

DBPD: A Dynamic Birthmark-based Software Plagiarism Detection Tool

Zhenzhou Tian, Qinghua Zheng, Ming Fan, Eryue Zhuang, Haijun Wang, Ting Liu*

Ministry of Education Key Lab For Intelligent Networks and Network Security

Department of Computer Science and Technology, Xi'an Jiaotong University, Xi'an, 710049, China

{zztian, fanming.911025, hjwang.china}@stu.xjtu.edu.cn; {qzhzheng, tingliu}@mail.xjtu.edu.cn; zhuang8225@126.com

Abstract: *With the burst of open source software, software plagiarism has been a serious threat to the software industry. In this paper, we present the demo tool DBPD: **D**ynamic **B**irthmark-based **S**oftware **P**lagiarism **D**etection. Major features of DBPD could be summarized as: 1) dynamic birthmark. The execution process of software is captured to generate the birthmark reflecting intrinsic properties of software; 2) high availability. It is available for cross-platform and binary executable's plagiarism detection; 3) customizable. The birthmarks, similarity calculation metrics and detection criteria are configurable. The DBPD is implemented using C++ and Java, and currently can work under both Windows and Linux system. Three dynamic birthmarks are implemented in DBPD to identify the software according to its instruction, stack operation and system call.*

Keywords: software plagiarism detection, dynamic birthmark

I. INTRODUCTION

Free and open source software projects allow users to use, change and distribute software under certain types of license such as the well-known GPL. However, driven by the huge commercial interests, some individuals and companies incorporate third party software or libraries into their own products without respecting the licensing terms. Recent incidents include the lawsuit against Verizon by Free Software Foundation for distributing Busybox in its FIOS wireless routers [1], and the crisis of Skype's VOIP service for the violation of licensing terms of Joltid. The unavailability of source code and the existence of powerful automated semantic-preserving code transformation tools, make the plagiarism an easy to implement but difficult to detect thing.

Software birthmark, a set of characteristics extracted from a program that reflect the program's intrinsic properties and that can be used to uniquely identify the program, is a promising way for solving the plagiarism detection problem. However, despite the tremendous progress of birthmark based plagiarism detection approaches, seldom tools are publically available. The rare few tools as far as we find are SandMark [2], Stigmata [3] and Birthmarking [4]. The former two are static birthmark based which are believed to be fragile faced with semantic-preserving code obfuscation techniques, and the last one is dynamic birthmark based which is believed to have better performance than the previous two static birthmarks, yet they all suffer the problem of language dependence, since they're

*Ting Liu is the corresponding author. The research was supported in part by National Science Foundation of China (91118005, 91218301, 61221063, 61203174), 863 Program (2012AA011003), The Ministry of Education Innovation Research Team (IRT13035), Key Projects in the National Science and Technology Pillar Program of China (2012BAH16F02).

only valid for java programs. Also, there are some mature tools such as the JPlag [5] that target at source code which is not always available, since plagiarists are not likely to provide their source code before certain evidences are collected. Thus more powerful and practical tools are in urgent needs to fill the gap of birthmark based plagiarism detection research and practice.

It is a generally accepted fact that dynamic birthmarks being abstractions of runtime behaviors are believed to be more accurate reflections of program semantics than static birthmarks. Therefore, we implement a demo tool DBPD for plagiarism detection using dynamic birthmark techniques. Three dynamic birthmarks are implemented in DBPD to identify the software, including DKISB (dynamic key instruction sequence birthmark) [6], SODB (call stack operation dynamic birthmark) [7] and SCSSB (system call short sequence birthmark) [8]. Since all of them can work directly on binary executables, DBPD can analyze various programming languages.

II. TOOL OVERVIEW AND IMPLEMENTATION

A. Tool Overview

Fig.1 shows the overview of the DBPD. It consists of three main modules: the dynamic analysis module, the birthmark generator, the similarity calculator and decision maker. The modular architecture qualifies it with good scalability of easily introducing new kinds of dynamic birthmark methods.

Given two binary executables the plaintiff (original program), the defendant (suspicious program) and a set of inputs, DBPD executes both programs with the same input one by one. Meantime, the dynamic analysis module monitors the executions, performs dynamic analysis and collects execution profiles containing three kinds of events: key instructions, stack operations, and system calls. After sequences of both plaintiff and defendant programs are available, they are fed into the birthmark generator where noises are filtered, valid

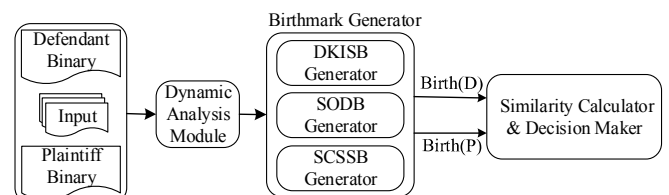


Fig. 1 Design overview of DBPD

