CODE CLONE DETECTION USING STRING BASED TREE MATCHING TECHNIQUE

WREARADILLA BINTI WAND

UNIVERSITI TEKNOLOGI MALAYSIA



UNIVERSITI TEKNOLOGI MALAYSIA PSZ 19:16 (Pind. 1/07)

DECLARATION OF THESIS / UNDERGRADUATE PROJECT PAPER AND COPYRIGHT			
	,		
Author's full name :	NORFA	RADILLA BINTI WAHID	
Date of birth :	<u>3 AUGI</u>	JST 1983	
Title :	CODE	CLONE DETECTION USING STRING BASED	
	STRING	G BASED TREE MATCHING TECHNIOUE	
Academic Session:	2008/ 20	009	
I declare that this the	sis is classifi	ed as :	
	NTIAL	(Contains confidential information under the Official Secret Act 1972)*	
		(Contains restricted information as specified by the organization where research was done)*	
		I agree that my thesis to be published as online open access (full text)	
l acknowledged that Universiti Teknologi Malaysia reserves the right as follows:			
 The thesis is the property of Universiti Teknologi Malaysia. The Library of Universiti Teknologi Malaysia has the right to make copies for the purpose of research only. The Library has the right to make copies of the thesis for academic exchange. 			
3. The library has the right to thake copies of the mess for academic exchange.			
SIGNATURE ASSOC PROF UP ALL BIN SELAMAT			
(NEW IC NO. /P.	ASSPORTN	O.) NAME OF SUPERVISOR	
Date : 31 001	OBER 20	Date: 31-10-08	
NET CONTRACTOR PORTION AND AND AND AND AND AND AND AND AND AN	• • • • • • • • • • • • • • • • • • • •		

NOTES : * If the thesis is CONFIDENTAL or RESTRICTED, please attach with the letter from the organization with period and reasons for confidentiality or restriction.

"I hereby declare that I have read this project report (course work) and in my opinion this project report is sufficient in terms of scope and quality for the award of the degree of Master of Science (Computer Science)"

> Signature Name of Supervisor Date

: ASSOC. PROF. DR ALI BIN SELAMAT : 3o - 1o - 08

CODE CLONE DETECTION USING STRING BASED TREE MATCHING TECHNIQUE

NORFARADILLA BINTI WAHID

A project report submitted in partial fulfillment of the requirements for the award of the degree of Master of Science (Computer Science)

Faculty of Computer Science and Information System Universiti Teknologi Malaysia

OCTOBER 2008

I declare that this project report entitled *"Code Clone Detection Using String Based Tree Matching Technique* "is the result of my own research except as cited in the references. The project report has not been accepted for any degree and is not concurrently submitted in candidature of any other degree.

Signature	:
Name	: NORFARADILLA BINTI WAHID
Date	. 31 ° C70 BEC 2608

To my beloved parents, fiancé and family.

ACKNOWLEDGEMENT

In preparing this report, I was in contact with many people, researchers, and academicians. They have contributed towards my understanding and thoughts. In particular, I wish to express my sincere appreciation to my project supervisor, Dr Ali Selamat, for encouragement, guidance, critics and friendship.

My fellow postgraduate students should also be recognized for their support. My sincere appreciation also extends to all my colleagues and others who have provided assistance at various occasions. Their views and tips are useful indeed. Without their continued support and interest, this thesis would not have been the same as presented here.

ABSTRAK

Pengklonan kod telah menjadi suatu isu sejak beberapa tahun kebelakangan ini selari dengan pertambahan jumlah aplikasi web dan perisian berdiri sendiri pada hari ini. Pengklonan memberi kesan yang sangat besar kepada fasa penyelenggaran sistem kerana secara tidak langsung peningkatan bilangan pengulangan kod yang sama di dalam sesebuah sistem akan menyebabkan kompleksiti sistem turut meningkat. Terdapat banyak teknik pengesanan klon telah dihasilkan pada hari ini dan secara umumnya ianya boleh dikategorikan kepada pengesanan berasaskan jujukan perkataan, token, pepohon dan semantik. Tujuan projek ini adalah untuk mengetahui kemungkinan untuk menggunakan suatu teknik dari pemetaan ontologi untuk menyelesaikan masalah ini, tetapi kami tidak menggunakan ontologi di dalam pengesanan klon. Telah dibuktikan di dalam eksperimen awalan bahawa ia mampu untuk mengesan klon. Di dalam tesis ini kami menggunakan dua aras pengesanan. Aras pertama menggunakan 'pelombong sub-pepohon terkerap' di mana ia mampu mengesan sub-pepohon yang sama antara fail yang berbeza. Kemudian sub-pepohon yang sama dinyatakan dalam bentuk ayat dan persamaan antara kedua-duanya dikira menggunakan 'metrik ayat'. Daripada eksperimen, kami mendapati bahawa sistem kami adalah tidak berganting kepada sebarang bahasa dah menghasilkan keputusan yang bagus dari segi *precision* tetapi tidak dari segi *recall*. Ia mampu mengesan klon serupa dan yang hamper sama.

ABSTRACT

Code cloning have been an issue in these few years as the number of available web application and stand alone software increase nowadays. The major consequences of cloning is that it would risk the maintenance process as there are many duplicated codes in the systems that practically increase the complexity of the system. There are many code clone detection techniques that can be found nowadays which generally can be group into string based, token based, tree based and semantic based. The aim of this project is to find out the possibility of using a technique of ontology mapping technique to solve the problem, but we are not using the real ontology for the clone detection. It has been prove that there is the possibility as it manages to detect clone code. In this thesis the clone detection is using two layers of detection; i.e. structural similarity and string based similarity. The structural similarity is by using subgraph miner where it capable to get the similar subtree between different files. And then we extract all elements of that particular subtree and treat the elements as a string. Two strings from different files then applied with similarity metric to know whether it is a clone pair. From the experimental result, we found that the system is language independent but the result is good in precision but not so good recall. It is also capable to detect two main types of clone, i.e identical clones and similar clones.

TABLE OF CONTENTS

TITLE

CHAPTER

	DEC	CLARATION	ii
	DED	DICATION	iii
	ACK	NOWLEDGEMENT	iv
	ABS	TRAK	V
	ABS	TRACT	vi
	TAB	vii	
	LIST	Х	
	LIST	F OF FIGURES	xi
	LIST	T OF ABBREVIATIONS	xiii
	LIST	T OF SYMBOLS	xiv
	LIST	F OF APPENDICES	XV
1	INT	RODUCTION	
	1.1	Overview	1
	1.2	Background of the Problem	2
	1.3	Problem Statement	5
	1.4	Objectives of the Project	6
	1.5	Scope of the Project	7
	1.7	Thesis outline	7
2	LIT	ERATURE REVIEW	
	2.1	Introduction	8
	2.2	Code Cloning	9
		2.2.1 Reasons of code cloning	11

PAGE

	2.2.2	Code cloning Consequences	14
	2.2.3	Code Cloning versus Plagiarism	15
	2.2.4	Code Cloning and the Software	
		Copyright Infringement Detection	16
2.3	Code	Cloning in web applications	17
	2.3.1	Definition of clones from web application research	'n
		View	19
	2.3.2	Source of Clones	19
2.4	Existi	ng Work of Code Cloning Detection	20
	2.4.1	String based	22
	2.4.2	Token based	23
	2.4.3	Tree based	24
	2.4.4	Semantic based	25
	2.4.5	Fingerprinting	25
	2.4.6	Analysis on Current Approaches	26
2.5	The S	emantic Web	28
	2.5.1	Architecture of the Semantic Web	29
	2.5.2	Web Ontology	30
	2.5.3	Web Ontology Description Languages	33
	2.5.4	Various Application of Ontology	34
	2.5.5	Ontology Mapping	36
	2.5.6	Ontology Mapping Approaches	39
	2.5.7	The Ontology Mapping Technique	40
		2.5.7.1 String Metrics	45
		2.5.7.2 Frequent Subgraph Mining	47
		2.5.7.3 MoFa, gSpan, FFSM, and Gaston	48
		2.5.7.4 Representing Web Programming as Tree	50
2.6	Clone	Detection Evaluation	52
2.7	Differ	ent with work by Jarzabek	54
	2.7.1	Clone Miner by Jarzabek	55
		2.7.1.1 Detection Of Simple Clones	56
		2.7.1.2 Finding Structural Clone	56
	2.7.2	Comparison of existing work and our	58
		proposed work.	

3 RESEARCH METHODOLOGY

3.1	Introduction	61
3.2	Proposed technique of clone detection	62
	3.2.1 Structural Tree Similarity	65
	3.2.2 String based tree matching	67
3.3	Preprocessing	70
3.4	Frequent subgraph mining	71
3.5	String based matching	73
3.6	Clone Detection Algorithm	75
3.7	Clone Detection Evaluation	75

4 EXPERIMENTAL RESULT AND DISCUSSION

4.1	Introduction	77
4.2	Data representation	78
	4.2.1 Original source program into XML format	79
	4.2.2 Subtree mining data representation	81
4.3	Frequent Subtree Mining	83
4.4	String metric computation	86
4.5	Experimental setup	87
4.6	Experimental results	88
4.7	Comparison of result using different parameters	96

5 CONCLUSION

5.1	Introduction	103
5.2	Future Works	104
5.3	Strength of the system	104

REFERENCES	105
Appendices A – C	112

LIST OF TABLES

TABLE NO.	TITLE	PAGE
2.1	A summary of code cloning and plagiarism detection	16
2.2	Brief description of ontology languages	33
2.3	List of string metric	45
3.1	Example of cross-table used to compare programs across two systems	71
3.2	Brief description of each frequent subgraph miner	72
4.1	Data for program testing	89
4.2	Experimental result using GSpan miner	91
4.3	Experiment using different parameter value	99

.

LIST OF FIGURE

FIGURE NO	TITLE	PAGE
1.1	A sample parse tree with generated characteristic vectors.	4
2.1	Example of a pair of cloned code in traditional program.	11
2.2	Tree-diagram for the Reasons for Cloning	13
2.3	Variation of clone detection research and the	
	classification of detection	22
2.4	Architecture of Semantic Web	29
2.5	Simple example of ontology	32
2.6	Simple example of mapping between two ontologies.	38
2.7	Illustration of ontology mapping approaches	40
2.8	Tree representation of an XML source code	48
2.9	Clones per file	57
2.10	Frequent clone pattern with file coverage	58

2.11	Similar node structure between two XML code fragments	59
2.12	Difference of work by Jarzabek and Basit(2005) and our proposed work	60
3.1	Mapping between concepts of O'_{α} and O'_{β}	65
3.2	Diagrammatic view of clone detection technique	69
3.3	Preprocessing phase	70
3.4	Illustration of detected clone within two trees	73
3.5	A pair of source code fragment classified as nearly identican nearly-identical	al. 74
3.6	Clone detection algorithm	75
4.1	Transformation of original PHP code into HTML code	80
4.2	XML form of the previous HTML code	81
4.3	A tree as list of nodes and edges	82
4.4	Example of tree as vertices and edges list	83
4.5	Frequent subtrees generated by graph miner.	85
4.6	Code fragment containing original frequent subtree.	87
4.7	Real output from the clone detection system	90
4.8	Recall and precision for GSpan-Jaro Winkler	93

4.9	Robustness of GSpan-Jaro Winkler	93
4.10	Computational time for GSpan-JaroWinkler	94
4.11	Recall and Precision for GSpan-Levenshtein Distance	94
4.12	Robustness for GSpan-Levenshtein Distance	95
4.13	Computational time for GSpan-Levenshtein Distance	95
4.14	Two close clones cannot be taken as a single clone	98
4.15	Precision result using different minimum support and threshold	100
4.16	Recall result using different minimum support and threshold	101

xiii

LIST OF ABBREVIATIONS

WA	Web Application
TS	Traditional Software
CCD	Code Clone Detection
PD	Plagiarism Detection

LIST OF SYMBOLS

<i>Sim(s</i> ₁ , <i>s</i> ₂)	-	similarity between two strings, s_1 and s_2
$Comm(s_1, s_2)$	-	commonality between s_1 and s_2
$Diff(s_1, s_2)$	-	difference between s_1 and s_2
$Winkler(s_1, s_2)$	-	improvement value to improve the result
max ComSubString	-	the sum of the lengths of common substring
$length(s_1)$	-	length of s ₁
$length(s_2)$	-	length of s ₂
uLen _{s1}	-	length of the unmatched substring from s_1
uLen ₂	-	length of the unmatched substring from s_2
р	-	a parameter of range 0 and ∞
θ	-	a threshold

LIST OF APPENDICES

APPENDIX	TITLE	PAGE
А	Project Activities	111
В	Existing Works of Code Clone Detection	113
С	Experimental result tables	117

CHAPTER 1

INTRODUCTION

1.1 Overview

As the world of computers is rapidly developing, there are tremendous needs of software development for different purposes. And as we can see today, the complexity of the software been developed are different between one and another. Sometimes, developers take easier way of implementation by copying some fragments of the existing programs and use the code in their work. This kind of work can be called as code cloning. Somehow the attitude of cloning can lead to the other issues of software development, for example the plagiarism and software copyright infringement (Roy and Cordy, 2007).

In most of the cases, in order to figure out the issues and to help better software maintenance, we need to detect the codes that have been cloned (Baker, 1995). In the web applications development, the chances of doing clones are bigger since there are too many open source software available in the Internet (Bailey and Burd, 2005). The applications are sometimes just a 'cosmetic' of another existing system. There are quite a number of researches in software code cloning detection, but not so particularly in the area of web based applications.

1.2 Background of the Problem

Software maintenance has been widely accepted as the most costly phase of a software lifecycle, with figures as high as 80% of the total development cost being reported (Baker, 1995). As cloning is one of the contributors towards this cost. the software clone detection and resolution has got considerable attention from the software engineering research community and many clone detection tools and techniques have been developed (Baker, 1995). However, when it comes to commercialization of the software codes, most of the software house developers tend to claim that their works are 100% done in house without using other codes copies form various sources. This has made a difficulty for the intellectual property copyright entities such as SIRIM and patent searching offices in finding the genuineity of the software source codes developed by the in house company. There is a need to identify the software source submitted for patent copyright application to be a genuine source code without having any copyright infringements. Besides that, the cloning is somehow raising the issue of plagiarism. The simplest example is in the academic area where students tend to copy their friends' works and submit the assignments with only slight modifications.

Usually, in software development process, there is a need for components reusability either in designing and coding. Reuse in object-oriented systems is made possible through different mechanisms such as inheritance, shared libraries, object composition, and so on. Still, programmers often need to reuse components which have not been designed for reuse. This may happen during the initial of systems development and also when the software systems go through the expansion phase and new requirements have to be satisfied. In these situations, the programmers usually follow the low cost copy-paste technique, instead of costly redesigning-the-system approach, hence causing clones. This type of code cloning is the most basic and widely used approach towards software reuse. Several studies suggest that as much as 20-30% of large software systems consist of cloned code (Krinke, 2001). The problem with code cloning is that errors in the original must be fixed in every copy. Other kinds of maintenance changes, for instance, extensions or

adaptations, must be applied multiple times. too. Yet. it is usually not documented where code was copied. In such cases, one needs to detect them. For large systems. detection is feasible only by automatic techniques. Consequently, several techniques have been proposed to detect clones automatically (Bellon et al., 2007).

There are quite a number of works that detect the similarity by representing the code in tree or graph representation and also some using string-based detection. and semantic-based detection. Almost all the clone detection technique had the tendency of detecting syntactic similarity and only some detect the semantic part of the clones. Baxter in his work (Baxter et al., 1998) proposes a technique to extract clone pairs of statements, declarations, or sequences of them from C source files. The tool parses source code to build an abstract syntax tree (AST) and compares its subtrees by characterization metrics (hash functions). The parser needs a "fullfledged" syntax analysis for C to build AST. Baxter's tool expands C macros (define. include, etc) to compare code portions written with macros. Its computation complexity is O(n), where n is the number of the subtree of the source files. The hash function enables one to do parameterized matching, to detect gapped clones, and to identify clones of code portions in which some statements are reordered. In AST approaches, it is able to transform the source tree to a regular form as we do in the transformation rules. However, the AST based transformation is generally expensive since it requires full syntax analysis and transformation.

In other work (Jiang et al. 2007) present an efficient algorithm for identifying similar subtrees and apply it to tree representations of source code. Their algorithm is based on a novel characterization of subtrees with numerical vectors in the Euclidean space Rⁿ and an efficient algorithm to cluster these vectors with respected to the Euclidean distance metric. Subtrees with vectors in one cluster are considered similar. They have implemented the tree similarity algorithm as a clone detection tool called DECKARD and evaluated it on large code bases written in C and Java including the Linux kernel and JDK. The experiments show that DECKARD is both scalable and accurate. It is also language independent, applicable to any language with a formally specified grammar.



Figure 1.1: A sample parse tree with generated characteristic vectors[14].

In (Krinke, 2001), Krinke presents an approach to identify similar code in programs based on finding similar subgraphs in attributed directed graphs. This approach is used on program dependence graphs and therefore considers not only the syntactic structure of programs but also the data flow within (as an abstraction of the semantics). As a result, it is said that no tradeoff between precision and recall- the approach is very good in both.

Kamiya in one of his work in (Kamiya et al., 2002) suggest the use of suffix tree. In the paper they have used a suffix-tree matching algorithm to compute tokenby token matching, in which the clone location information is represented as a tree with sharing nodes for leading identical subsequences and the clone detection is performed by searching the leading nodes on the tree. Their token-by token matching is more expensive than line-by-line matching in terms of computing complexity since a single line is usually composed of several tokens. They proposed several optimization techniques especially designed for the token-by-token matching algorithm, which enable the algorithm to be practically useful for large software.

Appendix B of this thesis, describe briefly some existing techniques of code clone detection and plagiarism. It also discusses the strength and weaknesses of each technique.