

Chapter 102

Open Source Software Development Challenges: A Systematic Literature Review on GitHub

Abdulkadir Seker

Sivas Cumhuriyet University, Turkey

Banu Diri

Yıldız Technical University, Turkey

Halil Arslan

Sivas Cumhuriyet University, Turkey

Mehmet Fatih Amasyalı

Yıldız Technical University, Turkey

ABSTRACT

GitHub is the most common code hosting and repository service for open-source software (OSS) projects. Thanks to the great variety of features, researchers benefit from GitHub to solve a wide range of OSS development challenges. In this context, the authors thought that was important to conduct a literature review on studies that used GitHub data. To reach these studies, they conducted this literature review based on a GitHub dataset source study instead of a keyword-based search in digital libraries. Since GHTorrent is the most widely known GitHub dataset according to the literature, they considered the studies that cite this dataset for the systematic literature review. In this study, they reviewed the selected 172 studies according to some criteria that used the dataset as a data source. They classified them within the scope of OSS development challenges thanks to the information they extract from the metadata of studies. They put forward some issues about the dataset and they offered the focused and attention-grabbing fields and open challenges that we encourage the researchers to study on them.

DOI: 10.4018/978-1-6684-3702-5.ch102

INTRODUCTION

Thanks to distributed version control systems such as Git, Mercurial, etc., open-source development platforms have reached a considerable number of users. The most common of these platforms is GitHub (based on git). GitHub has become the world's largest code server with more than 40 million developers hosting and collaborating over 100 million repositories.

On platforms such as GitHub, the development process is distributed. Developers can participate in a project, contribute, discuss bugs with each other, and write comments about code from various locations. In this way, a considerable amount of textual, numerical and network or collaboration-based features about the projects and developers are extracted from the platform. Besides, GitHub includes many social relations among users or projects. GitHub is the most common code hosting and repository service for open-source software projects. For the researchers that focus on software engineering, the content of this platform provides many valuable sources. Most of the studies about this domain use GitHub as a data source because of easy to access, amount of data, and diversity of features. In this context, we think that is important to conduct a literature review on studies that used GitHub data.

There are several options to reach GitHub data. In a survey study which is given the usage rates of GitHub dataset, they addressed that the most used dataset is GHTorrent (34%) in the articles that are reviewed according to the certain criteria (Cosentino, Luis, & Cabot, 2016). In Cosentino's systematic mapping study, the GHTorrent dataset is in the lead with a 41% use rate (Badashian, Shah, & Stroulia, 2015; Cosentino, Canovas Izquierdo, & Cabot, 2017). In the another study, GHTorrent is the most cited dataset (Kotti & Spinellis, 2019). The GHTorrent dataset was developed by Georgios Gousios in the software engineering department at Delft University of Technology (Gousios, 2013). The dataset is generated by systematically crawling with the GitHub API and includes information about all public projects and users on the platform. GHTorrent stores some information about repositories, projects, issue descriptions, comments, and pull request (PR) conversations in 26 relational tables totally.

We saw from other systematic literature review (SLR) papers that some studies can be missed when reviewing with a text-based (keyword) search from search engines or digital libraries. Because of that, to reach the studies, we conducted this literature review based on a GitHub dataset source study instead of a keyword-based search in digital libraries. Due to GHTorrent is the most widely known and used GitHub dataset according to the literature, we considered the studies which cite this dataset for the systematic literature review.

In this study, we offered to find out the topics of all studies and classified them. We focused on the studies with the context of open-source software development. We divided the studies into some categories and challenges. Besides, some distributions (type, venue, year, method, data, topic) have been obtained from the studies that used the dataset. We show which challenges are mentioned in the studies and how each study is using the dataset. Thus, we hope the study guided the researchers who interest in software engineering challenges with open-source systems. We formed this review following these research questions:

RQ1: What are the trends of open-source software development challenges?

RQ2: What are the handicaps/cons of GHTorrent?

RQ3: What are the open challenges that have not yet been studied with this dataset?

29 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/open-source-software-development-challenges/294562

Related Content

Vojta-Therapy: A Vision-Based Framework to Recognize the Movement Patterns

Muhammad Hassan Khan and Marcin Grzegorzek (2017). *International Journal of Software Innovation* (pp. 18-32).

www.irma-international.org/article/vojta-therapy/182534

Big Data: The Path to Maturity

Stephen H. Kaisler, William H. Money, Frank Armour and J. Alberto Espinosa (2017). *International Journal of Systems and Service-Oriented Engineering* (pp. 1-23).

www.irma-international.org/article/big-data/190410

A Novel Sentence Completion System for Punjabi Using Deep Neural Networks

Gurjot Singh Mahi and Amandeep Verma (2022). *International Journal of Software Innovation* (pp. 1-25).

www.irma-international.org/article/a-novel-sentence-completion-system-for-punjabi-using-deep-neural-networks/293271

Mathematical Models of Video-Sequences of Digital Half-Tone Images

E.P. Petrov, I.S. Trubin, E.V. Medvedeva and S.M. Smolskiy (2013). *Integrated Models for Information Communication Systems and Networks: Design and Development* (pp. 207-241).

www.irma-international.org/chapter/mathematical-models-of-video-sequences-of-digital-half-tone-images-/79666

Malware Forensics: An Application of Scientific Knowledge to Cyber Attacks

C. V. Suresh Babu, G. Suruthi and C. Indhumathi (2023). *Malware Analysis and Intrusion Detection in Cyber-Physical Systems* (pp. 285-312).

www.irma-international.org/chapter/malware-forensics/331309