

Open Source Software Computed Risk Framework

Jon Chapman, Hari Venugopalan
University of California, Davis
Department of Computer Science
One Shields Avenue Davis, CA 95616
jwchapman@ucdavis.edu, hvenugopalan@ucdavis.edu

Abstract—The increased dissemination of open source software to a broader audience has led to a proportional increase in the dissemination of vulnerabilities. These vulnerabilities are introduced by developers, some intentionally or negligently. In this paper, we work to quantify the relative risk that a given developer represents to a software project. We propose using empirical software engineering based analysis on the vast data made available by GitHub to create a Developer Risk Score (DRS) for prolific contributors on GitHub. The DRS can then be aggregated across a project as a derived vulnerability assessment, we call this the Computational Vulnerability Assessment Score (CVAS). The CVAS represents the correlation between the Developer Risk score across projects and vulnerabilities attributed to those projects. We believe this to be a contribution in trying to quantify risk introduced by specific developers across open source projects. Both of the risk scores, those for contributors and projects, are derived from an amalgamation of data, both from GitHub and outside GitHub. We seek to provide this risk metric as a force multiplier for the project maintainers that are responsible for reviewing code contributions. We hope this will lead to a reduction in the number of introduced vulnerabilities for projects in the Open Source ecosystem.

Index Terms—Big Data, Computer Security, Prediction Methods, Data Analysis

I. INTRODUCTION

Advanced software that controls much of everyday life has increasingly been created in an open source, collaborative manner. Many in the open source software world believe that given enough eyeballs, all bugs (and vulnerabilities) are shallow. This idea begins to break down at the scale of modern software projects. We see this in the fact that despite there being unlimited access to the code base of major software projects such as Apache, the Linux Kernel, cURL, etc, vulnerabilities still proliferate. The rapid growth in the amount of code requires a different paradigm to proactively prevent vulnerabilities before bad actors are able to exploit them.

In the past decade, GitHub has emerged as the preeminent platform for social software engineering [1], both open and closed source. GitHub provides an online collaborative platform for software development projects, acting primarily as the central repository and dissemination point for the projects' code.

This also means that GitHub has unwittingly become a platform for the spread of security vulnerabilities. A lack of security awareness among developers [3], [4] has led to the inclusion of vulnerabilities. Even security conscious

developers can be responsible for creating vulnerable projects owing to their use of third-party libraries which could have vulnerabilities. [5] The lack of attention given to security, coupled with the difficulty in managing large scale projects that may have hundreds (or thousands!) of contributors increases the opportunity for the inclusion of vulnerabilities. There are also the cases of large scale projects that often underpin the working of the internet may only have a handful of developers working on them, but consist of many thousands of lines of code spread across hundreds of files, such as OpenSSL. These various challenges could even set the conditions for malicious users to purposely include vulnerable code in a project that has a large audience in order to exploit later, such as we saw with NPM. [8]

The creation of a Developer Risk Score (DRS), combined with the Computational Vulnerability Assessment Score (CVAS) for projects are in support of two closely related research questions:

- **RQ1:** Can an accurate predictive model be created to give insights into future risk of a given contributor?
- **RQ2:** Can a derived, composite risk index be used to accurately predict which open source software projects will be assigned a CVE in the near future?

II. PRIOR WORK

This is a brief overview that should serve to motivate our work as well as provide some context to our efforts. As we have seen no evidence of other researchers attempting to quantify the risk of specific developers using a multi-source data pooling approach we feel that this work is new and novel.

There are some tools out there that have similar aims to our project. They seek to prevent vulnerabilities before being distributed to users. These primarily center around static analysis techniques. Some advocate using static analysis tools [9] and others recommend increasing the skill and experience level of those doing code review [10]

Security in the software development process is often an afterthought [11] [12]. In many cases, security measures are taken after an exploit has been carried out, resulting in large amounts of damage [13]. Even for the tools that attempt to mitigate the vulnerability before it effects users require code integration. There is little in the way of predictive tools that try to get a head of the vulnerability before it makes it into the code base.

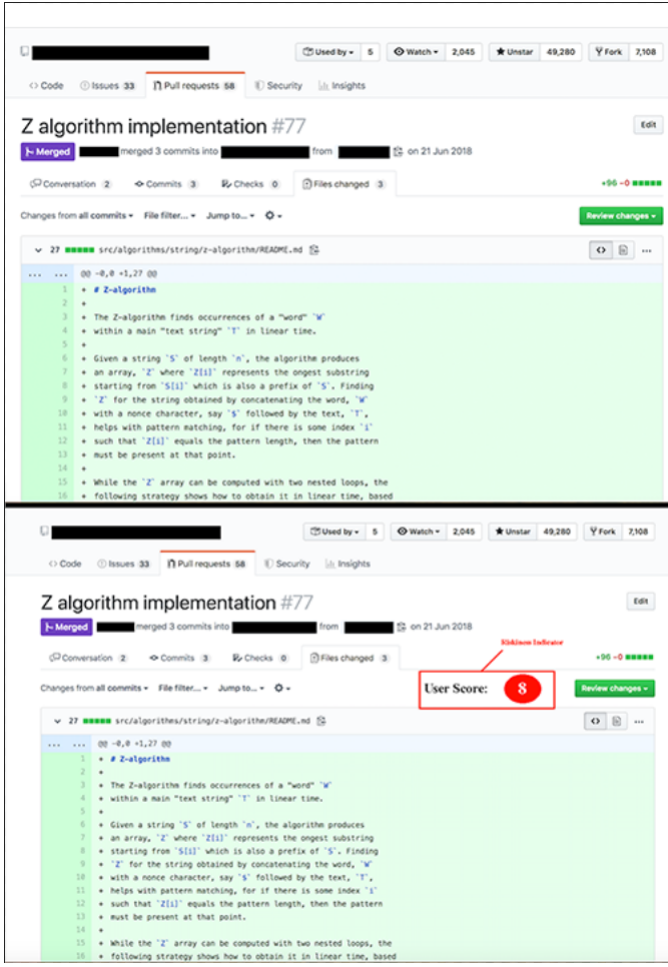


Fig. 1. The image on top shows the current pull request screen in GitHub. As can be seen, the pull request screen informs the user of only the contributors code. A link to the contributor’s GitHub profile is also part of the pull request which the reviewer can visit to get additional information about the contributor if required. However, doing so is tedious, and does not quantify the risk associated with the contributor in any way. The image below shows our proposed new pull request review screen with an associated user riskiness score. This score serves to indicate the riskiness associated with the contributor, which the reviewer can now use to decide the extent to which review is required for the pull request in the context of better security.

There are differing opinions of why this situation persists, from lack of awareness [14], to poor education or lack of security culture [15], and everything in between. Providing additional information to developers, and making it easy to create an assessment of code to be incorporated into a project will allow developers to make more secure choices.

Our approach is designed to give simple metrics by which decisions to integrate (or not integrate) code, or at least conduct a more thorough review of the code, can be made. This results in a lowering of the “cognitive load” of the project maintainers in the pursuit of increasing security. [16]

III. METHOD

There exists a bipartite many-to-many relationship between contributors and projects, since a particular contributor can contribute to multiple projects, and since multiple contributors

can contribute to the same project. We use this mapping, as shown in Figure 2 as the starting point of our analysis.

Each contributor is assigned a score for each project based on the number of vulnerabilities in the project weighted by their relative contribution to the project. This score is then aggregated for all projects the contributor has been involved with to calculate an effective score for the given contributor.

The sum of the scores of all contributors involved in a project would serve as the base score for the project itself. In future work we discuss this as a possible point of extension. In our study, we define the relative contribution of a developer to a project at a given point in time as the ratio of the number of commits the contributor has been involved in at that point in time (either as the author or as the committer) to the total number of commits that have been made to the project at that point.

We chose to go with this definition for the study, since the number of commits do serve as a direct proxy for the extent of a contributor’s involvement with the project and also since the commit data is readily available on GHTorrent. Including the time component in this definition is important and helps make the score more accurate since it ensures that vulnerabilities that existed in a project before a given developer started contributing to it are not attributed to the contributor. An example of this is shown in Figure 3. In the figure, Developer A is the sole contributor to the project till point t_2 when Developer B comes in, and they both contribute equally to the project from t_2 onwards. As per our metric, at any point in time up to t_3 , the riskiness score of Developer A to the project will be 10, and that of Developer B will be 0. After point t_3 , Developer A has contributed to 75% of the project on the whole, and Developer B has contributed 25%. As a result, after t_3 , Developer A’s riskiness score would be 17.5 and that of B would be 2.5

With the understanding so far, and an acknowledgement that an iterative design process to refine this scoring method is still required, we propose a risk scoring metric that is tentatively defined as follows:

Let,

$c_i(u_i, t_i, p_i)$ = Relative contribution of Contributor u_i , to Project p_i at time t_i

$r_i(u_i, t_i, p_i)$ = Partial Risk Score of Contributor u_i , derived from Project p_i at time t_i

Based on our definitions:

$$c_i(u_i, t_i, p_i) = \sum_{t=0}^{t_i} \frac{\text{num_commits}(u_i, p_i, t)}{\sum_{u \in p_i \text{ at time } t} \text{num_commits}(u, p_i, t)} \quad (1)$$

$$r_i(u_i, t_i, p_i) = c_i(u_i, t_i, P_i) \times (\text{cve_score}(p_i, t_i)) \quad (2)$$

Here, $\text{num_commits}(u_i, p_i, t_i)$ denotes the number of commits that have been made by contributor u_i to project p_i

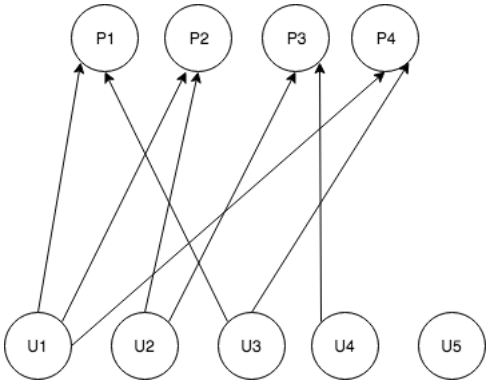


Fig. 2. Graph showing mapping between users and projects. The bottom nodes denote users, and those on top denote projects. Each user contributes to multiple projects and each project contains contributions from multiple users.

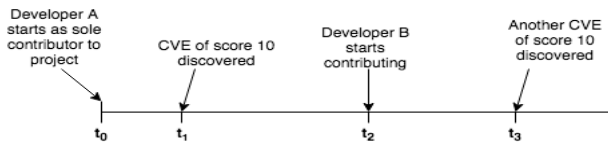


Fig. 3. This representative timeline demonstrates an example life cycle of a software development project on GitHub. This shows when CVEs were assigned to a project, and the points in time when different developers contributed to a particular project. As can be seen, assigning a score to Developer B for the CVE discovered before his/her contribution to the project would be inappropriate. Our model fully takes into account these sort of chronological considerations to ensure appropriate attribution. In this toy example, our model would assign Developer A with a risk score of 17.5 and Developer B with a risk score of 2.5.

at time t_i , and $cve_score(p_i, t_i)$ is the CVE score assigned to the project p_i at time t_i .

We can use this to construct an overall risk score of a given contributor at a specific time, in general terms:

$$R_i(u_i, t_i) = \text{Risk Score of Contributor } u_i \text{ at Time } t_i$$

We then use the prior definition to construct the formal definition of the individual developer risk score at a given point in time as the ratio of the sum of partial risks for the user across all their projects at that point in time to the total number of commits made by them at that point. More formally this can be expressed as:

$$R_i(u_i, t_i) = \frac{\sum_{\text{all } p_i \text{ that } u_i \text{ has contributed to}} r_i(u_i, t_i, p_i)}{\sum_{\text{all } p_i \text{ that } u_i \text{ has contributed to}} [\sum_{t=0}^{t_i} num_commits(u_i, p_i, t)]} \quad (3)$$

This risk score essentially conveys the risk associated with each commit made by the user, with a higher score indicating possibility of greater risk.

IV. RESULTS

In this section we describe our experiments and preliminary results. We first describe and report the results of a simple

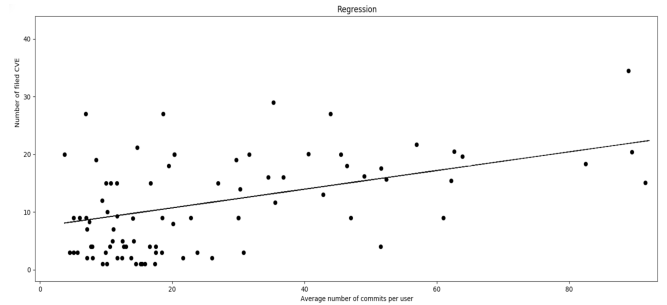


Fig. 4. Simple regression plot of the average number of commits per user per project against the number of CVEs present in the project. The plot works in line with our intuition that the more responsibility individually taken up by a developer, the larger the potential to introduce vulnerabilities. While the plot was heavily cleaned up to remove outliers, it serves to motivate that meaningful inferences can be drawn from the GitHub data towards analyzing risk.

experiment we ran, where we analyzed the correlation between the average number of commits made per user per project against the number of CVEs in the project. We worked on this with the hypothesis that the more number the number of commits made per user on average, the more responsibility each user takes up individually, resulting in less collaboration, thereby leading to more vulnerabilities. The resulting plot is shown in Fig 4. While we had to do significant cleanup of the plot to remove outliers, the plot does seem to indicate the presence of a correlation. We use this plot as motivation to inform ourselves of the possibility of drawing meaningful inferences from data. Our further, more detailed experiment to quantifying risk is described below. While the results of that experiment do seem to work towards our goal towards quantifying risk, we used them to further refine our metric to present that mentioned in Section III.

A. Experiment

We first curated a list of approximately the top 1000 users, these were the users with the most contributions on GitHub, which we obtained from [19]. We then obtained a list of all projects these users had contributed to from GHTorrent and queried them against [20] to obtain their associated CVEs. Of the 57322 projects we obtained from the 1000 users, 618 projects had CVE scores associated with them. We then computed riskiness scores for each of the 1000 contributors with a simpler version of our riskiness metrics that did not incorporate the time aspect, and then computed the relative contribution as the ratio of the total number of commits made by the user to the total number of commits made to the project. The resulting distribution of user scores is shown in Figure ??.

B. Results

We saw a distinct correlation between higher risk scores from development teams and projects having a vulnerability at a later point in our data set. This was an approximate .63 correlation between increasing score and vulnerability. This correlation increases as the risk score drives higher as seen in figure 7. Our thought is that this correlation would be even

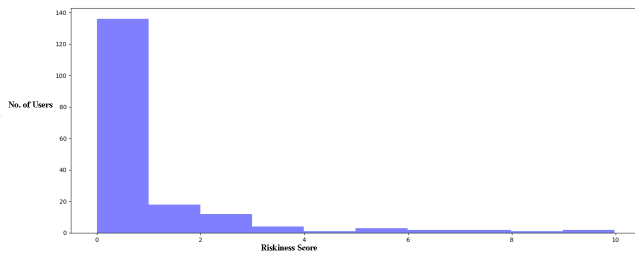


Fig. 5. Histogram showing the distribution of user scores. As can be seen from the figure, most user scores seem to cluster between 0 and 1. The plot only shows the counts for those users who had a non-zero score, and also does not show scores for a tiny portion users with abnormally large scores (over 300).

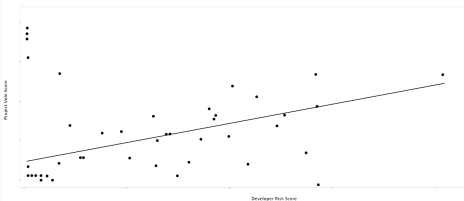


Fig. 6. Regression plot showing the user risk scores vs the vuln scores for a project over time.

higher with cleaner data. It was difficult to parse the disparate data sets and ensure the project was correlated to the NVD.

We anticipate greater correlation on future iterations of this work.

V. ACKNOWLEDGEMENTS

Jon Chapman would like to thank his advisor Dr. Shyhtsun Felix Wu, and Hari Venugopalan would like to thank his advisor, Dr. Samuel T. King, and both of us would like to thank Dr. Vladimir Filkov in helping to make our research possible, providing feedback, and inspiration for our work.

REFERENCES

- [1] E. Kalliamvakou, et al. "The promises and perils of mining GitHub." In Proceedings of the 11th Working Conference on Mining Software Repositories. Association for Computing Machinery, New York, NY, USA, 92–101. 2014. <https://doi.org/10.1145/2597073.2597074>
- [2] Chadni Islam, M. Ali Babar, Roland Croft, Helge Janicke, "SmartValidator: A framework for automatic identification and classification of cyber threat data". Journal of Network and Computer Applications, Volume 202, 103370, ISSN 1084-8045, 2022 <https://doi.org/10.1016/j.jnca.2022.103370>.
- [3] K. Watkins and S. M. Kywe, "Unsecured Firebase Databases: Exposing Sensitive Data via Thousands of Mobile Apps," Appthority, Tech. Rep., 2018.
- [4] Zuo, Chaoshun, Zhiqiang Lin, and Yinqian Zhang. "Why does your data leak? uncovering the data leakage in cloud from mobile apps." 2019 IEEE Symposium on Security and Privacy (SP). IEEE, 2019.
- [5] Alexandre Decan, Tom Mens, and Eleni Constantinou. "On the impact of security vulnerabilities in the npm package dependency network." In Proceedings of the 15th International Conference on Mining Software Repositories. Association for Computing Machinery, New York, NY, USA, 181–191. 2018. <https://doi.org/10.1145/3196398.3196401>
- [6] Common Vulnerabilities and Exposures <https://cve.mitre.org/>, retrieved Sept 2022
- [7] National Institutes of Standards and Technology National Vulnerability Database <https://nvd.nist.gov>, retrieved Sept 2022
- [8] Github Commitment to NPM ecosystem security <https://github.blog/2021-11-15-githubs-commitment-to-npm-ecosystem-security/>, retrieved Sept 2022
- [9] National Institutes of Standards and Technology Source Code Security Analyzers <https://www.nist.gov/itl/ssd/software-quality-group/source-code-security-analyzers>, retrieved Sept 2022
- [10] S. E. Ponta, H. Plate, A. Sabetta, M. Bezzi and C. Dangremont, "A Manually-Curated Dataset of Fixes to Vulnerabilities of Open-Source Software," 2019 IEEE/ACM 16th International Conference on Mining Software Repositories, pp. 383-387, 2019. doi: 10.1109/MSR.2019.00064.
- [11] R. A. Khan, S. U. Khan, H. U. Khan and M. Ilyas, "Systematic Literature Review on Security Risks and its Practices in Secure Software Development," in IEEE Access, vol. 10, pp. 5456-5481, 2022, doi: 10.1109/ACCESS.2022.3140181.
- [12] H. Al-Matouq, S. Mahmood, M. Alshayeb and M. Niazi, "A Maturity Model for Secure Software Design: A Multivocal Study," in IEEE Access, vol. 8, pp. 215758-215776, 2020, doi: 10.1109/ACCESS.2020.3040220.
- [13] Bana, Sarah, et. al. "Human Capital Acquisition in Response to Data Breaches" 2022. <http://dx.doi.org/10.2139/ssrn.3806060>
- [14] Zwilling, M., et al. "Cyber Security Awareness, Knowledge and Behavior: A Comparative Study." Journal of Computer Information Systems. 62. 82-97. 2022. 10.1080/08874417.2020.1712269.
- [15] N. Tomas, J. Li and H. Huang, "An Empirical Study on Culture, Automation, Measurement, and Sharing of DevSecOps," 2019 International Conference on Cyber Security and Protection of Digital Services (Cyber Security), pp. 1-8, 2019. doi: 10.1109/CyberSecPODS.2019.8884935.
- [16] Paul, CL, and J. Dykstra. "Understanding Operator Fatigue, Frustration, and Cognitive Workload in Tactical Cybersecurity Operations." Journal of Information Warfare, vol. 16, no. 2, pp. 1–11, 2017. JSTOR, <https://www.jstor.org/stable/26502752>.
- [17] Yamaguchi, Fabian & Golde, Nico & Arp, Daniel & Rieck, Konrad. "Modeling and Discovering Vulnerabilities with Code Property Graphs". 2014. Proceedings - IEEE Symposium on Security and Privacy. 10.1109/SP.2014.44.
- [18] Georgios Gousios. "The GHTorrent dataset and tool suite." In Proceedings of the 10th Working Conference on Mining Software Repositories (MSR '13). IEEE Press, 233–236. 2013. doi: 10.1109/MSR.2013.6624034.
- [19] Most active GitHub users <https://gist.github.com/paulmillr/2657075>, retrieved May 2022
- [20] CVE Details <https://www.cvedetails.com/>, retrieved June 2022
- [21] CIRCL CVE SEARCH <https://cve.circl.lu/api>, retrieved May 2022