

# Mining Software Repositories for Contextual Information of Code Clones

Simon Baars<sup>1,2</sup>[0000–0001–7905–1027] and Ana Oprescu<sup>1,3</sup>[0000–0001–6376–0750]

<sup>1</sup> University of Amsterdam

<sup>2</sup> [simon.j.baars@gmail.com](mailto:simon.j.baars@gmail.com)

<sup>3</sup> [A.M.Oprescu@uva.nl](mailto:A.M.Oprescu@uva.nl)

**Abstract.** The abstract should briefly summarize the contents of the paper in 15–250 words.

**Keywords:** Code Clones · MSR · Clone Relation · Inheritance · Object-Oriented Programming.

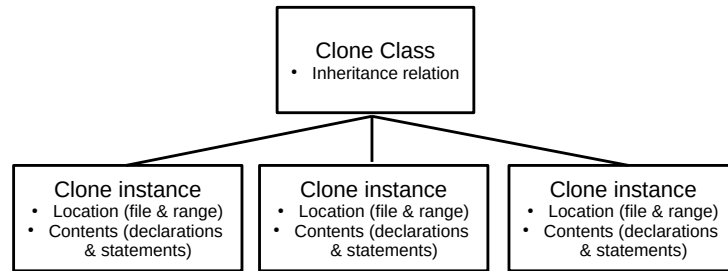
## 1 Introduction

## 2 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.



**Fig. 1.** Abstract representation of clone classes and clone instances.

Fig. 1 shows an abstract representation of clone classes and clone instances. The relation of clones through inheritance is determined for each clone class. The location and contents of clones are determined for each clone instance.

## 2.1 Relation

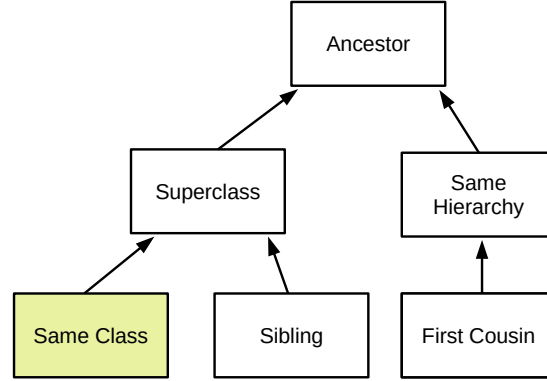
When merging code clones in object-oriented languages, it is important to consider the inheritance relation between clone instances. This relation has a big impact on how a clone should be refactored.

Fontana et al. [1] 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 are 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 Direct Interface:** All instances of the clone class are in a class or interface implement the same interface.
2. **Same Indirect Interface:** All instances of the clone class are in a class or interface that have a common interface anywhere in their inheritance hierarchy.
3. **No Direct Superclass:** All instances of the clone class are in a class that does not have any superclass.
4. **No Indirect Superclass:** All instances of the clone class are in a class that does not have any external classes in its inheritance hierarchy.
5. **External Ancestor:** All instances of the clone class are in a class that does not have any external classes in its inheritance hierarchy.



**Fig. 2.** 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 Direct Interface, Same Indirect 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 the above list that the clone class applies to. For instance: all “Superclass” clones also apply to “Same Hierarchy”, but because “Superclass” is earlier in the above list they will get the “Superclass” relation. This is because the items earlier in the above list denote a more favorable refactoring.

When CloneRefactor applies automated refactorings, it uses this inheritance relation to see where it must place the refactored code. We explain this for each relation category over the following sections.

**Common Class** The *Same method* and *Same class* relations share a common refactoring opportunity. Clones of both these categories, 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.

**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 [2].

**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 stricter visibility. Because of that, we favor all “Common Hierarchy” refactoring opportunities over “Common Interface”.

**Unrelated** Clones are unrelated if they share no common class or interface in their inheritance structure. These clones are least favorable for 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 classes with a *No Direct Superclass* relation mark the opportunity for creating a superclass abstraction and placing the extracted method there. For clone classes with a *No Indirect Superclass* relation, CloneRefactor creates such an abstraction for the ancestor that does not have a parent. Clone classes with a *External Superclass* or *External Ancestor* relation obstruct the possibility of creating a superclass abstraction. In such a case, CloneRefactor creates an interface abstraction to make their relation explicit.

## 2.2 Location

A paper by Lozano et al. [3] 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/Enumeration Level:** A clone that is (a part of) an interface or 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 the 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*.

**Method/Constructor Level Clones** Method/Constructor Level clones denote clones that are found in either a method or constructor. A constructor is a special method that is called when an object is instantiated. Most modern clone refactoring studies only focus on clones at method level [4–14]. This is because most clones reside at those places [1, 3] and most of those clones can be refactored with a relatively simple set of refactoring techniques [1, 6].

**Class/Interface/Enumeration Level Clones** Class/Interface/Enumeration Level clone instances are found inside the body of one of these declarations and optionally include the declaration itself. It can also be a clone instance that exceeds the boundaries of a single method. These clone instances can contain fields, (abstract) methods, inner classes, enumeration fields, etc. These types of clones require various refactoring techniques to refactor. For instance, we might have to move fields in an inheritance hierarchy. Or, we might have to perform a refactoring on more of an architectural level, if a large set of methods is cloned.

### 2.3 Contents

Finally, we looked at what nodes individual clone instances span. We selected a set of categories based on empirical evaluation of a set of clones in our dataset. We selected the following categories to be relevant for refactoring:

1. **Full Method/Constructor/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 part of) the body of a method/constructor. The declaration itself is not included.
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. **Other:** Anything that does not match with above-stated categories.

**Full Method/Constructor/Class/Interface/Enumeration** These categories denote that a full declaration, including its body, is cloned with another declaration. These categories often denote redundancy and are often easy to refactor: one of both declarations is redundant and should be removed. All usages of the removed declaration should be redirected to the clone instance that was not removed. Sometimes, the declaration should be moved to a location that is accessible by all usages.

**Partial Method/Constructor** These categories describe clone instances which are found in the body of a method or constructor. These clones can often be refactored by extracting a new method out of the cloned code.

**Several Methods** Several methods cloned in a single class is a strong indication of implicit dependencies between two classes. This increases the chance that these classes are missing some form of abstraction, or their abstraction is used inadequately.

**Only Fields** This category denotes that the clone spans over only global variables/fields that are declared outside of a method. This indicates data redundancy: pieces of data have an implicit dependency. In such cases, these fields may have to be encapsulated in a new object. Or, the fields should be somewhere in the inheritance structure where all objects containing the clone can access them.

**Other** The “Other” category denotes all configurations of clone contents that do not fall into above categories. Often, these are combinations of the above stated concepts. For instance, a combination of constructors and methods or a combination of fields and methods is cloned. Such clones indicate, like the “Several Methods”, the requirement of performing a more architectural-level refactoring. These are often more complicated to refactor, especially when aiming to automate this process.

### 3 Experimental Setup

All our experiments are quantitative experiments, measured over a corpus using CloneRefactor. In this chapter, we first explain the corpus we use. We then explain how we calculate the impact of refactorings on the software maintainability.

#### 3.1 The Corpus

For our experiments we use a large corpus of open-source projects assembled by Allamanis et al. [15]<sup>4</sup>. This corpus contains a set of Java projects from GitHub,

<sup>4</sup> The corpus can be downloaded at: <http://groups.inf.ed.ac.uk/cup/javaGithub/>

selected by the number of forks. The projects and files in this corpus were de-duplicated manually. 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.

Because CloneRefactor requires all dependencies for the projects it analyses, we created a set of scripts<sup>5</sup> to filter the corpus for all projects for which we can obtain all dependencies using Maven.

### 3.2 Resulting corpus

This procedure results in 2,267 Java projects including all their dependencies<sup>6</sup>. The projects vary in size and quality. The total size of all projects is 14,210,357 lines (11,315,484 when excluding whitespace) over a total of 99,586 Java files. This is an average of 6,268 lines over an average of 44 files per project, 141 lines on average per file. The largest project in the corpus is *VisAD* with 502,052 lines over 1,527 files.

### 3.3 Tool Validation

We have validated the correctness of CloneRefactor through unit tests and empirical validation. First, we created a set of 57 control projects<sup>7</sup> to verify the correctness in many (edge) cases. These projects test each identified relation, location, contents and refactorability category (see Section 2 and ??), to see whether they are correctly identified. Next, we run the tool over the corpus and manually verify samples of the acquired results. This way, we check the correctness of the identified clones, their context, and their proposed refactoring.

We also test the correctness of the resulting code after refactoring. For this, we use a project named JFreeChart<sup>8</sup>. JFreeChart has high test coverage and working tests, which allows us to test the correctness of the program after running CloneRefactor. We run CloneRefactor manually over the JFreeChart project, run its test cases and check the correctness of the proposed refactorings.

## 4 Results

To determine the refactoring method(s) that can be used to refactor most clones, we perform statistical analysis on the context of clones (see Section 2).

<sup>5</sup> All scripts to prepare the corpus are available on GitHub: <https://github.com/SimonBaars/GitHub-Java-Corpus-Scripts>

<sup>6</sup> The full list of projects is in the following file in our GitHub repository: [https://github.com/SimonBaars/GitHub-Java-Corpus-Scripts/blob/master/filtered\\_projects.txt](https://github.com/SimonBaars/GitHub-Java-Corpus-Scripts/blob/master/filtered_projects.txt)

<sup>7</sup> Control projects for testing CloneRefactor: <https://github.com/SimonBaars/CloneRefactor/tree/master/src/test/resources>

<sup>8</sup> JFreeChart is available on GitHub: <https://github.com/jfree/jfreechart>

#### 4.1 Relation

Table 1 displays the number of clone classes found for the entire corpus for different relations (see Section 2.1).

<i>Category</i>	<i>Relation</i>	<i>Clone Classes</i>	<i>%</i>	<i>Total</i>	<i>%</i>
Common Class	Same Class	22,893	26.8%	31,848	37.2%
	Same Method	8,955	10.5%		
Common Hierarchy	Sibling	15,588	18.2%	20,342	23.8%
	Superclass	2,616	3.1%		
	First Cousin	1,219	1.4%		
	Common Hierarchy	720	0.8%		
	Ancestor	199	0.2%		
	No Direct Superclass	10,677	12.5%		
Unrelated	External Superclass	4,525	5.3%	20,314	23.7%
	External Ancestor	3,347	3.9%		
	No Indirect Superclass	1,765	2.1%		
Common Interface	Same Direct Interface	7,522	8.8%	13,074	15.3%
	Same Indirect Interface	5,552	6.5%		

**Table 1.** Number of clone classes per clone relation.

Our results show that most clones (37%) are in a common class. 24% of clones are in a common hierarchy. Another 24% of clones are unrelated. 15% of clones are in an interface.

#### 4.2 Location

Table 3 displays the number of clone classes found for the entire corpus for different location categories (see Section 2.2). We can see from these results that nearly 80% of clones are found at method level. 17% of clones are found at class level, meaning they exceed the boundaries of a single method (or do not span methods at all). Constructors account for approximately 3% of clones. In interfaces, only 1% of clones is found.



<i>Category</i>	<i>Clone instances</i>	<i>%</i>
Method Level	232,545	78.43%
Class Level	50,402	17.00%
Constructor Level	10,039	3.39%
Interface Level	2,693	0.91%
Enum Level	788	0.27%

**Table 2.** Amount of clone instances with a per location category.

### 4.3 Contents

Table 3 displays the number of clone classes found for the entire corpus for different content categories (see Section 2.3).

<i>Category</i>	<i>Contents</i>	<i>Clone instances</i>	<i>Total</i>
Partial	Method Body	219,540 74.05%	229,521 77.42%
	Constructor Body	9,981 3.37%	
Other	Several Methods	22,749 7.67%	53,773 18.14%
	Only Fields	17,700 5.97%	
	Other	13,324 4.49%	
Full	Full Method	12,990 4.38%	13,173 4.44%
	Full Interface	64 0.02%	
	Full Constructor	58 0.02%	
	Full Class	37 0.01%	
	Full Enum	24 0.01%	

**Table 3.** Number of clone instances for clone contents categories

From these results, we see that 74% of clones span part of a method body (77% if we include constructors). 8% of clones span several methods. 6% of clones span only global variables. Only 4% of clones span a full declaration (method, class, constructor, etc.).

## 5 Discussion

Regarding clone context, our results indicate that most clones (37%) are in a common class. This is favorable for refactoring because the extracted method does not have to be moved after extraction. 24% of clones are in a common hierarchy. These refactorings are also often favorable. Another 24% of clones are

unrelated, which is often unfavorable because it often requires a more comprehensive refactoring. 15% of clones are in an interface.

Regarding clone contents, 74% of clones span part of a method body (77% if we include constructors). 8% of clones span several methods, which often require refactorings on a more architectural level. 6% of clones span only global variables, requiring an abstraction to encapsulate these data declarations. Only 4% of clones span a full declaration (method, class, constructor, etc.).

## 6 Conclusion

## References

- [1] F. A. Fontana and M. Zanoni, “A duplicated code refactoring advisor,” in *International Conference on Agile Software Development*, Springer, 2015, pp. 3–14.
- [2] M. Fowler, *Refactoring: improving the design of existing code*, Second. Addison-Wesley Professional, 2018.
- [3] 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.
- [4] E. Choi, N. Yoshida, T. Ishio, K. Inoue, and T. Sano, “Extracting code clones for refactoring using combinations of clone metrics,” in *Proceedings of the 5th International Workshop on Software Clones*, ACM, 2011, pp. 7–13.
- [5] R. Yue, Z. Gao, N. Meng, Y. Xiong, X. Wang, and J. D. Morgenthaler, “Automatic clone recommendation for refactoring based on the present and the past,” in *2018 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, IEEE, 2018, pp. 115–126.
- [6] 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.
- [7] F. Arcelli Fontana, M. Zanoni, A. Ranchetti, and D. Ranchetti, “Software clone detection and refactoring,” *ISRN Software Engineering*, vol. 2013, 2013.
- [8] Y. Lin, Z. Xing, X. Peng, Y. Liu, J. Sun, W. Zhao, and J. Dong, “Clonepedia: Summarizing code clones by common syntactic context for software maintenance,” in *Software Maintenance and Evolution (ICSME), 2014 IEEE International Conference on*, IEEE, 2014, pp. 341–350.
- [9] M. Mandal, C. K. Roy, and K. A. Schneider, “Automatic ranking of clones for refactoring through mining association rules,” in *Software Maintenance, Reengineering and Reverse Engineering (CSMR-WCRE), 2014 Software Evolution Week-IEEE Conference on*, IEEE, 2014, pp. 114–123.

- [10] M. Balazinska, E. Merlo, M. Dagenais, B. Lague, and K. Kontogiannis, “Advanced clone-analysis to support object-oriented system refactoring,” in *Reverse Engineering, 2000. Proceedings. Seventh Working Conference on*, IEEE, 2000, pp. 98–107.
- [11] Y. Yongting, L. Dongsheng, and Z. Liping, “Detection technology and application of clone refactoring,” in *Proceedings of the 2018 2nd International Conference on Management Engineering, Software Engineering and Service Sciences*, ACM, 2018, pp. 128–133.
- [12] S. Bouktif, G. Antoniol, E. Merlo, and M. Neteler, “A novel approach to optimize clone refactoring activity,” in *Proceedings of the 8th annual conference on Genetic and evolutionary computation*, ACM, 2006, pp. 1885–1892.
- [13] M. Fanqi, “Using self organized mapping to seek refactorable code clone,” in *Communication Systems and Network Technologies (CSNT), 2014 Fourth International Conference on*, IEEE, 2014, pp. 851–855.
- [14] U. Devi, A. Sharma, and N. Kesswani, “A study on the nature of code clone occurrence predominantly in feature oriented programming and the prospects of refactoring,” *International Journal of Computer Applications*, vol. 141, no. 8, 2016.
- [15] 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.