

Graph Theoretical Algorithms For Structural Comparison Of Java Source And Byte Code

Submitted By
Artem Garishin



FB2: Faculty of Computer Science and Engineering

*This thesis presented for the degree of
Master of Science
in the*

High Integrity Systems

Research Supervisor: Prof. Dr. Sergej Alekseev

Co-Supervisor: Prof. Dr. Matthias Wagner

September 2014

Legal Declaration

I declare that this thesis document is completely my own work and all used references have been clearly cited. I have not submitted this assignment in the context of an examination to any other examination board or person.

Signature:

Location, Date:

Abstract

TODO:

This should be a 1-page (maximum) summary of your work. What environment for development has been used, experimental results. An abstract is a summary in your own words of the Thesis. It is not evaluative and must not include your personal opinions. The purpose of an abstract is to give a reader sufficient information for him or her to decide whether it would be worthwhile reading the entire article or book. An abstract should aim at giving as much information as possible in as few words as possible.

Goal of this work is to search out the most optimal ways to compare different pieces of code. So far there are two techniques for code comparison: a normal text comparison and visual compare. Normal text-compare can be not sufficient to analyse two pieces of code, or to find a similarity between them. For that reason a structural/graph compare opens a doors to discover more possibilities of comparison.

Acknowledgments

I would like to take this time to thank Frankfurt University of Applied Sciences for all of the resources which they provided me in order to pursuing my master study in computer science and make this thesis possible.

I would like to express my sincere gratitude to Prof. Dr. Sergej Alekseev and Prof. Dr. Matthias Wagner for their patient guidance, encouragement and advice which they provided me throughout this thesis work.

Contents

1	Introduction	1
2	Description of problem	3
3	Graphs comparisons algorithms	6
3.1	An algorithm for Top-down maximum common sub-tree isomorphism	6
3.2	An algorithm for Bottom-Up maximum common sub-tree isomorphism	6
4	Code compare experiments	7
4.1	Introduction to experiments	7
4.2	Experiments on Java source code	9
4.3	Experiments using Abstract Syntax Tree graphs	11
5	Graph transformation algorithms	14
5.1	Introduction to the graph transformation	14
5.2	Techniques to build a tree from code	14
5.3	Convert graph to tree	14
5.4	Possible ideas	14
6	Existing Comparison methods	15
6.1	Plagiarism detection methods	15
7	Conclusion	17
	Bibliography	18

List of Figures

4.1	Text to text comparison example	13
4.2	Graph comparison on similar AST trees	13

List of Tables

Abbreviations

TDMC	T op D own M ax C ommon
BUMC	B ottom U p M ax C ommon
BCV	B yte C ode V isualizer
SCV	S ource C ode V isualizer
CFGF	C ontrol F low G raph F actory
TC	T ext C ompare

Chapter. 1

Introduction

Modern methods to compare of programming pieces of code are used to analyze code's changing, to explore development process and so on. Basically in current tools or plug-ins only text compare methods are used, that is not full sufficient to define code compare. Sometimes another techniques can be very helpful for such purposes. One of them is a structural code compare, based on building a trees, and methods to compare any similar or same structures.

TODO START: You can't write a good introduction until you know what the body of the paper says. Consider writing the introductory section(s) after you have completed the rest of the paper, rather than before. Be sure to include a hook at the beginning of the introduction. This is a statement of something sufficiently interesting to motivate your reader to read the rest of the paper, it is an important/interesting scientific problem that your paper either solves or addresses. You should draw the reader in and make them want to read the rest of the paper.

REDO: Code duplication or copying a code fragment and then reuse by pasting with or without any modifications is a well known code smell in software maintenance. Several studies show that about 5 to 20 percent of a software systems

can contain duplicated code, which is basically the results of copying existing code fragments and using them by pasting with or without minor modifications. One of the major shortcomings of such duplicated fragments is that if a bug is detected in a code fragment, all the other fragments similar to it should be investigated to check the possible existence of the same bug in the similar fragments. Refactoring of the duplicated code is another prime issue in software maintenance although several studies claim that refactoring of certain clones are not desirable and there is a risk of removing them. However, it is also widely agreed that clones should at least be detected. REDO

Tips: A statement of the goal of the paper: why the study was undertaken, or why the paper was written. Do not repeat the abstract.

Chapter. 2

Description of problem

In this chapter an issue of the work is being explained. As usual a compare of two codes we can consider them as classes, functions or methods. Thereby a compare can be counted as examinations of two pieces of code, in the best case a methods or functions. They can have a similar implementation or alike syntax, however these two pieces of code are different.

There are many purposes to compare a code, to find out a similarity or determine a difference between them. One of the option is to search for plagiarism in case a code can be taken from external source and a variables have been changed. In addition to general search can be improved to look out a similar code in big projects.

If two graphs are being compared with each other, then it's NP - complete problem and takes much times and efforts to be done. This type of graph comparison is very expensive from a computational point of view and thus, and some action must be taken into account to reduce the domain of comparison before performing the actual comparison.

Therefore this problem must to be reduced. Luckily, a tree comparison is able to executed in polynomial time and moreover there are some algorithms to com-

pare them.

TODO: write what is NP-complete, write a problems with time execution

TODO: create a pieces of code and compare them

Thus there are no deterministic algorithms to compare them because of loops in graphs. In this case the input code can be transformed into graph firstly, after the graph creation, it must be converted into tree, using simple techniques removing back edges. The back edges are edges point out the same node.

After described code transformation a three questions can be asked: 1. how to transform code into tree optimally 2. how to compare these trees 3. how to reference code pieces and nodes(how to put the code difference)

Regards to the first question, the concept of idea is described in section "Graph Transformation". The second question comprehends existing algorithm and their combination and improvements in chapter "Existing algorithms". The last issue is about how to lead back the result of the code.

TODO: look in internet what exists already(code compare)

Possible result of this thesis is development of concept to find out code difference using graph theory, in the best case a tool in Project Dr. Garbage [1] . can be implemented that highlights similarity/difference of input code snippet and represented respective graph.

To get started searching optimal way to find similarity/difference two AST trees, compare them, find a logic how to link nodes with the code. Small example demonstrating, what kind of result gives text compare(just string compare)and AST:

TODO: highlight this code

```
public void ast1() if(i < 1) i++;
```

```
public void ast2() if(i < 1) i++;
```

In this example, two pieces of code there one enter symbol after line (i < 1),

therefore the codes look different. Using simple text compare approach, only these gap will be found, however this difference does not play any role regards business-logic. In Abstract Syntax View, these two graphs will be same, and no discrepancy will have been discovered.

TODO: create an example of AST and text compare

Since this article includes comparison not only of source code, from where an Abstract Syntax Tree can be easily build, but also Java byte code, where there is no syntax. Basically to have a look at byte code example, there are no bounds to hold a functions or methods. Based on this, control flow graph can be derived from byte code, that represents a graph, but not a spanning tree. Every node has a reference to byte code address.

TODO: is any idea how to optimize spanning tree(node has a byte code address, topologically sort)

TODO: find java bytecode example and his tree

IDEA: text difference highlight and in parallel nodes in Tree marks and Matching (TODO: make some experiments, if it's ok, better, worse or same compare)

To investigate code comparison, two algorithms of structural compare are required, in fact Top down maximum common subtree and Bottom Up maximum common subtree algorithms. To make a contribution into development of structured code compare, the following tasks should be explored:

1. The existing algorithms must be investigated (The text-compare method is not sufficient to find a similarity in code)
2. The algorithms for the structured compare(Abstract syntax trees, Control flow graphs) must be explored
3. New methods and algorithms find a place to tried out. A prototypes of combination text-compare and structure-compare can be implemented.
4. Experimental results of compare must be derived.

Chapter. 3

Graphs comparisons algorithms

3.1 An algorithm for Top-down maximum common subtree isomorphism

TODO: explain briefly what this technique, examples from Valiente

3.2 An algorithm for Bottom-Up maximum common subtree isomorphism

TODO: explain briefly what this technique, examples from Valiente

Chapter. 4

Code compare experiments

4.1 Introduction to experiments

For better understanding of possible "code compare" concept, possible ideas of implementation and following development an amount of experiments are required. In order to build a proper tool or at least a concept, in Eclipse plugins at Dr. Garbage Community® some hand experiments in code should be fulfilled.

All these test cases are performed in Eclipse IDE [6] and divided into blocks. These steps can be approached by following:

1. Research on Java source code using existing methods to compare:
 - (a) Normal text compare
 - (b) Spanning trees transformed from control flow graphs
2. Research on Java source code using existing methods to compare:
 - (a) Normal text compare
 - (b) Abstract syntax trees
3. Research on Java byte code using existing methods to compare:

- (a) Normal text compare
- (b) Control flow graphs

All test set are investigated under Java methods and functions. Playing around with the patterns of code changing variables, names, sequences of commands, adding loops or conditions and apply the simple "text to text" compare. This "text compare" is already implemented in Eclipse IDE [6], so-called command "compare with each other by member". This type of comparison provides a pop-up window, where two pieces of code are compared, line by line.

TODO: Add Text compare picture

In parallel a control flow graphs or source graphs from the functions are being created and compared using implemented algorithms Top-Down and Bottom-Up(following called: TD& BU). The further task is to figure out the difference/similarity from graphical visual comparison. Consequently these both results must be matched and recorded for succeeding research.

The derived results from can be as follows:

1. Text compare and TD & BU have same difference
2. Text compare and TD & BU give similar difference
3. Text compare and TD & BU five full difference

Hence, as it was declared in section description of problem, based on these results can be decided what kind of tool in Dr. Garbage Eclipse plug-ins can be built. In case similar or full difference results, a combination of both methods can be used for optimal comparison.

Small example can be demonstrated: there are two functions that look very similar but nevertheless they have different number of string and different functionality. The abstract results can be following:

1. Text compare shows that strings 1 and 5 are different
2. Graph compare shows that string 7 is different

Conclusion: a combination of two methods can explicit that strings 1,5 and 7 are distinguished and much more distinction has been found. Thus it provides an optimal way of investigation.

TODO: explain briefly how TEXT compare works(search in internet), and sometimes it's not enough to observe the difference

4.2 Experiments on Java source code

TODO: write here how source code generated with dr. Garbage project

In this section a sequence of actions and following statistic are represented. The steps are carried through sequence of action:

- Write two similar functions in Eclipse IDE
- Apply for them text-to-text comparison
- Declare the statistic, respectively how many lines are different
- Create a source-code graph for both
- Apply TD& BU algorithms to get structural difference
- Declare the statistic, respectively how many nodes are different

TODO start: write how to evaluate percentage difference in code

In this section all experiments have been executed by hand using plug-in tool "Graph Comparison" in Eclipse from dr. Garbage project. Before starting, a several rules how to evaluate code difference in text and in graph must be established. It one of the crucial moment, because followed data statistic are used in further tool development.

Once java methods are created be compared, the rules to measure percentage are above. For business logic difference, that means how different methods from logic perspective. As a small instance, **if condition** is contrariwise, then the logic is being changed by 100%.

- Each java operator is considered like simple block percentage. For example,

```
int i = 1;
if(k == 0) { ... }
```

have two block to estimate the difference.
- Changing a conditions of block on the contrary issues 100% difference.
- Availability of extra variables in second piece of code is calculated by division of number of extra variables to amount of all variables.

Mostly all calculation are performed by roughly, because there are many criteria how to evaluate the logic difference. But from this perspective these rules are enough to examine graph's similarity.

For the text difference found via Eclipse tool some rules can be added:

- Percentage is computed thus division of the maximum number of lines to lines where the difference was found.
- If in one line of code only one symbol has been covered as found, then it is division of one to amount of symbols in this line.

TODO end

Conclusion about java source code stream line compare: After finishing experiments comparing java source code, generated by time line, the following outcomes have been derived:

- If the business logic is totally different (for example **if conditions**) Java Text Compare finds the difference and the graphical compare TD& BU are not able to see it
- In opposite said above, the graphical compare is quite useful to investigate for a code structure, for example if another third party person has changed local variables, the structure remains same, thus TD& BU (especially TD) find high level of similarity.

- The graph is totally bound with lines of code. If one brace is shifted, then it's considered one more block in the graph(Dr. Garbage Source Code Visualizer [1] generates extra node for each code operator). Hence TD cannot find following nodes.

4.3 Experiments using Abstract Syntax Tree graphs

In this section the abstract syntax trees are generated from Java source code using Dr. Garbage plugins [1]. The most notable advantage of building AST trees is a direct converting Java source code into AST tree, thereby avoiding graphs with cycles. Thus there is no need to delete back edges(see Spanning tree algorithms).

TODO: short explanation about AST how it looks like;

In this section all experiments have been executed by hand using plug-in tool "Graph Comparison" in Eclipse from dr. Garbage project [1]. Before starting, a several rules how to evaluate code difference in text and in graph must be established. It one of the crucial moment, because followed data statistic are used in further tool development.

TODO: after this compares write a conclusion what is better to compare

IDEA: compare AST AST optimization, for each node the label must be equal to the corresponding value in the code. go through the nodes, and compare the values in the nodes. Optimization: substitute some AST construction like "i++" to the $i = i + 1$, helps to find similarity AST native trees will more abstract (almost the same as JAVA compiler converts any code to). ATTACH tree AST and AST native

ARTICLE: Parsing: In case of parse tree-based approaches, the entire source code base is parsed to build parse tree or (annotated) abstract syntax tree (AST). In such representation, the source unit and comparison units are represented as subtrees of the parse tree or AST. Comparison algorithm then uses these subtrees to find clones [31, 213, 222]. Metrics-based approaches may also use such representation of code to calculate of the subtrees and find clones based on the metrics

values [146, 178].

```
public void test1(){
    int frameGroupLine = 10;
    for(int Cnt = 1; Cnt < frameGroupLine; Cnt += 2)
    {
        if(Cnt*4 != 2){
            frameGroupLine++;
        }
    }
}

public void test2(){
    int frameGroupLine = 10;
    for(int Counter = 1; Counter < frameGroupLine; Counter += 2)
    {
        if(Counter*4 != 2){
            frameGroupLine = frameGroupLine + 1;
        }
    }
}
```

If this line "frameGroupLine = frameGroupLine + 1;" will be converted into one format of sub-tree, that indicates the same as "frameGroupLine++;" then code

```
public void test3() int frameTeamLine = 10; for(int i = 1; i < frameTeamLine;
i += 2) if(i*4 != 2) frameTeamLine ++;
```

The code fragments test1() and test2() have the same application logic but different variables. However in the second one the increment of variable "frameTeamLine" is differently written. From text to text compare the functions are different at this point. If the text code similarity will be calculated, then these two fragments are not same(probably 90% similarity). Using Abstract Syntax Trees Optimization, this types of structure can be converted in the same sub-tree of whole AST tree. Thereby these two different text structures are represented as same sub-tree.

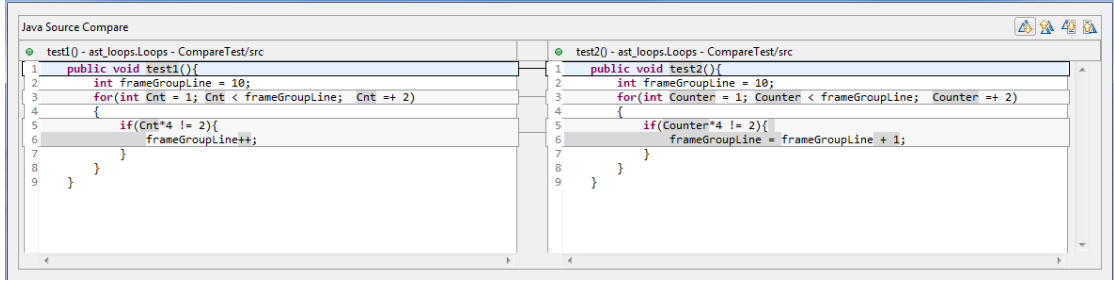


Figure 4.1: Example of standard text-to-text comparison of Java code

On figure 4.1 the example demonstrated how these two functions are compared using Eclipse Text comparator. After creation and comparison these two AST trees, it can be easily seen that sub-trees (the increment) are different.

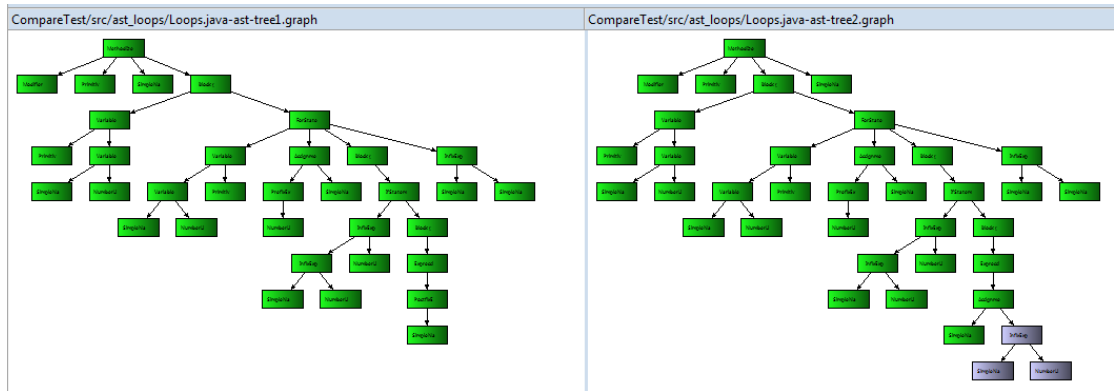


Figure 4.2: Graph comparison of function test1() and test2() using TDMC algorithm 3.1

From the figure 4.2 TDMC algorithm finds incomplete code similarity since not all nodes have been covered. From statistical point of view using simple math, the percentage of similarity is figured out: 37 nodes in the test2() and 3 of them are not covered. Thus, the calculation indicates $(1 - (\frac{3}{37})) \cdot 100 = 91\%$ code similarity according to AST trees and applied TDMC algorithm.

This mismatch can be optimized during AST tree production. These two lines of code `frameGroupLine++;` and `frameGroupLine = frameGroupLine + 1;` A simple replacement

Chapter. 5

Graph transformation algorithms

5.1 Introduction to the graph transformation

TODO: explain how graph is generated, which libraries are used, plugins

5.2 Techniques to build a tree from code

TODO: what is AST, examples;

TODO: what is BasicBlock, examples;

5.3 Convert graph to tree

TODO: describe here how to remove edges in order to get spanning tree

5.4 Possible ideas

Chapter. 6

Existing Comparison methods

6.1 Plagiarism detection methods

Task of plagiarism detection is an identification of text's similarity. Thus a research of existing methods is useful for this work regards code's comparison.

Plagiarism detection is the process of locating instances of plagiarism within a work or document. The widespread use of computers and the advent of the Internet has made it easier to plagiarize the work of others. Most cases of plagiarism are found in academia, where documents are typically essays or reports. However, plagiarism can be found in virtually any field, including scientific papers, art designs, and source code [4].

Mostly the task of plagiarism detection is considered for many fields, like text documents, software and source code. In this chapter source code plagiarism is being reviewed.

According to the article of Chanchal Kumar Roy and James R. Cordy [5], source-code similarity detection algorithms can be classified :

- Text-based Techniques
- Token-based Techniques

- Tree-based Techniques
- PDG-based Techniques
- Tree-based Techniques
- Metrics-based Techniques

Chapter. 7

Conclusion

TEMP based on experimental results and own opinion, write here what results were derived From time to time write here combined conclusions or improvements

Bibliography

- [1] The Dr. Garbage Tools Project® 2014, Sergej Alekseev, Peter Palaga and Sebastian Reschke, URL: <http://www.drgarbage.com>
- [2] Sergej Alekseev. *Graph theoretical algorithms for control flow graph comparison*, 2013.
- [3] Gabriel Valiente, *Algorithms on Trees and Graphs*, Berlin: Springer-Verlag, 2002.
- [4] Free content Internet encyclopedia - Wikipedia: Plagiarism detection, URL: http://en.wikipedia.org/wiki/Plagiarism_detection
- [5] Roy, Chanchal Kumar; Cordy, James R. (September 26, 2007). "*A Survey on Software Clone Detection Research*". School of Computing, Queen's University, Canada.
- [6] Eclipse documentation, URL: <http://help.eclipse.org/luna/index.jsp>
- [7] Sample Author, NAME, (Berlin: Springer-Verlag, 2002).