X-COBOL: A Dataset of COBOL Repositories

Mir Sameed Ali, Nikhil Manjunath, Sridhar Chimalakonda*
Research in Intelligent Software & Human Analytics (RISHA) Lab
Department of Computer Science & Engineering
Indian Institute of Technology Tirupati, India
{cs18b021,cs18b041,ch}@iittp.ac.in

ABSTRACT

Despite being proposed as early as 1959, COBOL (Common Business-Oriented Language) still predominantly acts as an integral part of the majority of operations of several financial, banking, and governmental organizations. To support the inevitable modernization and maintenance of legacy systems written in COBOL, it is essential for organizations, researchers, and developers to understand the nature and source code of COBOL programs. However, to the best of our knowledge, we are unaware of any dataset that provides data on COBOL software projects, motivating the need for the dataset. Thus, to aid empirical research on comprehending COBOL in open-source repositories, we constructed a dataset of 84 COBOL repositories mined from GitHub, containing rich metadata on the development cycle of the projects. We envision that researchers can utilize our dataset to study COBOL projects' evolution, code properties and develop tools to support their development. Our dataset also provides 1255 COBOL files present inside the mined repositories. The dataset and artifacts are available at https://doi.org/10.5281/zenodo.7968845.

CCS CONCEPTS

• Software and its engineering \rightarrow General programming languages; Software libraries and repositories.

KEYWORDS

COBOL, Dataset, GitHub, Mining Software Repositories

1 INTRODUCTION

In StackOverflow Developer Survey 2022, programming languages such as Javascript, Python, Java, and C# were found to be the most popular among developers. On the other hand, COBOL was the second least popular language¹. Although COBOL is generally considered outdated, it is still being used in 80% of financial services transactions, 95% of ATM swipes, and there are about 220 billion lines of code, with 1.5 billion lines written every year [22]. COBOL is still processing about USD 3 trillion in commerce every day [4]. Apart from the financial sector, legacy systems written in COBOL are used in other major sectors such as healthcare and governmental institutions [22].

Considering that legacy systems written in COBOL are still integral to vital sectors, it is essential to either maintain them or migrate them to modern systems. Migrating COBOL legacy systems to systems with modern technology is an effort-intensive work associated with an enormous amount of cost and risk. When Commonwealth Bank of Australia migrated their COBOL systems, it took them five

Datasets have been created to support empirical research for different programming languages and to understand and address challenges in several software engineering areas. Eskandani et al. constructed a dataset on open-source Serverless applications mined from GitHub [10]. In order to support empirical research in the domain of game engines, Vagavolu et al. presented a dataset of 526 game engine repositories mined from GitHub [23]. While there are several datasets to support empirical research in other software engineering areas such as docker [11, 17], android application development [13], program equivalence [2], and so on, to the best of our knowledge, there exists no dataset that caters to COBOL projects. National Computing Centre of UK provides COBOL85 test suite², which is a set of COBOL programs containing different features. Although the test suite can be utilized to construct COBOL parsers and compilers, it cannot support extensive analysis of the development of COBOL projects.

Hence, to facilitate the empirical research in COBOL, we present a curated dataset of 84 COBOL projects mined from GitHub, consisting of 4420 *commits*, 241 *pull requests*, and 727 *issues*. The resultant projects have been manually analyzed and selected based on multiple parameters. Along with the projects, we also provide 1255 COBOL program files extracted from the selected repositories.

The remainder of this paper is organized as follows. Section 2 presents an overview of the data extraction methodology. Next, we present the dataset schema and the dataset statistics in Section 3. Then, we describe a set of research applications of our dataset in Section 4. Finally, we list a few limitations of our dataset and discuss the scope of future work in Section 5.

years and cost them \$749.9 million [22]. Maintaining COBOL systems is less viable than maintaining modern systems, primarily due to a shortage of experienced COBOL programmers [9]. Despite the enormous complications present in the maintenance and migration of COBOL projects, software engineering research on COBOL is limited. Most efforts are towards restructuring COBOL code [18], extracting knowledge and business rules [19], and supporting the migration of legacy COBOL systems and industrial case studies [8, 16, 20, 21]. Ciborowska et al. [7] surveyed differences in defects and defect location strategies in COBOL and modern programming languages by interviewing 30 COBOL and 74 modern programming language developers. The survey showed significant differences in defect types present in COBOL and modern programming language projects, with similar defect location strategies employed by the developers in both kinds of projects [7]. Opdebeeck et al. [14] presented an approach to mine library usage patterns in COBOL code, which can assist the migration process.

^{*}Corresponding Author

¹https://survey.stackoverflow.co/2022/

 $^{^2} https://www.itl.nist.gov/div897/ctg/cobol_form.htm$

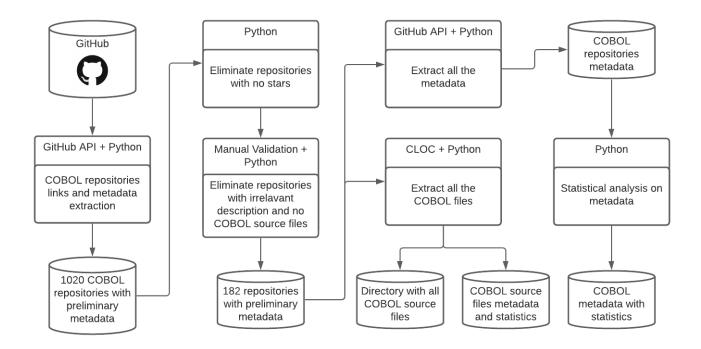


Figure 1: Data Collection Methodology

2 DATA COLLECTION

This section details the data extraction process, and the implementation details are shown in Figure 1. We have followed a similar methodology to the one that was used by Vagavolu et al. to collect a dataset on game engines [23]. We extract all the necessary repository data required for the preliminary elimination of repositories. We then eliminate the unwanted repositories using manual evaluation and a *python* script. We finally mine all the resultant repositories for metadata using the GitHub API and *python*, and COBOL files using *python* and *CLOC*³.

The following sections elaborate on the dataset creation process.

2.1 Preliminary Repository Metadata Extraction and Filtration

We accumulated the links of all the GitHub repositories whose primary language is COBOL to be included in the dataset. We extracted the names and metadata of all these repositories using GitHub API. We ended up with 1020 repositories. In order to have a quality dataset, we eliminated all the repositories with less than one star. We also removed the repositories whose topics include the following keywords ["list", "course", "resource", "tutorial", "learning", "exercise", "example"]. At this stage, we eliminated 628 repositories and had 392 repositories to be further processed. After eliminating repositories using GitHub metrics, some repositories containing undesirable data, such as no COBOL source files and irrelevant descriptions, still existed in the dataset. One such example was SaymanNsk/Sprinter200x, which contained compiled binaries of

COBOL instead of COBOL source code. We eliminated all the false positives by a semi-automated technique validating every repository after the first elimination using the type of code and repository description. We consider the repositories related to tutorials or basic COBOL learning irrelevant for this dataset and manually filtered out 210 repositories. The final resultant repository links and names are stored in a CSV file for further data extraction. We finally have 84 repositories in the dataset.

2.2 Extraction of metadata

After eliminating false positives in two levels, we decided to use *python* and the GitHub API to extract all the repositories' metadata using multiple GitHub access keys. For each repository, we extract metadata of 7 different categories as shown in Figure 2. We also performed statistical analysis on the collected repositories calculating measures such as *Max*, *Mean*, and *Total* for a set of metrics. The results of the analysis are listed in Table 1.

All the different category results were stored in separate CSV files. The number of repositories and the response for different API requests were less, which motivated us to dump data into simple CSV files instead of any database solutions. Detailed information about different metadata sections is mentioned in Database Schema.

2.3 Extraction of COBOL source files

We extracted all the COBOL files present in the selected 84 GitHub repositories⁴. We implemented a script that iterates over the selected repositories. During an iteration, a repository is cloned, and

³http://cloc.sourceforge.net

⁴The dataset can be found at https://doi.org/10.5281/zenodo.7968845

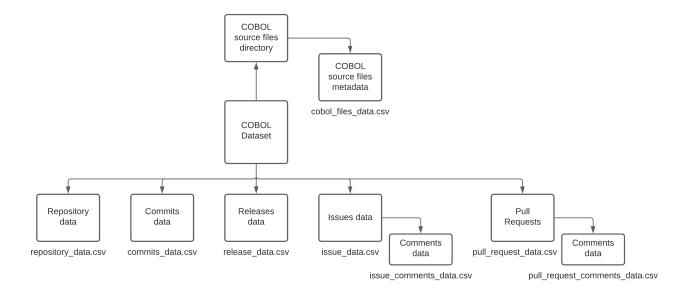


Figure 2: Dataset Schema

CLOC is used to identify all COBOL files present in the cloned repository. CLOC also provides metadata (File path, number of Blank lines, number of Comment lines, and number of Code lines) for the identified COBOL files with the extensions CBL, cbl, ccp, COB, cob, cobol, and cpy. All the recognized COBOL files and the CLOC result of the cloned repository are stored in a directory named AuthorName_RepositoryName in the COBOL files directory. We also retrieve the number of commits made for each COBOL file using the GitHub API.

3 THE DATASET

3.1 Dataset Schema

The dataset contains eight CSV files, capturing different properties of the dataset, and a directory named *COBOL_Files* containing the extracted COBOL files. The overall schema demonstrating the type of data present in the eight CSV files is shown in Figure 2.

The *repository_data.csv* provides the overall activity of the repositories using metrics such as commit frequency, merge frequency, commits, and others. repository_data.csv also contains repository metadata such as forks, stars, and size, along with metrics capturing the overall COBOL files content present in the repository such as total COBOL files, total blank lines, and total code lines. commits_data.csv has the data on all the commits made to the selected repositories providing information such as commit date, commit message, and commit author. Information on the releases made in the selected repositories is stored in release_data.csv. issue_data.csv includes data on the issues created for the chosen repositories. The comments made on the issues are stored in issue comments data.csv. pull request data.csv and pull request comments data.csv contain data on the pull requests made and their trailing comments, respectively. cobol_files_data.csv contains the metadata on the COBOL files stored in the COBOL_Files directory and has information such

Table 1: Statistics of Metadata of the X-COBOL Dataset

Metric Name	Max	Mean	Total
Forks	40	4.31	362
Size(KB)	203750	5313.13	446303
Releases	12	0.24	20
Open issues	205	3.39	285
Closed issues	138	5.26	442
Open pull requests	3	0.15	13
Closed pull requests	52	2.71	228
Commits	504	52.62	4420
Commit frequency	29.0	3.41	286.08
Committer frequency	2	0.49	41.18
Integration frequency	38.0	3.69	309.73
Integrator frequency	2	0.54	45.11
Merge frequency	3.0	0.20	17.11
Number of cobol files	117	14.94	1255
Total comment lines	9091	500.68	42057
Total code lines	23646	2397.37	201379
Comment lines per cobol file	664	37.98	3190.16
Code lines per cobol file	3320	241.38	20275.70
Commits per cobol file	302	10.63	893.61

as *blank lines*, *comment lines*, *code lines*, and *commits* made to the file. We see that this information could be potentially leveraged for analyzing COBOL projects based on varying metadata.

3.2 Dataset Statistics

The dataset contains metadata and COBOL files of 84 GitHub repositories. The number of other languages used along with COBOL in the curated dataset is 61, with *shell*, *makefile*, and *C* most used. The

average number of commits made for a repository is 40, with 504 commits being the highest for *debinix/openjensen*. The total number of COBOL files present in the dataset is 1255. On average, there are 21 COBOL files present in a repository. *abrignoli/COBSOFT* has the maximum number of COBOL files which is 117. The average number of comments and code lines present in a single COBOL file is 729 and 1602, respectively.

4 POTENTIAL DATASET USAGE OR APPLICATIONS

In this section, we list some of the use cases of the X-COBOL dataset.

4.1 Usability of COBOL Constructs

The usability of standard constructs of a programming language is critical and has been studied in the context of programming language design. Recently, Peng et al. performed an empirical study on the usability of different features of *Python* [15]. Al-Jarrah et al. conducted a similar study by analyzing 340 COBOL programs in 1979 [1]. We believe that an equivalent study on a more extensive set of COBOL programs can be performed using the COBOL source files present in the X-COBOL dataset. Using the results, we can understand the usage patterns of different COBOL constructs, which can aid in optimizations in compilation and migration systems. Further, the inexperienced COBOL developers can benefit from this study by adopting the commonly used patterns of popular COBOL constructs.

4.2 Understanding Bugs, Issues, Commits, Pull Requests in COBOL

Metadata of open source projects such as *bug reports, issue reports, commit messages*, and *pull requests* have been analyzed to understand the causes and attributes of bugs present in several application environments and investigate the strategies employed to detect them [5, 25]. An analogous study can be conducted using X-COBOL to comprehend the characteristics of bugs present in COBOL projects. Litecky et al. examined the types of errors and their occurrence frequencies using COBOL programs in 1976 [12]. This study can be extended by incorporating the project metadata such as *issues, commits*, and *pull request messages* present in X-COBOL to detect the cause of bugs identified. Furthermore, techniques used to locate and rectify the identified bugs can be recognized. This analysis can benefit the COBOL developers in debugging and maintaining COBOL projects.

4.3 COBOL Open Source Analysis

R. van Wendel de Joode et al. [24] have shown that the reliability of open-source software increases with an increase in software usage, with a transparent flow of information between developers and the popularity of the software. The X-COBOL dataset could be used to analyze the reliability of the open-source COBOL software in the current age, considering that COBOL is a legacy programming language. Also, following from the analysis of Choudhary et al. [6], it is essential to analyze the developer collaboration to measure productivity and evolution of software. Researchers can analyze these patterns in open-source COBOL software using the X-COBOL dataset.

Further, this dataset can also measure the open-source community support, interest, and growth concerning a legacy programming language, COBOL.

In addition, we see immense scope for conducting empirical research on COBOL language and legacy systems in similar lines to other languages like *Python*, *Java* and different kinds of systems such as deep learning, games and so on.

5 LIMITATIONS AND FUTURE WORK

We have used GitHub stars count to determine and get quality of preliminary repositories, but stargazers count, and fork count are not the only metrics to obtain a quality dataset [3]. Also, due to the manual nature of the evaluation, there might be repositories that contain data which might not be up to the intended quality. We plan to include other good quality repositories among 210 repositories eliminated and remove irrelevant repositories in the near future. Due to the small dataset size, we currently provide the dataset in different csv files and the COBOL source files as a zipped file. We plan to provide the dataset in different formats, including storing the metadata and source files in a database. Furthermore, we also intend to improve the dataset quality by evaluating the COBOL files by executing them.

6 CONCLUSION

We have presented a curated dataset containing structured metadata about the development cycle of 84 COBOL projects mined from GitHub. The dataset includes metadata on the *commits, issues, pull requests,* and *releases* of the mined repositories, along with the COBOL source files present in them. Additionally, we provide the metadata of COBOL files extracted. We expect the research community to utilize the dataset on COBOL projects to conduct empirical studies on code quality, COBOL software development practices, error analysis, security, and so on. Also, the dataset can aid in research studies and tools supporting the maintenance and migration of COBOL projects. Finally, the extracted COBOL source files can be used by researchers to perform static code analysis to develop tools that support the development of COBOL projects.

7 ACKNOWLEDGEMENTS

We would like to thank Exafluence Inc. for supporting us throughout this project.

REFERENCES

- M. M. Al-Jarrah and I. S. Torsun. 1979. An Empirical Analysis of COBOL Programs. Softw. Pract. Exp. 9, 5 (1979), 341–359. https://doi.org/10.1002/spe.4380090502
- [2] Sahar Badihi, Yi Li, and Julia Rubin. 2021. EqBench: A Dataset of Equivalent and Non-equivalent Program Pairs. In 2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR). 610–614. https://doi.org/10.1109/MSR52588. 2021.00084
- [3] Hudson Borges and Marco Tulio Valente. 2018. What's in a GitHub Star? Understanding Repository Starring Practices in a Social Coding Platform. Journal of Systems and Software 146 (Dec 2018), 112–129. https://doi.org/10.1016/j.jss.2018. 09.016
- [4] David Cassel. 2017. COBOL Is Everywhere. Who Will Maintain It? Retrieved January 13, 2022 from https://thenewstack.io/cobol-everywhere-will-maintain/
- [5] Gemma Catolino, Fabio Palomba, Andy Zaidman, and Filomena Ferrucci. 2019. Not All Bugs Are the Same: Understanding, Characterizing, and Classifying the Root Cause of Bugs. arXiv:1907.11031 [cs.SE]
- [6] Samridhi Choudhary, Christopher Bogart, Carolyn Rose, and Jim Herbsleb. 2020. Using Productive Collaboration Bursts to Analyze Open Source Collaboration

- Effectiveness. In 2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER). 400–410. https://doi.org/10.1109/SANER48275. 2020.9054852
- [7] Agnieszka Ciborowska, Aleksandar Chakarov, and Rahul Pandita. 2021. Contemporary COBOL: Developers' Perspectives on Defects and Defect Location. 2021 IEEE International Conference on Software Maintenance and Evolution (ICSME) (Sep 2021). https://doi.org/10.1109/icsme52107.2021.00027
- [8] Alessandro De Marco, Valentin Iancu, and Ira Asinofsky. 2018. COBOL to Java and Newspapers Still Get Delivered. 2018 IEEE International Conference on Software Maintenance and Evolution (ICSME) (Sep 2018). https://doi.org/10.1109/icsme. 2018.00055
- [9] John Delaney. 2020. COBOL Programmers are Back In Demand. Seriously. Retrieved January 13, 2022 from https://cacm.acm.org/news/244370-cobol-programmers-are-back-in-demand-seriously/fulltext
- [10] Nafise Eskandani and Guido Salvaneschi. 2021. The Wonderless Dataset for Serverless Computing. In 2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR). 565–569. https://doi.org/10.1109/MSR52588.2021. 00075
- [11] Jordan Henkel, Christian Bird, Shuvendu K. Lahiri, and Thomas Reps. 2020. A Dataset of Dockerfiles. Proceedings of the 17th International Conference on Mining Software Repositories (Jun 2020). https://doi.org/10.1145/3379597.3387498
- [12] Charles R. Litecky and Gordon B. Davis. 1976. A Study of Errors, Error-Proneness, and Error Diagnosis in Cobol. Commun. ACM 19, 1 (jan 1976), 33–38. https://doi.org/10.1145/359970.359991
- [13] Pei Liu, Li Li, Yanjie Zhao, Xiaoyu Sun, and John Grundy. 2020. AndroZooOpen: Collecting Large-Scale Open Source Android Apps for the Research Community. Association for Computing Machinery, New York, NY, USA, 548–552. https://doi.org/10.1145/3379597.3387503
- [14] Ruben Opdebeeck, Johan Fabry, Tim Molderez, Jonas De Bleser, and Coen De Roover. 2021. Mining for Graph-Based Library Usage Patterns in COBOL Systems. In 2021 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER). 595–599. https://doi.org/10.1109/SANER50967.2021. 00072
- [15] Yun Peng, Yu Zhang, and Mingzhe Hu. 2021. An Empirical Study for Common Language Features Used in Python Projects. In 2021 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER). 24–35. https://doi. org/10.1109/SANER50967.2021.00012

- [16] J. Rodriguez, M. Crasso, C. Mateos, A. Zunino, and M. Campo. 2013. Bottom-Up and Top-Down Cobol System Migration to Web Services. *IEEE Internet Computing* 17, 02 (mar 2013), 44–51. https://doi.org/10.1109/MIC.2011.162
- [17] Gerald Schermann, Sali Zumberi, and Jürgen Cito. 2018. Structured Information on State and Evolution of Dockerfiles on GitHub. In 2018 IEEE/ACM 15th International Conference on Mining Software Repositories (MSR). 26–29.
- [18] Alex Sellink, Harry Sneed, and Chris Verhoef. 2002. Restructuring of COBOL/CICS legacy systems. Science of Computer Programming 45, 2 (2002), 193– 243. https://doi.org/10.1016/S0167-6423(02)00061-8 Special Issue on Software Maintenance and Reengineering (CSMR 99).
- [19] H.M. Sneed. 2001. Extracting business logic from existing COBOL programs as a basis for redevelopment. In *Proceedings 9th International Workshop on Program Comprehension. IWPC 2001.* 167–175. https://doi.org/10.1109/WPC.2001.921728
- [20] Harry M. Sneed. 2010. Migrating from COBOL to Java. In 2010 IEEE International Conference on Software Maintenance. 1–7. https://doi.org/10.1109/ICSM.2010. 5609583
- [21] Harry M. Sneed and Katalin Erdoes. 2013. Migrating AS400-COBOL to Java: A Report from the Field. In 2013 17th European Conference on Software Maintenance and Reengineering. 231–240. https://doi.org/10.1109/CSMR.2013.32
- [22] Tom Taulli. 2020. COBOL Language: Call It A Comeback? Retrieved January 13, 2022 from https://www.forbes.com/sites/tomtaulli/2020/07/13/cobol-languagecall-it-a-comeback/?sh=7ed22be77d0f
- [23] Dheeraj Vagavolu, Vartika Agrahari, Sridhar Chimalakonda, and Akhila Sri Manasa Venigalla. 2021. GE526: A Dataset of Open-Source Game Engines. In 2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR). 605–609. https://doi.org/10.1109/MSR52588.2021.00083
- [24] R. van Wendel de Joode and M. de Bruijne. 2006. The Organization of Open Source Communities: Towards a Framework to Analyze the Relationship between Openness and Reliability. In Proceedings of the 39th Annual Hawaii International Conference on System Sciences (HICSS'06), Vol. 6. 118b–118b. https://doi.org/10. 1109/HICSS.2006.477
- [25] Yuhao Zhang, Yifan Chen, Shing-Chi Cheung, Yingfei Xiong, and Lu Zhang. 2018. An Empirical Study on TensorFlow Program Bugs. In Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis (Amsterdam, Netherlands) (ISSTA 2018). Association for Computing Machinery, New York, NY, USA, 129-140. https://doi.org/10.1145/3213846.3213866