

# Chapter 1.9

## FLOSSmole:

### A Collaborative Repository for FLOSS Research Data and Analyses

**James Howison**

*Syracuse University, USA*

**Megan Conklin**

*Elon University, USA*

**Kevin Crowston**

*Syracuse University, USA*

#### ABSTRACT

This article introduces and expands on previous work on a collaborative project, called FLOSSmole (formerly OSSmole), designed to gather, share, and store comparable data and analyses of free, libre, and open source software (FLOSS) development for academic research. The project draws on the ongoing collection and analysis efforts of many research groups, reducing duplication, and promoting compatibility both across sources of FLOSS data and across research groups and analyses. The article outlines current difficulties with the current typical quantitative FLOSS research process and uses these to develop requirements and presents the design of the system.

#### INTRODUCTION

This article introduces a collaborative project called FLOSSmole,<sup>1</sup> designed to gather, share, and store comparable data and analyses of free and open source software development for academic research. The project draws on the ongoing collection and analysis efforts of many research groups. Our intent in developing FLOSSmole is to reduce duplication, and to promote compatibility both across sources of FLOSS data and across research groups and analyses.

Creating a collaborative data and analysis repository for research on FLOSS is important because research should be as reproducible, extendable, and comparable as possible. Research with these characteristics creates the opportunity to employ meta-analyses, exploiting the diversity of existing research by comparing and contrasting results to expand our knowledge. Unfortunately, the current typical FLOSS research project pro-

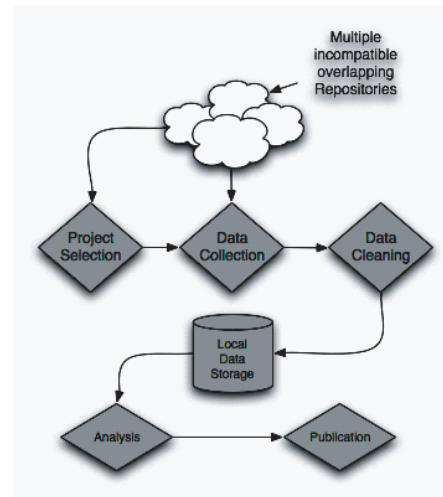
ceeds in a way that does not necessarily achieve these goals. These goals require detailed communal knowledge of the many choices made throughout a research project. Traditional publication prioritizes results, but masks or discards much of the information needed to understand and exploit the differences in our data collection and analysis methodologies. FLOSSmole was originally designed to provide resources and support to academics seeking to prepare the next generation of FLOSS research. Since its inception, FLOSSmole has also been a valuable resource for nonacademics who are also seeking good data about development practices in the open source software industry.

## BACKGROUND OF PROBLEM

Obtaining data on FLOSS projects is both easy and difficult. It is easy because FLOSS development utilizes computer-mediated communications heavily for both development team interactions and for storing artifacts such as code and documentation. This way of developing software leaves a freely available and, in theory at least, highly accessible trail of data upon which many academics have built interesting analyses about optimal organization of development teams, economics of building software in the commons, and the like. Yet, despite this presumed plethora of data, researchers often face significant practical challenges in using this data to construct a collaborative and deliberative research discourse. In Figure 1, we outline the research process we believe is followed in much of the quantitative literature on FLOSS.

The first step in collecting online FLOSS data is selecting which projects and which attributes to study, two techniques often used in estimation and selection are census and sampling. (Case studies are also used but these will not be discussed in this article.)

*Figure 1. The typical quantitative FLOSS research process (notice its noncyclical and noncollaborative nature)*



Conducting a census means to examine all cases of a phenomena, taking the measures of interest to build up an entire accurate picture. Taking a census is difficult in FLOSS for a number of reasons. First, it is hard to know how many FLOSS projects there are “out there,” and it is hard to know which projects should actually be included. For example, are corporate-sponsored projects part of the phenomenon or not? Do single-person projects count? What about school projects?

Second, the projects themselves, and the records they leave, are scattered across a surprisingly large number of locations. It is true that many are located in the major general repositories, such as Sourceforge<sup>2</sup> and GNU Savannah.<sup>3</sup> It is also true, however, that there are a number of other repositories of varying sizes and focuses (e.g., CodeHaus,<sup>4</sup> CPAN<sup>5</sup>), and that many projects, including the well-known and much-studied Apache and Linux projects, prefer to use their own repositories and their own tools. This diversity of location effectively hides significant portions of the FLOSS world from attempts at census. Even if a full listing of projects and their locations

8 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:

[www.igi-global.com/chapter/flossmole-collaborative-repository-floss-research/29381](http://www.igi-global.com/chapter/flossmole-collaborative-repository-floss-research/29381)

## Related Content

---

### Machine Learning for Designing an Automated Medical Diagnostic System

Ahsan H. Khandoker and Rezaul K. Begg (2009). *Handbook of Research on Modern Systems Analysis and Design Technologies and Applications* (pp. 544-559).

[www.irma-international.org/chapter/machine-learning-designing-automated-medical/21087](http://www.irma-international.org/chapter/machine-learning-designing-automated-medical/21087)

### SentiNeg: Algorithm to Process Negations at Sentence Level in Sentiment Analysis

Sandhya R. Savanur and R. Sumathi (2023). *International Journal of Software Innovation* (pp. 1-27).

[www.irma-international.org/article/sentineg/315741](http://www.irma-international.org/article/sentineg/315741)

### A Systematic Literature Review on Test Case Prioritization Techniques

Harendra Singh, Laxman Singh and Shailesh Tiwari (2022). *International Journal of Software Innovation* (pp. 1-36).

[www.irma-international.org/article/a-systematic-literature-review-on-test-case-prioritization-techniques/312263](http://www.irma-international.org/article/a-systematic-literature-review-on-test-case-prioritization-techniques/312263)

### Impulse Noise Detection and Removal Method Based on Modified Weighted Median

Ashpreet and Mantosh Biswas (2020). *International Journal of Software Innovation* (pp. 38-53).

[www.irma-international.org/article/impulse-noise-detection-and-removal-method-based-on-modified-weighted-median/248529](http://www.irma-international.org/article/impulse-noise-detection-and-removal-method-based-on-modified-weighted-median/248529)

### The Role of Compliance and Conformance in Software Engineering

José C. Delgado (2014). *Handbook of Research on Emerging Advancements and Technologies in Software Engineering* (pp. 392-420).

[www.irma-international.org/chapter/the-role-of-compliance-and-conformance-in-software-engineering/108627](http://www.irma-international.org/chapter/the-role-of-compliance-and-conformance-in-software-engineering/108627)