# Chapter 36 Using Kolmogorov Complexity to Study the Coevolution of Header Files and Source Files of C-Alike Programs

#### Liguo Yu

Indiana University South Bend, Computer Science Department, South Bend, USA

### ABSTRACT

In C-alike programs, the source code is separated into header files and source files. During the software evolution process, both these two kinds of files need to adapt to changing requirement and changing environment. This paper studies the coevolution of header files and source files of C-alike programs. Using normalized compression distance that is derived from Kolmogorov complexity, we measure the header file difference and source file difference between versions of an evolving software product. Header files distance and source files distance are compared to understand their difference in pace of evolution. Mantel tests are performed to investigate the correlation of header file evolution and source file evolution. The study is performed on the source code of Apache HTTP web server.

#### INTRODUCTION

In C, C++, or Objective-C programs, there are two types of code files: header files and source files. Header files contain the declarations of classes, functions, global variables, and identifiers. Source files contain the implementations of classes and functions and the usages of variables and identifiers. The purpose of this organization is to separate program structure and interface from specific implementations (Stroustrup, 2013). On the other hand, software evolution is inevitable. Software products need to adapt to new requirements and new environments (Bendifallah, 1987; Gall et al., 1997; Kemerer & Slaughter, 1999; Lehman, 1980; Yu & Ramaswamy, 2009b; Yu & Ramaswamy, 2006). In this evolution process, both the program structure (interface) and the program implementation will need to be changed.

DOI: 10.4018/978-1-7998-3016-0.ch036

Because header files contain the program structure and interface, they are expected to be more robust than source files to requirement changes and environment changes. However, there has been no empirical research to validate this speculation. The objective of this study is to understand and compare the evolution of header files and source files. Using normalized compression distance that is derived from Kolmogorov complexity, we measure the difference of versions of header files and difference of versions of source files. Our study is performed on the source code of Apache HTTP, a web server written in C. Through this study, our objective is to answer the following two questions:

- 1. In the source code of Apache HTTP, do header files and source files have similar paces of evolution?
- 2. In the source code of Apache HTTP, is the evolution of header files correlated with the evolution of source files?

To answer these two questions, the remainder of the paper is organized as follows. Section 2 reviews Kolmogorov complexity and normalized compression distance. Section 3 describes our research method and research data. Section 4 presents the results of this study. Conclusions appear in Section 5.

# KOLMOGOROV COMPLEXITY AND NORMALIZED COMPRESSION DISTANCE

In algorithmic information theory, the Kolmogorov complexity of an object, such as a string (a piece of text), is the length of the string's shortest description in some fixed universal description language (Li & Vitányi, 2009). The Kolmogorov complexity of a binary string x can be represented by K(x) denoting the shortest length of string x described using a universal language.

Based on Kolmogorov complexity, Bennett et al. (1998) defined information distance between two binary strings x and y as the following.

$$E(x,y) = \max\{K(x|y), K(y|x)\}$$
(1)

In Equation 1, K(x|y) denotes the conditional Kolmogorov complexity of string x given string y. Information distance E(x,y) measures the absolute distance between two strings x and y and thus does not reflect the relative difference between the two strings. Li et al. [10] then defined normalized information distance between two binary strings x and y.

$$d(x,y) = \frac{\max\left\{ (K(x|y), K(y|x)\right\}}{\max\left\{ K(x), k(y)\right\}}$$
(2)

Equation 2 is based on Kolmogorov complexity, which is non-computable in practice. To make the information distance more practical to use, Li et al. (2004) further substituted Kolmogorov complexity K(x) of a string x with the optimal compressed length C(x) of string x. The normalized compression distance between two binary strings x and y is defined below (Li et al., 2004).

9 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: <u>www.igi-global.com/chapter/using-kolmogorov-complexity-to-study-the-</u> coevolution-of-header-files-and-source-files-of-c-alike-programs/261055

# **Related Content**

#### Design Patterns for Social Intelligent Agent Architectures Implementation

Manuel Kolp, Yves Wauteletand Samedi Heng (2021). Research Anthology on Recent Trends, Tools, and Implications of Computer Programming (pp. 294-319).

www.irma-international.org/chapter/design-patterns-for-social-intelligent-agent-architectures-implementation/261032

# Software Engineering and New Emerging Technologies: The Involvement of Users for Development Applications for Tablets

Sergio Ricardo Mazini (2018). Computer Systems and Software Engineering: Concepts, Methodologies, Tools, and Applications (pp. 1586-1609).

www.irma-international.org/chapter/software-engineering-and-new-emerging-technologies/192937

#### Quality of Service in SDN Technology

Ankur Dumka (2018). Innovations in Software-Defined Networking and Network Functions Virtualization (pp. 195-216).

www.irma-international.org/chapter/quality-of-service-in-sdn-technology/198199

# CSE as Epistemic Technologies: Computer Modeling and Disciplinary Difference in the Humanities

Matt Ratto (2012). Handbook of Research on Computational Science and Engineering: Theory and Practice (pp. 567-586).

www.irma-international.org/chapter/cse-epistemic-technologies/60375

#### Architecture and Implementation Issues

Ajantha Dahanayake (2001). Computer-Aided Method Engineering: Designing CASE Repositories for the 21st Century (pp. 95-137).

www.irma-international.org/chapter/architecture-implementation-issues/6876