# The Software IP Detective's
# HANDBOOK

## Measurement, Comparison, and Infringement Detection

# BOB ZEIDMAN

# Praise for *The Software IP Detective's Handbook*

"As a frequent expert witness myself, I found Bob's book to be important and well written. Intellectual property and software plagiarism are complicated subjects, as are patents and copyrights. This book explains the key elements better than anything else I have seen. The book is highly recommended to anyone who develops software and also to those who need to protect proprietary software algorithms. The book should also be useful to attorneys who are involved with intellectual property litigation."

—*Capers Jones, president, Capers Jones & Associates, LLC*

"Intellectual property [IP] is an engine of growth for our high-tech world and a valuable commodity traded in its own right. Bob Zeidman is a leading authority on software IP, and in this book he shares with us his expertise. The book is comprehensive. It contains clear explanations of many difficult subjects. Businesspeople who study it will learn how to protect their IP. Lawyers will use it to understand the specifics of how software embodies IP. Judges will cite it in their decisions on IP litigation."

—*Abraham Sofaer, George P. Shultz Senior Fellow in Foreign Policy and National Security Affairs, Hoover Institution, Stanford University*

"Bob has done a fantastic job in making computer science forensics understandable to mere mortals: attorneys, engineers, and managers. This is the ultimate handbook for expert witnesses, due diligence execution, and developing a baseline for software valuation. Buy it before your competitors do!"

—*Don Shafer, CSDP, chief technology officer, Athens Group, LLC*

"Bob has considerable experience in dealing with issues associated with unauthorized use of software code. His insights in this book are particularly helpful for those seeking to provide expert analysis with respect to software code."

—*Neel I. Chatterjee, partner, Orrick, Herrington & Sutcliffe, LLC*

"This readable book perfectly bridges the jargon divide between software engineers and IP attorneys. It helps each group finally understand exactly what the other is talking about. As a software developer and expert witness I will definitely keep a copy handy and recommend it to others on my team."

—*Michael Barr, president, Netrino, LLC*

"This book makes intellectual property law understandable and accessible to programmers by combining discussions of the law with discussions of computer science, and interweaving case studies to elucidate the intersection of these disciplines."

—*Robert C. Seacord, Secure Coding Manager, Software Engineering Institute*

*This page intentionally left blank*

# THE SOFTWARE IP
# DETECTIVE'S HANDBOOK

*This page intentionally left blank*

# The Software IP Detective's Handbook

## MEASUREMENT, COMPARISON, AND INFRINGEMENT DETECTION

Robert Zeidman

Software Analysis and Forensic Engineering Corporation

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and the publisher was aware of a trademark claim, the designations have been printed with initial capital letters or in all capitals.

The author and publisher have taken care in the preparation of this book, but make no expressed or implied warranty of any kind and assume no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

The publisher offers excellent discounts on this book when ordered in quantity for bulk purchases or special sales, which may include electronic versions and/or custom covers and content particular to your business, training goals, marketing focus, and branding interests. For more information, please contact:

> U.S. Corporate and Government Sales
> (800) 382-3419
> corpsales@pearsontechgroup.com

For sales outside the United States please contact:

> International Sales
> international@pearson.com

Visit us on the Web: informit.com/ph

*Library of Congress Cataloging-in-Publication Data is on file with the Library of Congress*

> Pearson Education, Inc.
> Rights and Contracts Department
> 501 Boylston Street, Suite 900
> Boston, MA 02116
> Fax: (617) 671-3447

*This book is dedicated to all those who attempt to do things that others say are wrong or impossible, and to all those who encourage them.*

*This page intentionally left blank*

# CONTENTS

# PREFACE

## WHAT IS THIS BOOK ABOUT?

This book is generally about software intellectual property and specifically about the field of software forensics. While you may never testify in a courtroom, attempt to deduce the true owner of a valuable program, or measure the intellectual property (IP) value of software source code, if you are at all involved with software you need to know how to protect your rights to that software. This book will give you an understanding of those rights, their limitations, how to protect those rights, and how to take action against someone or some organization that you believe has infringed on those rights.

Unlike digital forensics, which studies the bits and bytes on a digital storage medium, such as a hard disk or DVD-ROM, without a deep understanding of what those ones and zeros represent, software forensics studies the software code that instructs a computer to perform operations. Software forensics discovers information about the history and usage of that software for presentation in a court of law. It combines information and techniques from computer science, mathematics, and law in a way that is unique and, I believe, particularly interesting.

## HOW IS THIS BOOK ORGANIZED?

This book contains overviews as well as in-depth information about law, mathematics, and computer science. The book is organized into chapters that

can be categorized as follows: those primarily about intellectual property law, those about computer science, those about mathematics, and those about business and business procedures. There is, of course, overlap. Chapter 1 describes the organization of the book in detail and allows you to choose those chapters that are most relevant to your career and your interests.

## WHO SHOULD READ THIS BOOK?

This book is for anyone interested in software intellectual property. Specifically, I believe the book will appeal to computer scientists, computer programmers, business managers, lawyers, engineering consultants, expert witnesses, and software entrepreneurs. While the focus is on software IP measurement, comparison, and infringement detection, the book also has useful information about many issues related to software IP, and I believe that many people involved with software will find the book valuable and interesting.

## SUPPORT AND COMMENTS

Thank you for purchasing my book. Please send me feedback on corrections and improvements. Of course, I also like to hear about the things I did right, the things you like about the book, and how it has helped you in some way.

Bob Zeidman
bob@SAFE-corp.biz
www.SAFE-corp.biz
Cupertino, California
March 2011

# ACKNOWLEDGMENTS

When I finished my last book, back in 2002, I swore I'd never write another one. They say that pregnancy is the same way. Ask a woman right after delivering a child whether she'd do it again and, if she's coherent, she'll give you an emphatic, "Never!" Yet within months, if not days, she forgets the pain of birth and only remembers the pleasure of spending time with her baby, then her infant, then her toddler, then her kid. (By the time the child is a teenager, some of the pain of birth may have returned to her memory, though.)

So it took me seven years to forget the pain of giving birth to my last book. Actually, I remembered that pain (though the memory had softened over the years) but really felt I had something worthwhile to say once again. I'm proud of all of my books (a parent isn't supposed to favor one child), but this one is different. Those other books were intended to explain areas of engineering that many others had invented, developed, and explored before me. I explained some of my own discoveries here and there, but the bulk of the creativity is correctly credited to the pioneers who preceded me. For this book, most of the work is original. The techniques, the mathematics, the algorithms, and the procedures in this book are being adopted fairly well, and that's exciting. I hope that others can take what I've done and build upon it. I continue to do that myself and am finding lots of unexplored areas that are being revealed, offering plenty of opportunities, I believe, for mathematicians, lawyers, programmers, computer scientists, and entrepreneurs.

Many people helped me with this book, some explicitly and some implicitly. First, I'll mention the lawyers who, despite their unbelievably full schedules, still found time to

# About the Author

**Bob Zeidman** is the president and founder of Zeidman Consulting (www.ZeidmanConsulting.com), a premiere contract research and development firm in Silicon Valley that focuses on engineering consulting to law firms handling intellectual property dispute cases. Since 1983, Bob has designed computer chips and circuit boards for RISC-based parallel processor systems, laser printers, network switches and routers, and other complex systems. His clients have included Apple Computer, Cisco Systems, Cadence Design Systems, Facebook, Intel, Symantec, Texas Instruments, and Zynga. Bob has worked on and testified in cases involving billions of dollars in disputed intellectual property.

Bob is also the president and founder of Software Analysis and Forensic Engineering Corporation (www.SAFE-corp.biz), the leading provider of software intellectual property analysis tools. Bob is considered a pioneer in the field of analyzing software source code, having created the CodeSuite program for detecting software intellectual property theft, which is sold by SAFE Corp, founded in 2007.

Previously, Bob was the president and founder of Zeidman Technologies (www.zeidman.biz) where he invented the patented SynthOS program for automatically generating real-time software. Before that, Bob was the president and

founder of The Chalkboard Network, an e-learning company that put high-end business and technology courses from well-known subject matter experts on the web. Prior to that, Bob invented the concept of remote backup and started Evault, the first remote backup company.

Bob is a prolific writer and instructor, giving seminars at conferences around the world. Among his publications are numerous articles on engineering and business as well as three textbooks—*Designing with FPGAs and CPLDs*, *Verilog Designer's Library*, and *Introduction to Verilog*. Bob holds numerous patents and earned two bachelor's degrees, in physics and electrical engineering, from Cornell University and a master's degree in electrical engineering from Stanford University.

Bob is also active in a number of nonprofits. He also enjoys writing novels and screenplays and has won a number of awards for this work.

# ABOUT THIS BOOK

This book crosses a number of different fields of computer science, mathematics, and law. Not all readers will want to delve into every chapter. This is the place to start, but from this point onward each reader's experience will be different. In this chapter I describe each of the parts and chapters of the book to help you determine which chapters will be useful and appealing for your specific needs and interests.

I should make clear that I am not a lawyer, have never been one, and have never even played one on TV. All of the issues I discuss in this book are my under-standing based on my technical consulting and expert witness work on nearly 100 intellectual property cases to date. My consulting company, Zeidman Consulting, has been growing over the years, and now the work is split between my employees and me. When I refer in the book to my experiences, in most cases that is firsthand information, but in other cases it may be information discovered and tested by an employee and related and explained to me.

In this book I also refer to forensic analysis tools that I have used to analyze software, in particular the CodeSuite tool that is produced and offered for sale by my software company, Software Analysis and Forensic Engineering Corporation (S.A.F.E. Corporation), and can be downloaded from the company website at www.SAFE-corp.biz. The CodeSuite set of tools currently consists of the following functions: BitMatch, CodeCLOC, CodeCross, CodeDiff, and SourceDetective. Functions are being continually added and updated. Each of these functions uses one or more of the algorithms described in later chapters.

Also, the CodeMeasure program uses the CLOC method to measure software evolution, which is explained in Chapter 12. It is also produced and sold by S.A.F.E. Corporation and can be downloaded from its own site at www.CodeMeasure.com.

Table 1.1 should help you determine which chapters will be the most helpful and relevant to you. Find your occupation at the top of the table and read downward to see the chapters that will be most relevant to your background and your job.

## PART I: INTRODUCTION

The introduction to the book is just that—an introduction, intended to give you a broad overview of the book and help you determine why you want to read it and which chapters you will find most in line with your own interests and needs. This part includes a description of the other parts and chapters in the book. It also gives information and statistics about intellectual property crime, to give you an understanding of why this book is useful and important.

## PART II: SOFTWARE

In this part I describe source code, object code, interpreted code, macros, and synthesis code, which are the blueprints for software. This part describes these important concepts, which are well known to computer scientists and programmers but may not be understood, or may not be understood in sufficient depth, by attorneys involved in software IP litigation. This part will be valuable for lawyers to help them understand how different kinds of software code relate to each other, and how these different kinds of software code can affect a software copyright infringement, software trade secret, or software patent case.

## PART III: INTELLECTUAL PROPERTY

In this part I describe intellectual property, in particular copyrights, patents, and trade secrets. I have found that many of these concepts are unclear or only partially understood by many computer scientists, programmers, and corporate managers. In this part I define these terms in ways that I believe will be comprehensible to those with little or no legal background.

**Table 1.1** Finding Your Way through This Book

| Chapter | Title | Computer scientist | Computer programmer | Manager | Lawyer | Consultant/ expert witness | Software entrepreneur |
|---|---|---|---|---|---|---|---|
| Part I | Introduction | X | X | X | X | X | X |
| Chapter 1 | About This Book | X | X | X | X | X | X |
| Chapter 2 | Intellectual Property Crime | X | X | X | X | X | X |
| Part II | Software | X | X | X | X | X | X |
| Chapter 3 | Source Code | | | X | X | | |
| Chapter 4 | Object Code and Assembly Code | | | X | X | | |
| Chapter 5 | Scripts, Intermediate Code, Macros, and Synthesis Primitives | | | X | X | | |
| Part III | Intellectual Property | X | X | X | X | X | X |
| Chapter 6 | Copyrights | X | X | X | | X | X |
| Chapter 7 | Patents | X | X | X | | X | X |
| Chapter 8 | Trade Secrets | X | X | X | | X | X |
| Chapter 9 | Software Forensics | X | X | X | X | X | X |
| Part IV | Source Code Differentiation | X | X | X | X | X | X |
| Chapter 10 | Theory | X | | | | X | X |
| Chapter 11 | Implementation | X | X | | | X | X |
| Chapter 12 | Applications | X | X | X | | X | X |
| Part V | Source Code Correlation | X | X | X | X | X | X |
| Chapter 13 | Plagiarism Detection | X | | | | X | X |
| Chapter 14 | Source Code Characterization | X | X | | X | X | X |

*Continues*

**Table 1.1**   Finding Your Way through This Book *(Continued)*

| Chapter | Title | Computer scientist | Computer programmer | Manager | Lawyer | Consultant/ expert witness | Software entrepreneur |
|---|---|---|---|---|---|---|---|
| Chapter 15 | Theory | X | | | | X | X |
| Chapter 16 | Implementation | X | X | | | X | X |
| Chapter 17 | Applications | X | X | X | | X | X |
| Part VI | Object and Source/Object Code Correlation | X | X | X | X | X | X |
| Chapter 18 | Theory | X | | | | X | X |
| Chapter 19 | Implementation | X | X | | | X | X |
| Chapter 20 | Applications | X | X | X | | X | X |
| Part VII | Source Code Cross-Correlation | X | X | X | X | X | X |
| Chapter 21 | Theory, Implementation, Application | X | X | X | | X | X |
| Part VIII | Detecting Software IP Theft and Infringement | X | X | X | X | X | X |
| Chapter 22 | Detecting Copyright Infringement | | | X | X | X | X |
| Chapter 23 | Detecting Patent Infringement | | | X | X | X | X |
| Chapter 24 | Detecting Trade Secret Theft | | | X | X | X | X |
| Part IX | Miscellaneous Topics | X | X | X | X | X | X |
| Chapter 25 | Implementing a Software Clean Room | | X | X | X | X | X |
| Chapter 26 | Open Source Software | | X | X | X | X | X |
| Chapter 27 | Digital Millennium Copyright Act | | X | X | X | X | X |
| Part X | Past, Present, and Future | X | X | X | | X | X |

I also define the field of software forensics in this part. When I am asked to work on a case, there is sometimes confusion about the fields of software forensics and digital forensics. In some cases, engineers practicing digital forensics claim to practice software forensics and sometimes use the tools of digital forensics to attempt to draw conclusions about software IP, yielding incorrect or inconclusive results. Software forensics requires the specialized tools of the field and expertise in the field to extract relevant information from the tools, reach appropriate conclusions, and opine on those conclusions. In this part I offer definitions of the two fields. In fact, the definition of software forensics has, to this point, been somewhat vague. My explanation in this part will clarify the practice of software forensics, show how it fits into the field of forensic science, and differentiate it from digital forensics.

## PART IV: SOURCE CODE DIFFERENTIATION

This part describes source code differentiation, a very basic method of comparing and measuring software source code. Source code differentiation is especially useful for finding code that has been directly copied from one program to another and for determining a percentage of direct copying. While there are many metrics for measuring qualities of software, source code differentiation has some unique abilities to measure development effort, software changes, and software intellectual property changes that are particularly useful for determining software intellectual property value for such applications as transfer pricing calculations.

In this part I introduce the mathematics of the theory of source code differentiation and explain implementations of source code differentiation for programmers who want to understand how to implement it. I also describe the "changing lines of code" or "CLOC" method of measuring software growth that is based on source code differentiation, and I compare it to traditional methods like "source lines of code" or "SLOC." I then discuss various applications of source code differentiation, though I believe that many more applications of this metric will be found in the future.

## PART V: SOURCE CODE CORRELATION

This part starts by exploring the various methods and algorithms for "software plagiarism detection" that have been developed over the last few decades.

I describe the origins of these methods and algorithms, and I explain their limitations. In particular, there have been no standard definitions and no supporting theory for this work, so I introduce the theory of source code correlation and definitions for characterizing source code. This characterization of software source code is practical for determining correlation and, ultimately, for determining whether copying occurred. While the theory and definitions are broad enough to be useful in various areas of computer science, they are particularly valuable in litigation.

In this part I also describe practical implementations of the theory for those programmers who want to understand how to implement the algorithms. Additionally, I describe applications of the theory in the real world. This part is highly mathematical, though the chapter on source code characterization will be useful for lawyers in understanding how elements of software source code can be categorized, how these various elements relate, and how the elements can affect a software copyright infringement, software trade secret, or software patent case.

## PART VI: OBJECT AND SOURCE/OBJECT CODE CORRELATION

In this part I introduce the theory and mathematics of object code correlation, which is used to compare object code to object code to find signs of copying. I also introduce the theory of source/object code correlation, which is used to compare source code to object code to find signs of copying. Both of these correlation measures are helpful before litigation when there is no access to source code from at least one party's software. I also describe practical implementations of the theory for those programmers who want to understand how to implement these correlation measures, and I describe applications of the methods and algorithms in the real world.

## PART VII: SOURCE CODE CROSS-CORRELATION

In this part I introduce the theory and mathematics of source code cross-correlation, which is specifically used to compare functional source code statements to nonfunctional source code comments to find signs of copying. This correlation measure is effective, in certain cases, for finding copied code that has been disguised enough to avoid detection with one of the other correlation measures. I describe some ways of effectively implementing code to

measure source code cross-correlation for those programmers who want to understand how to implement this measure, and I describe applications of source code cross-correlation in the real world.

## PART VIII: DETECTING SOFTWARE IP THEFT AND INFRINGEMENT

All of the correlation measures described in previous parts are useful for detecting software intellectual property theft; however, expert review is still required. Previously developed algorithms often produced a measure that claimed to show whether code was copied or not. In reality, a mathematical measure in and of itself is not enough to make this determination, and that is one of the problems with previous work in this area. In this part I describe detailed, precise steps to be taken once correlation has been calculated. These steps are as important to the standardization and objectivity required for determining intellectual property theft and infringement as are the various correlation measurements described in the previous parts.

## PART IX: MISCELLANEOUS TOPICS

This part covers areas that have come up in my involvement with intellectual property litigation. These subjects were also suggested by some of the experienced reviewers of this book who felt they deserved discussion. The issues described in this part often arise in software intellectual property cases and are also important for code developers and managers to understand. In particular, I discuss procedures for implementing a software clean room, I explain open source code, and I describe the Digital Millennium Copyright Act.

## PART X: PAST, PRESENT, AND FUTURE

The topics discussed in this book are cutting-edge, and I find them to be very interesting and exciting. A lot of work remains to be done, including extending the theories, advancing the mathematics, standardizing the definitions, and promoting the methodologies. In this part I discuss what has been done to date, speculate on areas of future research that build on the concepts in this book, and look toward new applications in various aspects of law and computer science.

*This page intentionally left blank*

# INDEX