group of) tokens. Each literal, method call or variable becomes a positive number, whereas each other token becomes a negative number. An example of two code fragments converted to equality arrays is displayed in figure 4.7.

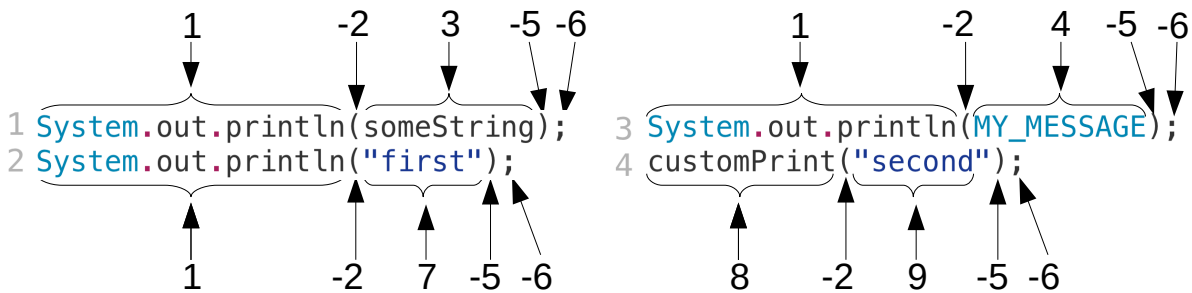


Figure 4.7: The conversion of code to an equality array.

The equality array for each of the lines in this example is shown in table 4.1.

Table 4.1: The equality arrays in our example figure 4.7.

Node (n)	Equality Array (E)
1	[1, -2, 3, -5, -6]
2	[1, -2, 7, -5, -6]
3	[1, -2, 4, -5, -6]
4	[8, -2, 9, -5, -6]

In the example of figure 4.7 the fragment on the left (consisting of n_1 and n_2) and the fragment on the right (consisting of n_3 and n_4) are two clone instances of a clone class. Table 4.1 shows their corresponding equality arrays.

4.3.3.2 Checking global thresholds

Using the equality arrays explained in the previous section, we can determine the variability threshold of any clone class. We calculate the variability of the example given in table 4.1 as follows:

$$\text{T2R Variability} = \frac{|\{x > 0 \wedge y > 0 \wedge x \neq y \Rightarrow (x, y) \mid x \in E_1 \cup E_2, y \in E_3 \cup E_4\}|}{|\{x > 0 \wedge y > 0 \Rightarrow (x, y) \mid x \in E_1 \cup E_2, y \in E_3 \cup E_4\}|} * 100 \quad (4.4)$$

Where E_x refers to the equality arrays shown in table 4.1 and x is the node number displayed in the same table. When we apply this equation to the example clone classes displayed in figure 4.7, we get the following sum:

$$\frac{|\{(3, 4), (1, 8), (7, 9)\}|}{|\{(1, 1), (3, 4), (1, 8), (7, 9)\}|} * 100 = \frac{3}{4} * 100 = 75\% \quad (4.5)$$

So for the example given in figure 4.7 we have a variability of 75%. In CloneRefactor, the maximum variability percentage is a threshold that is entered in a configuration file. If a clone adheres to this threshold, it will stay in the set of found clones. However, if it does not adhere to the thresholds, a problem arises. Because a clone does not adhere to the thresholds, it does not yet mean it has to be discarded. This is because an invalid clone class can still contain valid clones that are a subset of the invalid clone (our definition of subsets of clones is given in section 4.3.2.1).

When a found clone class does not adhere to the global thresholds as explained in the previous section, we need to determine whether it contains any valid subclones. Below, we explain a few cases of valid subclones that may exist within an invalid clone class.

1 doA(a, b, c);	1 doB(d, e, f);
2 doA();	2 doA();
3 doB();	3 doB();
4 doC();	4 doC();

Figure 4.8: One node in a type 2R clone has a high variability.

The first line of the cloned fragment shown in figure 4.8 has a high variability with its cloned fragment. However, the rest of this method does not have any variability. The global thresholds could indicate a too high variability and thus render this clone invalid. However, in this case, it might still have a valid subclone.

1 doA();	1 doA();	1 doD();
2 doB();	2 doB();	2 doE();
3 doC();	3 doC();	3 doF();

Figure 4.9: One clone instance in a type 2R clone has a high variability.

Another example of high variability between clones, in which a valid subclone can be found, is displayed in figure 4.9. In this case, one clone instance has such a high variability that it shouldn't be refactored. In this case, the clone instance with high variability should be removed.

1	doA ();	1	doD ();	1	doE ();
2	doB ();	2	doA ();	2	doE ();
3	doC ();	3	doB ();	3	doA ();
4	doC ();	4	doD ();	4	doB ();

Figure 4.10: A small subset of nodes has a high variability.

Figure 4.10 shows an example of a clone class where the valid sequence inside the clone does not align.

If the check for global thresholds fails, we have to seek for valid clones within the clone class. In a single invalid clone class can be zero to many subclones. This requires an extensive search for such clones. This problem is very related to the problem of clone detection, except now it is within the boundaries of a single clone class except for an entire codebase. Because of that, just like the clone detection process, we build a clone graph and detect clones on it. This process is explained over the following sections.

4.3.3.3 Convert Equality Array to Graph

After the global threshold check has failed, we build a clone graph on basis of the equality arrays. The process used for building this graph is the same as the process described in 4.3.1.2.

1	doA (); // {1, -2}	1	doD (); // {5, -2}	1	doE (); // {6, -2}
2	doB (); // {3, -2}	2	doA (); // {1, -2}	2	doE (); // {6, -2}
3	doC (); // {4, -2}	3	doB (); // {3, -2}	3	doA (); // {1, -2}
4	doC (); // {4, -2}	4	doD (); // {5, -2}	4	doB (); // {3, -2}

Figure 4.11: Cloned code fragments from figure 4.10 together with their equality arrays.

In figure 4.11 we display clone classes and their (simplified) equality arrays. We can convert this to a graph, using a similar process as the graph creation process used for clone detection. In this case, duplicate relations represent nodes of which their equality arrays are within the variability threshold. In figure 4.11 all relations are between equal equality arrays, but this does not always have to be the case. Large equality arrays, denoting statements that consist of many tokens, may have some variability and still be considered duplicates in this graph.



Figure 4.12: Graph representation of the code example displayed in figure 4.11.

In this graph, the duplicate relations (green arrows) do not necessarily denote an exact match. The duplicate relations allow for variability in the equality arrays as long as it falls under the variability threshold. However, negative numbers in the equality array may never vary. This is because negative numbers are the tokens/expressions that we do not allow variability in because they cannot be refactored as easily:

$$E_1 \text{ can only be cloned with } E_2 \text{ iff } |E_1| = |E_2| \wedge \forall (x \in E_1, y \in E_2) \ x < 0 \vee y < 0 \Rightarrow x = y \quad (4.6)$$

4.3.3.4 Detect Clones on Graph

Similar to the clone detection process, we detect clones on basis of the graph described in the previous section. The only difference is that, in this step, we do not remove clones that do not meet the thresholds (as explained in section 4.2). This is because the next step, explained in the next section, could potentially expand the found clones and thus make clones that would currently be invalid by the thresholds valid.

Looking back at the example of the previous section, it would result in the following clone class collection:

1	doA(); // C ₁ I ₃	1	doD(); // C ₃ I ₂	1	doE(); // C ₂ I ₂
2	doB(); // C ₁ I ₃	2	doA(); // C ₁ I ₂	2	doE(); // C ₂ I ₁
3	doC(); // C ₄ I ₂	3	doB(); // C ₁ I ₂	3	doA(); // C ₁ I ₁
4	doC(); // C ₄ I ₁	4	doD(); // C ₃ I ₁	4	doB(); // C ₁ I ₁

Figure 4.13: Clone instances in figure 4.11 as determined .

Where C denotes the clone class and I denotes the clone instance. The process by which these clones are found is equal to the clone detection process as described in section 4.3.2. In this example, we see the four clones found on the graph of figure 4.12. Three of these clone classes have a two clone instances, each consisting of a single nodes. The other clone class has three clone instances, each consisting of two nodes.

4.3.3.5 Try to Expand

The problem with the graph clone detection technique is that only single nodes that are within the variability threshold are considered duplicates of each other. However, single nodes that are not within this threshold could still be part of a clone class if the other cloned nodes are more towards the lower bound of the threshold.

1	doA(a); // {1, -2, 7, -3, -4}	1	doA(a); // {1, -2, 7, -3, -4}
2	doB(a); // {5, -2, 7, -3, -4}	2	doC(a); // {6, -2, 7, -3, -4}
3	doD(); // {8, -2, -3, -4}	3	doD(); // {8, -2, -3, -4}
4	doE(); // {9, -2, -3, -4}	4	doE(); // {9, -2, -3, -4}
5	doF(); // {10, -2, -3, -4}	5	doF(); // {10, -2, -3, -4}
6	doG(a); // {11, -2, 7, -3, -4}	6	doH(b); // {12, -2, 13, -3, -4}

Figure 4.14: Type 2R clones with their equality arrays.

As an example, consider the clone fragment in figure 4.14. In this example, we have 2 clone classes. In the “Try To Expand” step, we will check whether we can create a clone class with larger clone instances on basis of the clones found in the previous step. We start with the largest clone class, measured in clone volume. We measure clone volume as the product of the amount of clone instances in a clone class and the amount of nodes in any of its instances.

In the example in figure 4.14 the largest instance has a volume of:

$$3 \text{ nodes} * 2 \text{ clone instances} = 6 \quad (4.7)$$

Starting from this clone, we will try to expand it. We will include the previous node in the clone class and verify its variability threshold (node N_2). In this case, using the formulas shown in section 4.3.3.2, we can check whether this clone class would be valid by the variability threshold:

$$\frac{|\{(5,6)\}|}{|\{(5,6), (7,7), (8,8), (9,9), (10,10)\}|} * 100 = \frac{1}{5} * 100 = 20\% \quad (4.8)$$

Dependent on the actual variability threshold, the clone class would get expanded with node N_2 . For this example, we assume the variability threshold is set to 15%. In such a case, including this node would not result in a valid clone class. However, we continue trying to expand this clone class until we have reached the first node of the originally cloned fragment (the cloned fragment that did not take the variability threshold into account). Thus, we now try to expand with $\{N_1, N_2\}$. We then get the following formula:

$$\frac{|\{(5,6)\}|}{|\{(1,1), (7,7), (5,6), (7,7), (8,8), (9,9), (10,10)\}|} * 100 = \frac{1}{7} * 100 = 14\% \quad (4.9)$$

This falls under the threshold of 15% so we expand the clone class. Thus, the clone class $\{N_3, N_4, N_5\}$ will become $\{N_1, N_2, N_3, N_4, N_5\}$. We cannot expand further because we have reached the end of the original cloned fragment. When we are done expanding to previous nodes, we will try to expand to next nodes. In this case, we check whether we can include N_6 into the clone class. If N_6 would be included, the variability threshold would be as follows:

$$\frac{|\{(5,6), (11,12), (7,13)\}|}{|\{(1,1), (7,7), (5,6), (7,7), (8,8), (9,9), (10,10), (11,12), (7,13)\}|} * 100 = \frac{3}{9} * 100 = 33\% \quad (4.10)$$

This does not fall within the variability threshold, so we do not include this node in the clone class. After we have checked a single clone class within the bigger cloned fragment for expansion opportunities, we go on with the next clone class (by volume). In this case, the other clone class within this fragment has become a subset of the clone class we just expanded. Because of that, this other clone class is discarded. For the example code fragment in figure 4.14 the resulting clone class is $\{N_1, N_2, N_3, N_4, N_5\}$.

4.3.4 Checking for type 3 opportunities

As described in section 3.2.3, we define type 3R clones as gapped clones. This means that two existing clones may have a gap of non-gapped clones. We check for every found clone class whether it is a type 3R clone with another clone:

$$\forall(c_1 \in S) \exists(c_2 \in S) x \neq y \wedge isType3R(c_1, c_2) \quad (4.11)$$

Where S is the set of all found clone classes. $isType3R(C_1, C_2)$ is a Boolean-valued function that returns whether two clone classes are type 3R clones of each other. Two clone classes are type 3R clones of each other if each of their instances are within a certain distance of each other:

$$isType3R(C_1, C_2) = \forall(i_1 \in C_1) \exists(i_2 \in C_2) Fi_1 = Fi_2 \wedge gap(i_1, i_2) < gap_threshold \quad (4.12)$$

Where F is the file in which the clone instance is located. $gap_threshold$ is the defined threshold percentage of the maximum gap size between two clone instances. $gap(I_1, I_2)$ is the function that calculates the gap between two clone instances. This gap is calculated by dividing the amount of statements in the gap by the amount of statements that both clone instances span together:

$$gap(i_1, i_2) = \frac{|\{(Rn > Ri_1 \wedge Rn < Ri_2) \vee (Rn > Ri_2 \wedge Rn < Ri_1) \mid n \in nodes(Fi_1)\}|}{|i_1| + |i_2|} * 100 \quad (4.13)$$

Where F is the file in which the clone instance is located and R is the range of a clone instance or node. A range denotes the line and column at which a code fragment starts and ends. We define the partial order relationship on ranges as follows:

$$r_1 < r_2 \text{ iff } ELr_1 < BLr_2 \vee (ELr_1 = BLr_2 \wedge ECr_1 < BCr_2) \quad (4.14)$$

$$r_1 > r_2 \text{ iff } BLr_1 > ELr_2 \vee (BLr_1 = ELr_2 \wedge BCr_1 > ECr_2) \quad (4.15)$$

Where BL is the begin line of a range, EL is the end line of a range, BC is the begin column of a range and EC is the end column of a range.

4.3.4.1 Example

In this section we show an example of the identification of a type 3R clone.

1	<code>// File1.java</code>	1	<code>// File2.java</code>
2	<code>doA(); // C₁ I₁</code>	2	<code>doA(); // C₁ I₂</code>
3	<code>doB(); // C₁ I₁</code>	3	<code>doB(); // C₁ I₂</code>
4	<code>doC();</code>	4	<code>doC(); // C₂ I₂</code>
5	<code>doC(); // C₂ I₁</code>	5	<code>doD(); // C₂ I₂</code>
6	<code>doD(); // C₂ I₁</code>		

Figure 4.15: Gapped clones.

In the example of figure 4.15 we see two clone classes, which are separated by a single node. Using equation 4.13 we can calculate the size of the gap as follows:

$$gap(C_1 I_1, C_2 I_1) = \frac{|\{n_4\}|}{|\{n_2, n_3\}| + |\{n_5, n_6\}|} * 100 = \frac{1}{4} = 25\% \quad (4.16)$$

The same calculation is conducted for the other clone instances:

$$gap(C_1 I_2, C_2 I_2) = \frac{|\emptyset|}{|\{n_2, n_3\}| + |\{n_4, n_5\}|} * 100 = \frac{0}{4} = 0\% \quad (4.17)$$

If all resulting percentages are under the threshold, this clone will be added to the collection of found clones and all of its subsets will be removed.

4.4 Context Analysis of Clones

To be able to refactor code clones, CloneRefactor first maps the context of code clones. We define the following aspects of the clone as its context:

1. **Relation:** The relation of clone instances among each other through inheritance.
2. **Location:** Where a clone instance occurs in the code.
3. **Contents:** The statements/declarations of a clone instance.

We define categories for each of these aspects and enable CloneRefactor to determine the categories of clones.

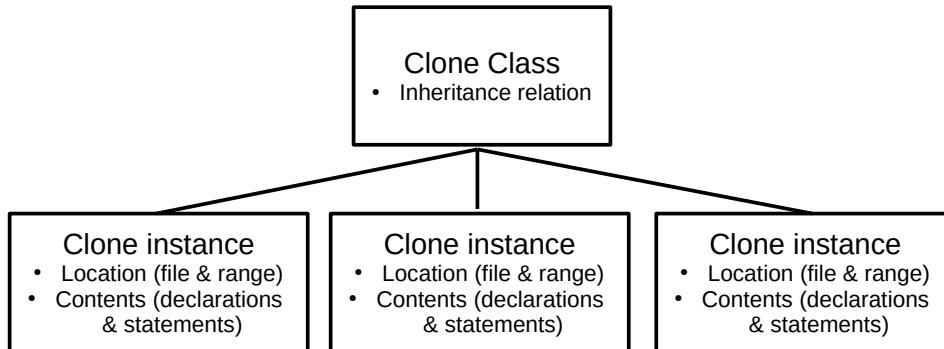


Figure 4.16: Abstract representation of clone classes and clone instances.

Fig. 4.16 shows an abstract representation of clone classes and clone instances. The relation of clones through inheritance is measured for each clone class. The location and contents of clones are measured for each clone instance.

4.4.1 Relation

When merging code clones in object-oriented languages, it is important to consider the inheritance relation between clone instances. This relation has a big impact on how a clone should be refactored.

Fontana et al. [16] describe measurements on 50 open source projects on the relation of clone instances to each other. To do this, they first define several categories for the relation between clone instances in object-oriented languages. These categories are as follows:

1. **Same Method:** All instances of the clone class are in the same method.
2. **Same Class:** All instances of the clone class are in the same class.
3. **Superclass:** All instances of the clone class are in a class that is child or parent of each other.
4. **Sibling Class:** All instances of the clone class have the same parent class.
5. **Ancestor Class:** All instances of the clone class are superclasses except for the direct superclass.
6. **First Cousin Class:** All instances of the clone class have the same grandparent class.
7. **Same Hierarchy Class:** All instances of the clone class belong to the same inheritance hierarchy, but do not belong to any of the other categories.
8. **Same External Superclass:** All instances of the clone class have the same superclass, but this superclass is not included in the project but part of a library.
9. **Unrelated class:** There is at least one instance in the clone class that is not in the same hierarchy.

We added the following categories, to gain more information about clones and reduce the number of unrelated clones:

1. **Same Interface:** All instances of the clone class are in a class or interface that have a common interface anywhere in their inheritance hierarchy.
2. **No Direct Superclass:** All instances of the clone class are in a class that does not have any superclass.
3. **No Indirect Superclass:** All instances of the clone class are in a class that does not have any external classes in its inheritance hierarchy.

4. **External Ancestor:** All instances of the clone class are in a class that does not have any external classes in its inheritance hierarchy.

Figure 4.16 shows an abstract representation of clone classes and clone instances. The relation of clones through inheritance is measured on clone class level: it involves all child clone instances. The location and contents of clones is measured on clone instance level. A clone’s location involves the file it resides in and the range it spans (for example: line 6 col 2 - line 7 col 50). A clone instance contents consists of a list of all statements and declarations it spans.

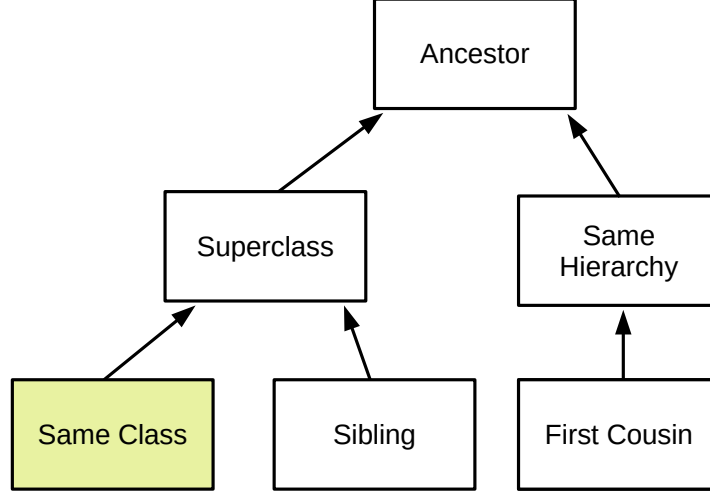


Figure 4.17: Abstract figure displaying relations of clone classes. Arrows represent superclass relations.

We separate these relations into the following categories, because of their related refactoring opportunities:

- **Common Class:** *Same Method, Same Class*
- **Common Hierarchy:** *Superclass, Sibling Class, Ancestor Class, First Cousin, Same Hierarchy*
- **Common Interface:** *Same Interface*
- **Unrelated:** *No Direct Superclass, No Indirect Superclass, External Superclass, External Ancestor*

Every clone class only has a single relation, which is the first relation from above list that the clone class applies to. For instance: all “Superclass” clones also apply to “Same Hierarchy”, but because “Superclass” is earlier in above list they will get the “Superclass” relation. This is because the items earlier in the list denote a more favorable refactoring.

To check a the relation for a given clone class, we compare each node in a clone instance with the respective nodes in each other clone instance within a clone class. So for a clone class with x clone instances and each clone instance has y nodes, we perform:

$$\frac{x * (x + 1)}{2} * y \quad (4.18)$$

comparisons to determine the relation. Each of these checks will result in one of the relations from the list above. Of these relations, the item lowest in the list is going to be the relation of the clone class. So for instance, if a clone class has 15 nodes that denote a *Superclass* relation but 3 nodes are *Unrelated*, the clone class becomes *Unrelated*.

4.4.1.1 Common Class

The *Same method* and *Same class* relations share a common refactoring opportunity. Clones of both these catagories, when extracted to a new method, can be placed in the same class. Both of these relations are most favorable for refactoring, as they require a minimal design tradeoff. Furthermore, global variables that are used in the class can be used without having to create method parameters.

Cloned nodes are flagged as *Same method* if these nodes are found in the same method. We define the method as the first method declaration that is encountered when following the parent nodes in the

ast for a given cloned node. Please note that a clone instance may not always be in a method, for which this predicate will fail.

Cloned nodes are flagged as *Same class* if these nodes are found in the same class. We define the class as the first class declaration that is encountered when following the parent nodes in the ast for a given cloned node.

4.4.1.2 Common Hierarchy

Clones that are in a common hierarchy can be refactored by using the “Extract Method” refactoring method followed by “Pull Up Method” until the method reaches a location that is accessible by all clone instances. However, the more often “Pull Up Method” has to be used, the more detrimental the effect is on system design. This is because putting a lot of functionality in classes higher up in an inheritance structure can result in the “God Object” anti-pattern. A god object is an object that knows too much or does too much.

Cloned nodes are flagged as *Superclass* if the classes in which these nodes are found are parent and child class of each other. Cloned nodes are flagged as *Siblings* if the classes in which these nodes are found all share the same parent and this parent class is not external. Cloned nodes are flagged as *First Cousin* if the classes in which these nodes are found all share the same grandparent and this grandparent class is not external.

Cloned nodes are flagged as *Ancestor* if the classes in which these nodes are found are recursively the parent of an antecedent. Cloned nodes are flagged as *Same Hierarchy* if the classes in which these nodes are found are all in the same inheritance hierarchy and not linked by any external classes.

4.4.1.3 Common Interface

Many object-oriented languages know the concept of “interfaces”, which are used to specify a behavior that classes must implement. As code clones describe functionality and interfaces originally did not allow for functionality, interfaces did not open up refactoring opportunities for duplicated code. However, many programming languages nowadays support default implementations in interfaces. Since Java 7 and C#8, these programming languages allow for functionality to be defined in interfaces. Many other object-oriented languages like Python allow this by nature, as they do not have a true notion of interfaces.

The greatest downside on system design of putting functionality in interfaces is that interfaces are per definition part of a classes’ public contract. That is, all functionality that is shared between classes via an interface cannot be hidden by setting a stricter visibility. Because of that, we favor all “Common Hierarchy” refactoring opportunities over “Common Interface”.

To check whether two cloned nodes have common interfaces, we recursively walk all types the class or interface of the cloned nodes implement and extend. The common interface that is closest to the cloned node in terms of depth of recursion will be flagged as the refactoring candidate.

4.4.1.4 Unrelated

Clones are unrelated if they share no common class or interface in their inheritance structure. These clones are least favorable when looking at refactoring, because their refactoring will almost always have a major impact on system design. We formulated four categories of unrelated clones to look into their refactoring opportunities.

Cloned nodes are flagged *No Direct Superclass* if they are in classes that do currently not have a parent. This marks the opportunity for creating a superclass abstraction and placing the extracted method there. Cloned nodes are flagged *No Indirect Superclass* if any of their ancestors does not have a parent. In such a case, it would be possible to create such an abstraction for the ancestor that does not have a parent.

Cloned nodes are flagged *External Superclass* if they are in a class which has an external parent. Cloned nodes are flagged *External Ancestor* if one of their ancestors has an external parent. Both of these relations obstruct the possibility of creating a superclass abstraction. In such a case, an interface abstraction could be created to make their relation explicit.

4.4.2 Location

A paper by Lozano et al. [24] discusses the harmfulness of cloning. The authors argue that 98% are produced at method-level. However, this claim is based on a small dataset and based on human copy-

paste behavior rather than static code analysis. We decided to measure the locations of clones through static analysis on our dataset. We chose the following categories:

1. **Method/Constructor Level:** A clone instance that does not exceed the boundaries of a single method or constructor (optionally including the declaration of the method or constructor itself).
2. **Class Level:** A clone instance in a class, that exceeds the boundaries of a single method or contains something else in the class (like field declarations, other methods, etc.).
3. **Interface Level:** A clone that is (a part of) an interface.
4. **Enumeration Level:** A clone that is (a part of) an enumeration.

We check the location of each clone instance for each of its nodes. If any node reports a different location from the others, we choose the location that is lowest in above list. So for instance, if a clone instance has 15 nodes that denote a *Method Level* location but 3 nodes are *Class Level*, the clone instance becomes *Class Level*.

4.4.2.1 Method/Constructor Level Clones

Method/Constructor Level clones denote clones that span are found in either a method or constructor. A constructor is a special method that is called when an object is instantiated. Current clone refactoring studies only focus on clones at method level [choi2011extracting, yue2018automatic, arcelli2013software, lin2014clonepedia, mandal2014automatic, balazinska2000advanced, yongting2018detect, bouktif2006novel, fanqi2014using, devi2016study, 23]. This is because most clones reside at those places [16, 24] and most of those clones can be refactored with a relatively simple set of refactoring techniques [16, 23].

To detect whether a given clone instance is at Method Level, we check for each node in the clone instance whether it has an ancestor node that is a method declaration:

a clone instance I is considered “Method Level” iff $\forall (n \in I) Tn = \text{“MethodDeclaration”} \vee \exists (p \in \text{ancestors}(n)) Tp = \text{“MethodDeclaration”}$

Where T is the type of a node. $\text{ancestors}(n)$ denotes all of ancestor nodes of a node. The ancestor nodes are recursively all the parents of a given node in the AST. Figure ?? shows an example of a method level clone and its recursive parents.

<pre> 1 public class Class1 { // p3 2 public void doA() { // p2 3 if(isA()) // p1 4 doB(); // n1 5 } 6 }</pre>	<pre> 1 \begin{javacode} 2 public class Class2 { // p3 3 public void doB() { // p2 4 if(isB()) // p1 5 doB(); // n2 6 } 7 }</pre>
--	---

Figure 4.18: Example of a method level clone and its parents.

4.4.3 Contents

Finally, we looked at what nodes individual clone instances span. We selected the following categories to be relevant for refactoring:

1. **Full Method/Class/Interface/Enumeration:** A clone that spans a full class, method, constructor, interface or enumeration, including its declaration.
2. **Partial Method/Constructor:** A clone that spans a method partially, optionally including its declaration.
3. **Several Methods:** A clone that spans over two or more methods, either fully or partially, but does not span anything but methods (so not fields or anything in between).
4. **Only Fields:** A clone that spans only global variables.

5. **Includes Fields/Constructor:** A clone that spans a combination of fields and other things, like methods.
6. **Method/Class/Interface/Enumeration Declaration:** A clone that contains the declaration (usually the first line) of a class, method, interface or enumeration.
7. **Other:** Anything that does not match with above-stated categories.

Explain
each cat-
egory in
more de-
tail

4.5 Method extraction opportunities

The most used technique to refactor clones is method extraction (creating a new method on basis of the contents of clones). However, method extraction cannot be applied in all cases. In these instances, more conditions may apply to be able to conduct a refactoring, if beneficial at all.

We measured the number of clones that can be refactored through method extraction (without additional transformations being required). We defined the following categories:

- **Can be extracted:** This clone is a fragment of code that can directly be extracted to a method. Then, based on the relation between the clone instances, further refactoring techniques can be used to refactor the extracted methods (for instance “pull up method” for clones in sibling classes).
- **Complex control flow:** This clone contains `break`, `continue` or `return` statements, obstructing the possibility of method extraction.
- **Spans part of a block:** This clone spans a part of a statement.
- **Is not a partial method:** If the clone does not fall in the “Partial method” category of Table 5.4, the “extract method” refactoring technique cannot be applied.

Explain
each cat-
egory in
more de-
tail

Chapter 5

Results

In this chapter, we present the results of our experiments.

5.1 The corpus

For our experiments we use a large corpus of open source projects [43]. This corpus has been assembled to contain relatively higher quality projects. Also, any duplicate projects were removed from this corpus. This results in a variety of Java projects that reflect the quality of average open source Java systems and are useful to perform measurements on.

As indicated in chapter 3.2.5 CloneRefactor requires all libraries of software projects we test. As these are not included in the used corpus [43], we decided to filter the corpus to only include Maven projects. Maven is a build automation tool used primarily for Java, and works on basis of an `pom.xml` file to describe the projects' dependencies. As no `pom.xml` files are included in the corpus, we cloned the latest version of each project in the corpus. We then removed each project that has no `pom.xml` file. As a final step, we collected all dependencies for each project by using the `mvn dependency:copy-dependencies -DoutputDirectory=lib` Maven command, and removed each project for which not all dependencies were available (due to non-Maven dependencies being used or unsatisfiable dependencies being referenced in the `pom.xml` file).

Some general data regarding this corpus is displayed in Table 5.1.

Table 5.1: General results for GitHub Java projects corpus [43].

Amount of projects	1,361
Amount of lines (excluding whitespace, comments and newlines.)	1,414,996
Amount of statements/declarations	1,212,189
Amount of tokens (excluding whitespace, comments and newlines.)	11,643,194

5.2 Clone types

In this section we display the differences between clone type 1-3 [12] and type 1R-3R as proposed in chapter 3. When running our clone detection script over the corpus, we get the results displayed in Fig. 5.1.

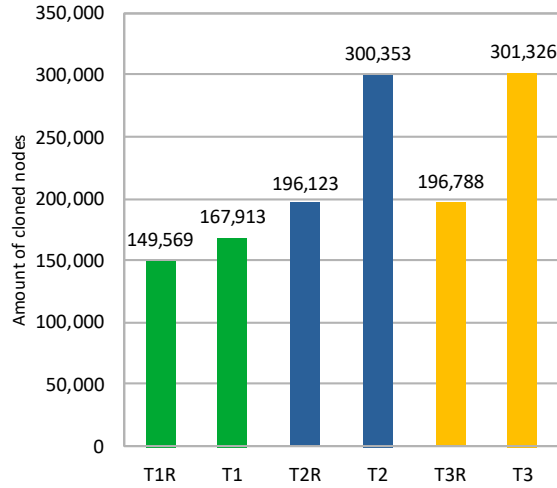


Figure 5.1: Number of cloned declarations/statements.

In this figure, the number of cloned nodes per clone type are displayed. The difference between T1R and T1 is small (10.9%), because most often textually equal code is also functionally equal. The difference between T2R and T2 is bigger (34.7%) because the T2R definition is more strict. T3R and T3 are similar to T2R and T2 because our dataset does not have so many gapped clones for the thresholds used.

We also measured the duration of finding clones by the different clone types. Figure 5.2 shows the duration of detecting all clones in the corpus using CloneRefactor for different clone types. Although this data is partly dependent on our implementation of the clone types, there is a notable difference between the refactoring-oriented clone types and the literature clone types. The reason for this is further explained in section 3.2.5.

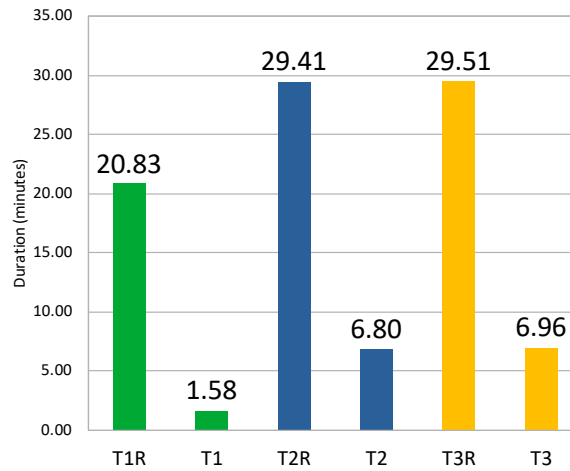


Figure 5.2: Duration in minutes of identifying clones for clone type definitions.

5.3 Thresholds

todo

5.4 Context Analysis of Clones

We analyzed the context of clones in a large corpus of open source projects. For these experiments, we used our CloneRefactor tool. These experiments follow the structure of the context: The relation between clone instances is explained, measured and discussed in chapter 5.4.1; the location of clone instances is

explained, measured and discussed in chapter 5.4.2; the content of clone instances are explained, measured and discussed in chapter 5.5.

5.4.1 Relations Between Clone Instances

In section 4.4.1 we introduced our experiments regarding relations between clones. Table 5.2 contains our results regarding the relations between clone instances.

update
these re-
sults

Table 5.2: Clone relations

Relation	#	%
Unrelated	12,368	34.88
Same Class	11,483	32.38
Same Method	5,056	14.26
Sibling	4,182	11.79
External Superclass	1,066	3.01
Superclass	558	1.57
First Cousin	489	1.38
Common Hierarchy	206	0.58
Ancestor	54	0.15

The most notable difference when comparing it to the results of Fontana et al. [16] is that in our results most of the clones are unrelated (34.44%), while for them it was only 15.70%. This is likely due to the fact that we consider clone classes rather than clone pairs, and mark the clone class “Unrelated” even if just one of the clone instances is outside a hierarchy. It could also be that the corpus which we use, as it has generally smaller projects, uses more classes from outside the project (which are marked “Unrelated” if they do not have a common external superclass). About a third of all clone classes have all instances in the same class, which is generally easy to refactor. On the third place come the clones that are in the same method, which are similarly easy to refactor.

5.4.2 Clone instance location

Measuring the clone locations categories defined in section 4.4.2 yields the results displayed in Table 5.3.

update
these re-
sults

Table 5.3: Clone instance locations

Location	#	%
Method Level	83,813	82.62
Class Level	12,534	12.35
Constructor Level	4,391	4.33
Interface Level	567	0.56
Enum Level	145	0.14

Our results indicate that around 58% of the clones are produced at method-level. About 39% of clones either span several methods/constructors or contain something like a field declaration. Another 3% of the clones are found in constructors. The amount of clones found in interfaces and enumerations is very low. Regarding the differences between type 1 and type 1R, it seems that there are relatively less method level clones and more class level clones for type 1R. This is probably due to that the main reason for variability between type 1 and type 1R is variable references, which occur more at method level than class level.

5.5 Clone instance contents

Measuring the clone contents categories defined in section 4.4.3 yields the results displayed in Table 5.4.

update
these re-
sults

Table 5.4: Clone instance contents

Contents	#	%
Partial Method	79,945	78.80
Several Methods	7,323	7.22
Partial Constructor	4,385	4.32
Full Method	3,868	3.81
Only Fields	3,039	3.00
Includes Constructor	1,522	1.50
Includes Field	616	0.61
Includes Class Declaration	376	0.37
Other Categories	376	0.37

Unsurprisingly, most clones span a part of a method. The most used refactoring technique for clones that span part of a method is “Extract Method”. Because of that, we focus our research efforts on refactoring such clones.

5.6 Merging duplicate code through refactoring

In section 4.5 we define categories regarding the extent to which clones can be refactored through method extraction. Our results are displayed in Table 5.5.

update
these re-
sults

Table 5.5: Refactorability through method extraction

	#	%
Can be extracted	14,664	41.35
Spans part of a block	10,801	30.46
Is not a partial method	8,074	22.77
Complex control flow	1,922	5.42

From Table 5.5, we can see that 41% of the clones can directly be refactored through method extraction (and possibly other refactoring techniques based on the relation of the clone instances). For the other clones, other techniques or transformations will be required.

5.7 Refactoring clones

TODO

Chapter 6

Discussion

TODO

In this chapter, we discuss the results of our experiment(s) on ...

Finding 1: Highlight like this an important finding of your analysis of the results.

Refer to Finding 1.

6.1 Clone Type Definitions

In this study, we proposed a set of clone type definitions for which a refactoring opportunity is known. We based these clone type definitions on the clone type definitions that are commonly used in literature. The design of these clone type definitions entails some decisions that have a large impact on the clones we considered for refactoring in this study. In this section we discuss our clone type definitions as proposed in section 3.2.

6.1.0.1 Type 2R clones

With type 2R clones we allow variability in some identifiers and literals such that the code can and should still be refactored. For type 2R clones we chose a set of expressions in which we allow variability and proposed a recommended refactoring strategy. We think however that type 2R could still use a lot of improvement to find more duplication patterns that can be refactored.

One method we think can be used to find more refactoring opportunities is to allow variability in expressions that have/return the same type. If expressions have/return the same type, they can be extracted to a parameter and the corresponding expression can be passed as a parameter. An example of this is displayed in figure 6.1. The only thing to watch out for is method that have side effects. Because methods may be executed in another point during execution, this might affect the functionality of the code.

<pre> 1 // Original 2 public void doStuff(){ 3 int numbers = 456; 4 doA(getTitle()); 5 doB(123); 6 doC(); 7 doA("456"); 8 doB(numbers); 9 } 10 11 public String getTitle(){ 12 return "123"; 13 } </pre>	<pre> 1 // Refactored 2 public void doStuff(){ 3 int numbers = 456; 4 doAandB(getTitle(), 123); 5 doC(); 6 doAandB("456", numbers); 7 } 8 9 public void doAandB(String var1, 10 int var2){ 11 doA(var1); 12 doB(var2); 13 } 14 15 public String getTitle(){ 16 return "123"; 17 } </pre>
--	--

Figure 6.1: Refactoring different expressions that have the same return type.

6.2 Clone Context Analysis

In section 4.4 we introduced the categories we defined for mapping the context of clones. In this section, we discuss this together with the related experiments.

6.2.1 Refactorability

In section 5.6 we introduced to what extent clones can be refactored through method extraction. Because we strived to get results fast, we excluded categories that could not not be directly refactored through method extraction. However, with a few transformations or further considerations it might be possible to make these clones refactorable. In this section we will highlight a few of these categories which we believe to be refactorable through method extraction with a bit more effort.

6.2.1.1 Partial block

We did not consider clones for refactoring that span a part of a block. Although it is indeed not possible to refactor such clones, there are possibilities to make such clones refactorable. For instance, if the programming language supports lambda expressions, we can move the difference of statements in the block in a lambda expression. Figure 6.2 shows an example of such a refactoring opportunity.

<pre> 1 // Original 2 public void doStuff() { 3 for(int i = 0; i<5; i++) { //Only 4 the declaration of this for 5 loop is cloned, but the loop 6 body is not. 7 System.out.println("hello!"); 8 } 9 } </pre>	<pre> 1 // Refactored 2 public void doStuff() { 3 doFiveTimes(() -> System.out. 4 println("hello!")); 5 doFiveTimes(() -> CoreController. 6 activateCore(i)); 7 } 8 9 public void doFiveTimes(Runnable 10 runnable){ 11 for(int i = 0; i<5; i++) { //Only 12 the declaration of this for 13 loop is cloned, but the loop 14 body is not. 15 runnable.run(); 16 } 17 } </pre>
---	---

Figure 6.2: Refactoring a method that is obstructed by a complex control flow.

6.2.1.2 Overlap within a clone class

We do not consider clone classes for refactoring that overlap in the clone class itself. That is because such clones might require other techniques to refactor them. We consider analyzing such refactoring methods future work. An example refactoring of clone instances that overlap within their clone class is displayed in figure 6.3. In this example one clone instance is found between line 3 and 6 and the other between line 4 and 7. To refactor such clones, other techniques may be required, like introducing a new for-loop to reduce the duplication.

<pre> 1 // Original 2 public void sayHello() { 3 System.out.println("hello!"); 4 System.out.println("hello!"); 5 System.out.println("hello!"); 6 System.out.println("hello!"); 7 System.out.println("hello!"); 8 } </pre>	<pre> 1 // Refactored 2 public void sayHello() { 3 for(int i = 0; i<5; i++){ 4 System.out.println("hello!"); 5 } 6 } </pre>
---	--

Figure 6.3: Refactoring a clone class that has overlap between clone instances.

6.2.1.3 Complex control flow

Break, continue and return statements can obstruct the possibility of performing method extraction. However, with some extra transformations, method extraction will still be possible in such cases. Figure 6.4 shows such a transformation. We can wrap the newly extracted method in a conditional to indicate whether the “control flow modifying statement” should be executed. In other cases, other methods might apply to refactor such clones.

```
1 | // Original
2 | public boolean doStuff() {
3 |     if(doA());
4 |         return false;
5 |     doB();
6 |     doC();
7 |     if(doA());
8 |         return false;
9 |     doB();
10 |     return true;
11 | }
```

```
1 | // Refactored
2 | public boolean doStuff() {
3 |     if(!doAandB())
4 |         return false;
5 |     doC();
6 |     return doAandB();
7 | }
8 |
9 | public boolean doAandB() {
10 |     if(doA())
11 |         return false;
12 |     doB();
13 |     return true;
14 | }
```

Figure 6.4: Refactoring a method that is obstructed by a complex control flow.

Chapter 7

Related work

The Duplicated Code Refactoring Advisor (DCRA) looks into refactoring opportunities for clone pairs [16, 42]. DCRA only focuses on refactoring clone pairs, with the authors arguing that *clone pairs are much easier to manage when considered singularly*. As intermediate steps, the authors measure a corpus of Java systems for some clone-related properties of the systems, like the relation (in terms of inheritance) between code fragments in a clone pair. We further look into these measurements in Sec. 5.4.1.

A tool named Aries [44, 45] focuses on the detection of refactorable clones. They do this based on the relation between clone instances through inheritance, similar to Fontana et al. [42]. This tool only proposes a refactoring opportunity and does not provide help in the process of applying the refactoring.

We investigated several research efforts that look into code clone refactoring [46–48]. However, all of these tools only support a subset of all harmful clones that are found. Also, these tools are limited to suggesting refactoring opportunities, rather than actually applying refactorings where suitable. Finally, all published approaches have limitations, such as false positives in their clone detection [47] or being limited to clone pairs [45].

Rewrite
this part
and add
a lot of
content

Chapter 8

Conclusion

In the research we have conducted so far we have made three novel contributions:

- We proposed a method with which we can detect clones that can/should be refactored.
- We mapped the context of clones in a large corpus of open source systems.
- We mapped the opportunities to perform method extraction on clones this corpus.

We have looked into existing definitions for different types of clones [12] and proposed solutions for problems that these types have with regards to automated refactoring. We propose that fully qualified identifiers of method call signatures and type references should be considered instead of their plain text representation, to ensure refactorability. Furthermore, we propose that one should define thresholds for variability in variables, literals and method calls, in order to limit the number of parameters that the merged unit shall have.

The research that we have conducted so far analyzes the context of different kinds of clones and prioritizes their refactoring. Firstly, we looked at the inheritance relation of clone instances in a clone class. We have found that more than a third of all clone classes are flagged unrelated, which means that they have at least one instance that has no relation through inheritance with the other instances. For about a fourth of the clone classes all of its instances are in the same class. About a sixth of the clone classes have clone instances that are siblings of each other (share the same superclass).

Secondly, we looked at the location of clone instances. Most clone instances (58 percent) are found at method level. About 37 percent of clone instances were found at class level. We defined “class level clones” as clones that exceed the boundaries of a single method or contain something else in the class (like field declarations, other methods, etc.). Thirdly, we looked at the contents of clone instances. Most clones span a part of a method (57 percent). About 26 percent of clones span over several methods.

We also looked into the refactorability of clones that span a part of a method. Over 10 percent of the clones can directly be refactored by extracting them to a new method (and calling the method at all usages using their relation). The main reason that most clones that span a part of a method cannot directly be refactored by method extraction, is that they contain `return`, `break` or `continue` statements.

8.1 Threats to validity

We noticed that, when doing measurements on a corpus of this size, the thresholds that we use for the clone detection have a big impact on the results. There does not seem to be one golden set of thresholds, some thresholds work in some situations but fail in others. We have chosen thresholds that, according to our manual assessment, seemed optimal. However, by using these, we definitely miss some harmful clones.

8.2 Future work

This study presents a foundation for research in a largely unexplored field of studies: analyzing maintainability through automated refactoring. However, we scratched just the tip of the iceberg regarding all research opportunities in this field. In this section we describe possible extensions to this research, as well as other research opportunities in this field of studies.

Move part
of future
work to
discussion

8.2.1 Automated Refactoring for more metrics

In this study we presented evidence regarding the value of applying automated refactoring to analyze the before- and after state of source code and refactored source code. Analyzing these states we were able to analyze the improvement in maintainability after applying certain refactorings. This allows us to better assess thresholds by which maintainability issues in source code are identified. We also get better insights in the costs and values of applying certain refactoring efforts.

For this study we chose to focus only on the automated refactoring of duplication in source code. However, software maintainability depends on more factors and can be measured by more metrics. These factors also have opportunities to automate their refactorings. We think that several similarly sized studies can be conducted to automate the refactoring of other maintainability metrics.

A study by Heitlager et al. [10] presents several metrics by which the maintainability of source code can be assessed. They propose thresholds that indicate issues with these metrics. Many of these metrics have automated refactoring opportunities. In this section we will focus on several of these metrics to outline their opportunities for automated refactoring.

8.2.1.1 Long parameter list

When multiple parameters are used together in a method, there is an implicit dependency between these parameters: the dependency of being required by that method. *If a lot of data hangs around really tight together, they should be made into their own object* [11, 49]. Guidelines describe to limit the number of parameters per method to at most 4 [49].

If a method has many parameters, we can group strongly related parameters into an object. This can be done automatically, but two things must be considered:

- How do we determine whether parameters are strongly related?
- At what other places is this data used in unison and should thus use the new abstraction?

To determine whether parameters are strongly related we must look into at what places they are used in the codebase. We must then define some threshold that denotes the percentage of usages of these variables in which they are used together. If this threshold exceeds a certain amount, we can group them into an object. We must then trace all places in which they are used together and replace the variables by the newly created abstractions.

8.2.1.2 Method complexity

Method complexity refers to the complexity of the logic in a method body. There are several methods to compute method complexity. The most used complexity metric is (McCabe) Cyclomatic Complexity [49], which refers to the amount of independent paths that can be taken through the source code. Another complexity metric that has recently become fairly popular is Cognitive Complexity [50], which attempts to measure the human perceived complexity. Both indicate an aspect of source code maintainability.

Dealing with method complexity can largely be done by method extraction. We extract a part of a complex method to a new method. This way we split the complexity of the original method into separately testable methods. Also, the methods become easier to read.

Refactoring complex methods can largely be done automatically. Also, many of the results of this study can be reused. To assess an automated refactoring opportunity for complex methods, we should assess which parts of methods can be split to end up with parts of similar complexity. For this, our research can be used to assess whether a given piece of code can be extracted to a new method (see section 5.6).

We recommend to extend our tool, CloneRefactor, to allow for such capabilities. CloneRefactor already contains a component that calculates the Cyclomatic Complexity of a given method. Using our automated refactored model, identified problems can relatively easily be refactored.

8.2.1.3 Method size

Method size has a strong relation, in terms of refactoring, with Cyclomatic Complexity. Although method size is often related to Cyclomatic Complexity, a study by Landman et al. [51] shows that Cyclomatic Complexity is not redundant to method size. However, they do share a similar method of refactoring. The automated refactoring opportunities described in the previous section also apply to method size.

8.2.1.4 Combining the metrics

Combining the automated refactoring models for each of these metrics can result in a model that ultimately provides significant improvement in the ease of writing well-maintainable source code. In this study, we presented a few tradeoffs that are the result of refactoring code clones. For instance, refactoring code clones with a lot of variability can result in long parameter lists. However, if we can combine this with automated refactoring of long parameter lists, this tradeoff can be mitigated. This way, it is possible to work towards an model of automated refactoring that reduces manual refactoring efforts significantly.

When programming, there are often trade-offs between technical debt and velocity. When a deadline comes near, often software quality is sacrificed to gain velocity [52–54]. Apart from that, a low programmer aptitude can result in low quality code [55]. This is because often forming appropriate abstractions requires time, effort and critical thinking. By introducing these abstractions automatically, this negative impact can (partly) be mitigated.

8.2.2 Type 4 clones

In section 2.2 we introduced the clone types that are used in literature. Of these clone types, we considered type 1, 2 and 3 for refactoring. We chose not to consider type 4 clones because they appear far less in source code and are more difficult to detect and refactor. However, that does not mean that these clones should not be refactored.

We think that there are relevant research opportunities in refactoring type 4 clones. Type 4 clones are functionally equal pieces of code that are implemented through different implementations. Although functionally equal, type 4 clones might still differ in other aspects. For instance, for type 4 clones, one alternative might be easier to maintain. Also, they can differ in performance.

There are interesting research opportunities in automatically choosing the better alternative among type 4 clone instances.

8.2.2.1 Proposing new clone type definitions

We think that current definitions of clone types still lack in the identification of all duplication issues that a developer should invest his/her refactoring efforts in. Current clone detection techniques still result in many false positives and false negatives. By proposing good definitions for code clones that mark the characteristics of harmful anti-patterns, we can create more accurate suggestions to developers.

8.2.3 Naming of refactored methods/classes/etc.

In this study, we took naming of refactored entities out of the scope. We applied automated refactoring to duplication, and gave all generated methods, classes and interfaces generated names. For our purposes that was not much of an issue, because of the validation methods used. We validated our approach using the SIG Maintainability model, which does not take naming quality into account. Because of that, the results of our experiments are in no way dependent on the names we give the generated methods, classes and interfaces.

When using our work with the purpose of refactoring assistance, these names will have to be manually provided. However, recent studies allow assistance in this process by generating a name that matches the body of a declaration. If this can be done in a reliable way, we can apply refactorings without any manual steps required.

A study by Allamanis et al. [56] proposes a machine learning model that can suggest accurate method and class names. This study shows promising results towards generating method and class names on basis of their body and context. However, the source code of this study is not available, making it harder to apply their findings to generate class and method names for any software project.

In a recent study by Alon et al. [57] they propose code2seq. Code2seq is a machine learning model that guesses the name of a method given a method body. This model has been trained on a large set of method bodies and names. The model already shows promising results. The source code of code2seq is publicly available, making it possible to embed this model in any application. Although this model is still far from perfect, combining it with our research could already greatly reduce the manual refactoring effort required. The main deficiency of this model lies in that its limited to the method body and does currently not consider anything else in a method body. This method is flawed, because the meaning of method bodies depends heavily on their context. For instance, the body of the method `get()` in `ArrayList` does

not make any sense in relation to its name without considering its context (class, inheritance hierarchy, etc.).

8.2.4 Looking into other languages/paradigms

In this study we describe duplicate code refactoring opportunities for object-oriented languages. We built a tool to refactor code clones in Java and used it to run our experiments.

Applying our experiments with other programming languages than just Java might result in valuable results. Refactoring opportunities are greatly dependent on coding conventions, which differ per language. Other languages might result in different results, which might result in different insights regarding the resolution of duplication problems found. We prioritized refactoring efforts based on Java, which might differ from the prioritization that can result from running our experiments with other programming languages.

In this study we focused only on the object-oriented paradigm because of the shared concepts. However, the problems that come with duplication in source code also appears in different programming languages. Because of that, more insights could be obtained when looking into automated refactoring opportunities for other paradigms. For instance, it might be valuable to look into the opportunities to reduce duplication in the functional domain. Dealing with code clones in the functional paradigm is pretty much an unexplored field and might hold valuable results.

8.2.5 Running our experiments with a different set of software projects

In our experiments we chose to use a GitHub based software corpus. This corpus contains a diversity of projects of many different sizes. We noticed that there is a lot of difference in the quality of these software systems: some systems had a lot of duplication and some did not have any duplication.

We think that there is value in running our experiments with different corpuses. We would, for instance, be interested in the results for industrial projects instead of open source software systems. We are curious to see whether the distributions are comparable, or whether they show large differences.

Furthermore, we think that there is value in running our experiments with a set of larger sized and more professionally developed open source systems. We noticed that in our corpus there were a lot of projects that have only few commits and contributors. We think it would be valuable to for instance perform our experiments with a set of larger open source systems, like the systems in the Apache ecosystem.

8.2.6 Automatically refactoring code that is duplicated with a library

One of the clone detection tools we analyzed in this study is Sisyphus [58]. This clone detection tool searches a codebase for methods that are part of the Java standard libraries. There are definite automated refactoring opportunities here, automatically switching to the library implementation if available.

This doesn't have to be limited to just the Java standard library. Although using a new library for a single cloned library, we can also search the libraries the project already uses. CloneRefactor is an excellent tool to expand for finding such clones with external libraries and refactoring the source code to use the library implementation if available. This is because CloneRefactor has scripts to gather all dependent libraries for Maven projects. It loads these libraries in order to be able to resolve symbols. If this would be expanded not only to resolve symbols but also to consider the clones between the analyzed project and its dependencies, we could find such refactoring opportunities.

Acknowledgements

Throughout the writing of this dissertation I have received a great deal of support and assistance. I would first like to thank my supervisor, Dr. A.M. Oprescu , whose expertise was invaluable in the academic writing and mathematical aspects of this thesis.

Next, I would like to thank my peer student Sander Meester for always being there when I was stuck. Whenever I needed to test some ideas you would always provide realistic insights. The long conversations about the problems I have been dealing with in my thesis would not have been solved in the same way if it wasn't for you. Your proofreading efforts helped take my writing to the next level.

Finally, I would like to acknowledge my colleagues from my internship at SIG for their wonderful collaboration. You supported me greatly and were always willing to help me.

Bibliography

- [1] M. Fowler, *Refactoring: improving the design of existing code*, Second. Addison-Wesley Professional, 2018.
- [2] W. C. Wake, *Refactoring workbook*. Addison-Wesley Professional, 2004.
- [3] B. P. Lientz, E. B. Swanson, and G. E. Tompkins, “Characteristics of application software maintenance”, *Communications of the ACM*, vol. 21, no. 6, pp. 466–471, 1978.
- [4] T. Mens and T. Tourwé, “A survey of software refactoring”, *IEEE Transactions on software engineering*, vol. 30, no. 2, pp. 126–139, 2004.
- [5] T. Mens, A. Van Deursen, *et al.*, “Refactoring: Emerging trends and open problems”, in *Proceedings First International Workshop on REFactoring: Achievements, Challenges, Effects (REFACE)*, 2003.
- [6] J. Ostberg and S. Wagner, “On automatically collectable metrics for software maintainability evaluation”, in *2014 Joint Conference of the International Workshop on Software Measurement and the International Conference on Software Process and Product Measurement*, Oct. 2014, pp. 32–37. DOI: 10.1109/IWSM.Mensura.2014.19.
- [7] M. Bruntink, A. Van Deursen, R. Van Engelen, and T. Tourwe, “On the use of clone detection for identifying crosscutting concern code”, *IEEE Transactions on Software Engineering*, vol. 31, no. 10, pp. 804–818, 2005.
- [8] T. Bakota, P. Hegedűs, G. Ladányi, P. Körtvélyesi, R. Ferenc, and T. Gyimóthy, “A cost model based on software maintainability”, in *2012 28th IEEE International Conference on Software Maintenance (ICSM)*, IEEE, 2012, pp. 316–325.
- [9] J. B. Munson, “Software maintainability: A practical concern for life-cycle costs”, in *The IEEE Computer Society’s Second International Computer Software and Applications Conference, 1978. COMPSAC’78.*, IEEE, 1978, pp. 54–59.
- [10] I. Heitlager, T. Kuipers, and J. Visser, “A practical model for measuring maintainability”, in *6th international conference on the quality of information and communications technology (QUATIC 2007)*, IEEE, 2007, pp. 30–39.
- [11] M. Fowler, *Refactoring: improving the design of existing code*. Addison-Wesley Professional, 1999.
- [12] C. K. Roy and J. R. Cordy, “A survey on software clone detection research”, *Queen’s School of Computing TR*, vol. 541, no. 115, pp. 64–68, 2007.
- [13] H. Sajnani, V. Saini, J. Svajlenko, C. K. Roy, and C. V. Lopes, “Sourcerercc: Scaling code clone detection to big-code”, in *2016 IEEE/ACM 38th International Conference on Software Engineering (ICSE)*, IEEE, 2016, pp. 1157–1168.
- [14] J. Svajlenko and C. K. Roy, “Bigcloneeval: A clone detection tool evaluation framework with bigclonebench”, in *2016 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, IEEE, 2016, pp. 596–600.
- [15] G. Bavota, B. De Carluccio, A. De Lucia, M. Di Penta, R. Oliveto, and O. Strollo, “When does a refactoring induce bugs? an empirical study”, in *2012 IEEE 12th International Working Conference on Source Code Analysis and Manipulation*, IEEE, 2012, pp. 104–113.
- [16] F. A. Fontana and M. Zanoni, “A duplicated code refactoring advisor”, in *International Conference on Agile Software Development*, Springer, 2015, pp. 3–14.
- [17] C. Kapser and M. W. Godfrey, ““cloning considered harmful” considered harmful”, in *Reverse Engineering, 2006. WCRE’06. 13th Working Conference on*, Citeseer, 2006, pp. 19–28.

- [18] L. Jiang, G. Misherghi, Z. Su, and S. Glondou, “Deckard: Scalable and accurate tree-based detection of code clones”, in *Proceedings of the 29th international conference on Software Engineering*, IEEE Computer Society, 2007, pp. 96–105.
- [19] S. Haeffliger, G. Von Krogh, and S. Spaeth, “Code reuse in open source software”, *Management science*, vol. 54, no. 1, pp. 180–193, 2008.
- [20] I. D. Baxter, A. Yahin, L. Moura, M. Sant’Anna, and L. Bier, “Clone detection using abstract syntax trees”, in *Proceedings. International Conference on Software Maintenance (Cat. No. 98CB36272)*, IEEE, 1998, pp. 368–377.
- [21] E. Kodhai, S. Kanmani, A. Kamatchi, R. Radhika, and B. V. Saranya, “Detection of type-1 and type-2 code clones using textual analysis and metrics”, in *2010 International Conference on Recent Trends in Information, Telecommunication and Computing*, IEEE, 2010, pp. 241–243.
- [22] B. van Bladel and S. Demeyer, “A novel approach for detecting type-iv clones in test code”, in *2019 IEEE 13th International Workshop on Software Clones (IWSC)*, IEEE, 2019, pp. 8–12.
- [23] E. Kodhai and S. Kanmani, “Method-level code clone modification using refactoring techniques for clone maintenance”, *Advanced Computing*, vol. 4, no. 2, p. 7, 2013.
- [24] A. Lozano, M. Wermelinger, and B. Nuseibeh, “Evaluating the harmfulness of cloning: A change based experiment”, in *Fourth International Workshop on Mining Software Repositories (MSR’07: ICSE Workshops 2007)*, IEEE, 2007, pp. 18–18.
- [25] M. White, M. Tufano, C. Vendome, and D. Poshyvanyk, “Deep learning code fragments for code clone detection”, in *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, ACM, 2016, pp. 87–98.
- [26] and L. Bergman, T. Lau, and D. Notkin, “An ethnographic study of copy and paste programming practices in oopl”, in *Proceedings. 2004 International Symposium on Empirical Software Engineering, 2004. ISESE ’04.*, Aug. 2004, pp. 83–92. DOI: 10.1109/ISESE.2004.1334896.
- [27] M. Mohnen, “Interfaces with default implementations in java”, in *Proceedings of the inaugural conference on the Principles and Practice of programming, 2002 and Proceedings of the second workshop on Intermediate representation engineering for virtual machines, 2002*, National University of Ireland, 2002, pp. 35–40.
- [28] V. Saini, F. Farmahinifarahani, Y. Lu, D. Yang, P. Martins, H. Sajjani, P. Baldi, and C. Lopes, “Towards automating precision studies of clone detectors”, *arXiv preprint arXiv:1812.05195*, 2018.
- [29] S. Livieri, Y. Higo, M. Matushita, and K. Inoue, “Very-large scale code clone analysis and visualization of open source programs using distributed ccfinder: D-ccfinder”, in *29th International Conference on Software Engineering (ICSE’07)*, IEEE, 2007, pp. 106–115.
- [30] T. Kamiya, S. Kusumoto, and K. Inoue, “Ccfinder: A multilingual token-based code clone detection system for large scale source code”, *IEEE Transactions on Software Engineering*, vol. 28, no. 7, pp. 654–670, 2002.
- [31] Y. Semura, N. Yoshida, E. Choi, and K. Inoue, “Ccfindersw: Clone detection tool with flexible multilingual tokenization”, in *2017 24th Asia-Pacific Software Engineering Conference (APSEC)*, IEEE, 2017, pp. 654–659.
- [32] C. K. Roy and J. R. Cordy, “Nicad: Accurate detection of near-miss intentional clones using flexible pretty-printing and code normalization”, in *2008 16th IEEE international conference on program comprehension*, IEEE, 2008, pp. 172–181.
- [33] J. Svajlenko and C. K. Roy, “Evaluating modern clone detection tools”, in *2014 IEEE International Conference on Software Maintenance and Evolution*, IEEE, 2014, pp. 321–330.
- [34] C. Ragkhitwetsagul and J. Krinke, “Siamese: Scalable and incremental code clone search via multiple code representations”, *Empirical Software Engineering*, pp. 1–49, 2019.
- [35] T. Lavoie and E. Merlo, “Automated type-3 clone oracle using levenshtein metric”, in *Proceedings of the 5th international workshop on software clones*, ACM, 2011, pp. 34–40.
- [36] C. Kamalpuriya and P. Singh, “Enhancing program dependency graph based clone detection using approximate subgraph matching”, in *2017 IEEE 11th International Workshop on Software Clones (IWSC)*, IEEE, 2017, pp. 1–7.

- [37] C. K. Roy, J. R. Cordy, and R. Koschke, “Comparison and evaluation of code clone detection techniques and tools: A qualitative approach”, *Science of computer programming*, vol. 74, no. 7, pp. 470–495, 2009.
- [38] J. R. Cordy and C. K. Roy, “The nicad clone detector”, in *2011 IEEE 19th International Conference on Program Comprehension*, IEEE, 2011, pp. 219–220.
- [39] V. Saini, F. Farmahinifarahani, Y. Lu, P. Baldi, and C. V. Lopes, “Oreo: Detection of clones in the twilight zone”, in *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ACM, 2018, pp. 354–365.
- [40] Y. Higo, H. Murakami, and S. Kusumoto, “Revisiting capability of pdg-based clone detection”, Citeseer, Tech. Rep., 2013.
- [41] N. Smith, D. van Bruggen, and F. Tomassetti, *Javaparser*, May 2018.
- [42] F. A. Fontana, M. Zانونi, and F. Zانونi, “Duplicated code refactoring advisor (dcra): A tool aimed at suggesting the best refactoring techniques of java code clones”, PhD thesis, UNIVERSITÀ DEGLI STUDI DI MILANO-BICOCCA, 2012.
- [43] M. Allamanis and C. Sutton, “Mining Source Code Repositories at Massive Scale using Language Modeling”, in *The 10th Working Conference on Mining Software Repositories*, IEEE, 2013, pp. 207–216.
- [44] Y. Higo, T. Kamiya, S. Kusumoto, K. Inoue, and K. Words, “Aries: Refactoring support environment based on code clone analysis.”, in *IASTED Conf. on Software Engineering and Applications*, 2004, pp. 222–229.
- [45] Y. Higo, S. Kusumoto, and K. Inoue, “A metric-based approach to identifying refactoring opportunities for merging code clones in a java software system”, *Journal of Software Maintenance and Evolution: Research and Practice*, vol. 20, no. 6, pp. 435–461, 2008.
- [46] A. Alwaqfi, “A refactoring technique for large groups of software clones”, PhD thesis, Concordia University, 2017.
- [47] Z. Chen, Y.-W. Kwon, and M. Song, “Clone refactoring inspection by summarizing clone refactorings and detecting inconsistent changes during software evolution”, *Journal of Software: Evolution and Process*, vol. 30, no. 10, e1951, 2018.
- [48] G. G. Koni-N’Sapu, “A scenario based approach for refactoring duplicated code in object oriented systems”, *Master’s thesis, University of Bern*, 2001.
- [49] J. Visser, S. Rigal, G. Wijnholds, P. van Eck, and R. van der Leek, *Building Maintainable Software, Java Edition: Ten Guidelines for Future-Proof Code*. " O’Reilly Media, Inc.", 2016.
- [50] G. A. Campbell, “Cognitive complexity-a new way of measuring understandability”, *Technical Report. SonarSource SA, Switzerland*, 2017.
- [51] D. Landman, A. Serebrenik, E. Bouwers, and J. J. Vinju, “Empirical analysis of the relationship between cc and sloc in a large corpus of java methods and c functions”, *Journal of Software: Evolution and Process*, vol. 28, no. 7, pp. 589–618, 2016.
- [52] S. H. Costello, “Software engineering under deadline pressure”, *ACM SIGSOFT Software Engineering Notes*, vol. 9, no. 5, pp. 15–19, 1984.
- [53] R. D. Austin, “The effects of time pressure on quality in software development: An agency model”, *Information systems research*, vol. 12, no. 2, pp. 195–207, 2001.
- [54] H. Shah, M. J. Harrold, and S. Sinha, “Global software testing under deadline pressure: Vendor-side experiences”, *Information and Software Technology*, vol. 56, no. 1, pp. 6–19, 2014.
- [55] P. H. Cheney, “Effects of individual characteristics, organizational factors and task characteristics on computer programmer productivity and job satisfaction”, *Information & management*, vol. 7, no. 4, pp. 209–214, 1984.
- [56] M. Allamanis, E. Barr, C. Bird, and C. Sutton, “Suggesting accurate method and class names”, Aug. 2015, pp. 38–49. DOI: 10.1145/2786805.2786849.
- [57] U. Alon, O. Levy, and E. Yahav, “Code2seq: Generating sequences from structured representations of code”, *arXiv preprint arXiv:1808.01400*, 2018.

- [58] J. Eremondi, R. Fabian, and I. Akhter, “Sisyphus - crawling up the library to prevent code duplication”, Mar. 2017.

Appendix A

Non-crucial information