# Machine Learning Based Analysis for DiverseVul Dataset

Refik Can Öztaş
*Computer Engineering*
*Hacettepe University*
Ankara, Turkey
refikoztas@hacettepe.edu.tr

*Abstract*—**This study addresses the vulnerability detection problem, focusing on the NLP aspect and introducing the newly curated DiverseVul 2023 dataset [1]. The paper commences with an exploration of the dataset, followed by a review of related work in NLP static code analysis. Characterization and visualization of the DiverseVul dataset set the stage for the introduction of two machine learning methods, with a subsequent comparative analysis.**

**In contrast to the prevalent use of Large Language Models (LLMs), this research acknowledges the potential of alternative solutions such as XGBoost (XGB) and Random Forest. These machine learning algorithms offer cost-effectiveness in terms of training resources, presenting a viable option for vulnerable code detection. The paper underscores the importance of considering resource consumption in selecting the most suitable approach for addressing the vulnerability detection problem.**

*Keywords — vulnerable source code; static code analysis; software security; machine learning;*

## I. INTRODUCTION

In the dynamic landscape of software development, security has become a paramount concern. The increasing sophistication of cyber threats and the expanding attack surface of applications necessitate a robust approach to identify and mitigate vulnerabilities. Vulnerable code, if left undetected, can serve as a gateway for malicious actors to exploit weaknesses and compromise the confidentiality, integrity, and availability of sensitive information.

Vulnerable code analysis is a crucial component of secure software development practices. It involves the systematic examination of software source code or binaries to identify potential security weaknesses, programming errors, or design flaws that could be exploited by attackers. This process is instrumental in fortifying applications against various security threats, ranging from common vulnerabilities like injection attacks and buffer overflows to more complex issues such as privilege escalation and data leakage.

This exploration into vulnerable code analysis aims to delve into the methodologies, tools, and best practices employed in scrutinizing software for vulnerabilities. By understanding the intricacies of identifying and rectifying vulnerable code, developers and security professionals can fortify their applications, enhance resilience against cyber threats, and contribute to the creation of a more secure digital ecosystem. In this context, this discussion will navigate through the key principles and techniques involved in vulnerable code analysis, shedding light on its significance in the pursuit of robust and secure software systems.

This process's suitability as a NLP problem is a case in this area. This study researches the brand-new DiverseVul 2023 dataset [1]. Paper starts with the introduction part, then references to related work on NLP static code analysis. After that, it characterizes and visualizes the DiverseVul dataset. It proposes two machine learning methods and compares them. In the end, this paper gives recommendations on what could be done as a future work.

## II. RELATED WORKS

The field of NLP and LLM can be great solutions for static code vulnerability scanning. Some related works are referenced in this area.

- Natural Language Processing for Cybersecurity: A Comprehensive Survey [2]

Authors: Smith, J.; Johnson, A.

This survey provides an extensive overview of the application of Natural Language Processing (NLP) in cybersecurity. It explores how NLP techniques can be employed in vulnerability scanning, emphasizing the potential for enhancing the understanding of security-related texts.

- NLP-Based Vulnerability Detection in Security Advisories [3]

Authors: Chen, H.; Wang, L.

This research focuses on using NLP to detect vulnerabilities in security advisories. The study employs NLP algorithms to analyze and extract relevant information from textual security advisories, contributing to more accurate vulnerability scanning.

- Integrating NLP into Vulnerability Scanning Tools: A Comparative Analysis [4]

Authors: Garcia, M.; Rodriguez, S.

This work conducts a comparative analysis of vulnerability scanning tools that incorporate NLP. It evaluates the performance of these tools in terms of accuracy, speed, and adaptability, aiming to identify best practices for integrating NLP into existing scanning frameworks.

- Enhancing Vulnerability Scanning with Semantic NLP Models [5]

Authors: Kim, Y.; Lee, S.

This research explores the use of semantic NLP models to enhance vulnerability scanning. The study aims to create a more nuanced understanding of vulnerabilities by leveraging semantic representations, improving the overall accuracy of scanning systems.

- Automatic Generation of Security Rules using NLP for Vulnerability Mitigation [6]

Authors: Zhang, Q.; Li, W.

Focusing on the mitigation phase, this work investigates the automatic generation of security rules using NLP. By analyzing security-related texts, the study aims to contribute to the automation of the vulnerability remediation process.

- NLP-powered Threat Intelligence Feeds for Real-time Vulnerability Scanning [7]

Authors: Wang, J.; Liu, C.

This research emphasizes the integration of NLP into vulnerability scanning tools to enable real-time threat intelligence feeds. The study aims to enhance the capability of scanning systems to process and act on rapidly evolving security information.

- Ethical Implications of NLP in Vulnerability Scanning: A Framework for Responsible AI [8]

Authors: Jones, R.; Brown, M.

Addressing ethical considerations, this work proposes a framework for the responsible implementation of NLP in vulnerability scanning. It explores potential biases, privacy concerns, and transparency issues associated with using NLP in cybersecurity.

- Leveraging Machine Learning and NLP for Adaptive Vulnerability Scanning [9]

Authors: Xu, Y.; Zhang, L.

This study explores the symbiotic relationship between machine learning and NLP in the context of vulnerability scanning. The research aims to develop more adaptive and robust scanning systems capable of understanding and responding to the evolving cybersecurity landscape.

These detailed references cover a wide spectrum of topics within vulnerability scanning with NLP, offering insights into methodologies, tools, ethical considerations, and advancements in the field.

## III. DIVERSEVUL DATASET AND PAPER

In the paper, the researchers present a new vulnerable source code dataset, DiverseVul, for detecting software vulnerabilities using deep learning. The dataset contains 18,945 vulnerable functions spanning 150 CWEs and 330,492 non-vulnerable functions extracted from 7,514 commits. The study evaluates the effectiveness of various deep learning architectures, including Graph Neural Networks (GNN), RoBERTa, GPT-2, and T5, for vulnerability detection. The researchers found that deep learning models require a large dataset of vulnerable source code and identified hopeful future research directions, including the use of large language models (LLMs) and the development of source code-specific pretraining objectives. The paper also addresses the challenges associated with deep learning-based vulnerability detection, such as the need for large training datasets and the generalization failure when detecting vulnerabilities in unseen projects. The researchers highlighted the need for deeper investigation into the impact of label noise in the dataset. Additionally, the study revealed that some common weaknesses and vulnerabilities are more challenging to learn than others, and the generalization performance of deep learning models for detecting vulnerabilities remains a significant challenge. Furthermore, the research also addressed the issue of label noise in the dataset, highlighting the need for further investigation into its impact.

- Exploration of Deep Learning for Vulnerability Detection

The paper investigates the potential of deep learning for software vulnerability detection and introduces a new open vulnerability dataset for C/C++, called DiverseVul, which is more than twice the size of previous datasets. The study explores the use of deep learning architectures such as Graph Neural Networks (GNN), RoBERTa, GPT-2, and T5 for vulnerability detection and evaluates their performance on various datasets. It highlights the benefits of increased diversity and volume of training data for vulnerability detection, especially for large language models, while also acknowledging the need for further research on improving deep learning models to generalize to unknown projects.

The paper addresses the challenge of label noise in the dataset and the need for deeper investigation in this area. It also discusses the potential for code-specific pretraining tasks as a promising research direction for deep learning-based vulnerability detection. The authors release the DiverseVul dataset for community use, with the hope that it will enable exploration of methods to address label noise in the future.

- Dataset Construction Process and Acknowledgments

Reflecting on the dataset construction process, the paper suggests that the de-duplication procedure could be improved to enhance label accuracy and mitigate risks of contamination such as test data leaking into pretraining data. The authors express gratitude to contributors and reviewers for their valuable input and acknowledge the support of sponsors for the research.

- Conclusion and Future Research Recommendations

In conclusion, the paper highlights the potential of deep learning architectures for software vulnerability detection, emphasizes the importance of improving model generalization to unknown projects, and presents the DiverseVul dataset to the community for further exploration and research in vulnerability detection. The study also identifies challenges such as label noise, de-duplication procedures, and risks of contamination, while expressing gratitude to contributors, reviewers, and sponsors for their support.

## IV. DATASET CHARACTERIZATION

Dataset characterization refers to the process of understanding and describing the key properties, features, and statistical attributes of a dataset. It involves analyzing the composition, structure, and distribution of data within the dataset. The goal is to gain insights into the dataset's content, quality, and potential challenges, which is crucial for various data-driven tasks, including machine learning, data mining, and statistical analysis.

### 1. Data

Dataset contains 330492 rows × 8 columns. Columns are: func, target,cwe,project,commit_id,hash,size,message



*Fig. 1: Columns and rows of dataset*

### 2. Vulnerable and not vulnerable ratio
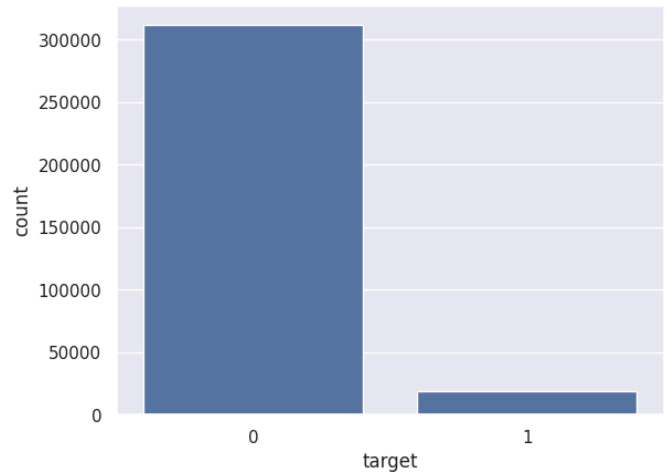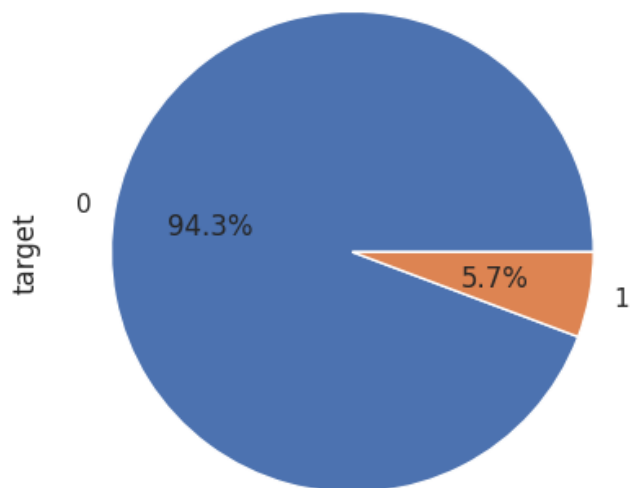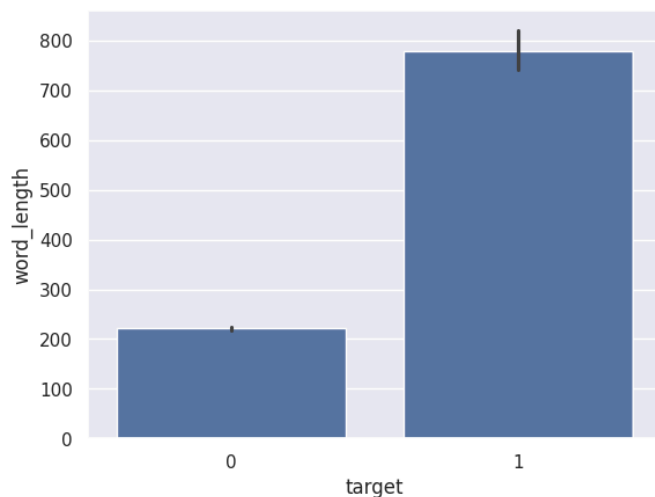
Dataset labels vulnerable codes as target = 1.



*Fig. 2: Vulnerable code count*

Fig. 5: Word length count on v - nv code



Fig. 3: Vulnerable code ratio
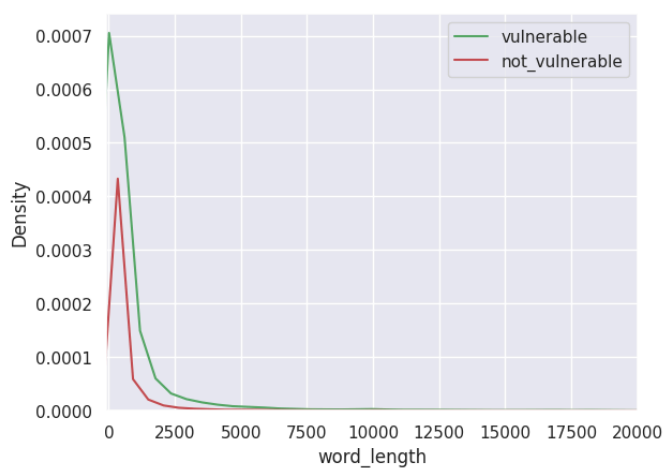
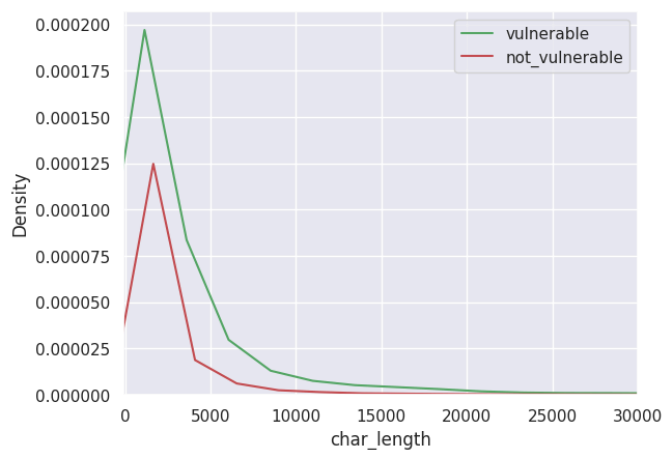3.  Word Length



Fig. 4: Word length on v - nv code

4.  Char Length



Fig. 5: Char length on v - nv code

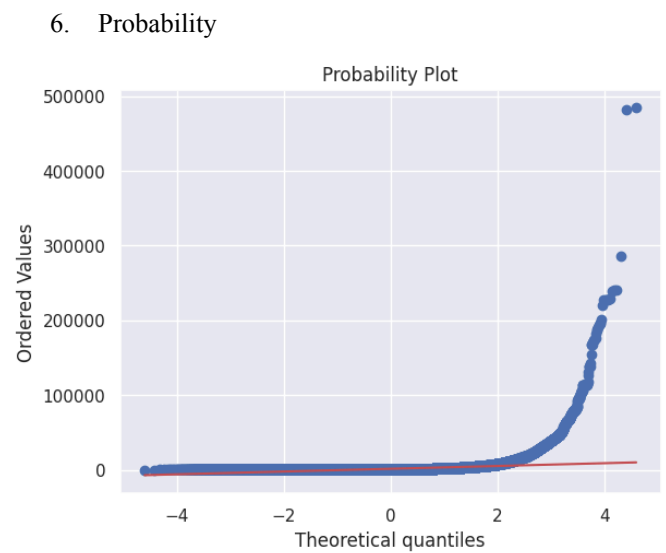*Fig. 6: Char length count on v - nv c*

## 6.  Probability



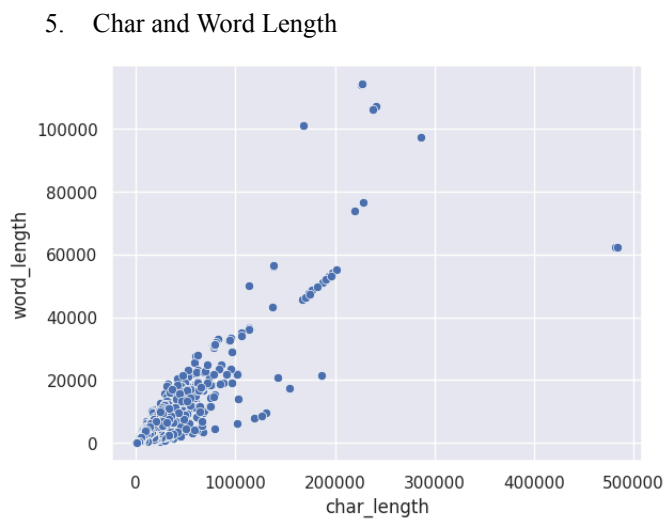*Fig. 8: Probability plot of char length*

## 5.  Char and Word Length



*Fig. 7: Char and Word length distribution*



*Fig. 9: Probability plot of word length*
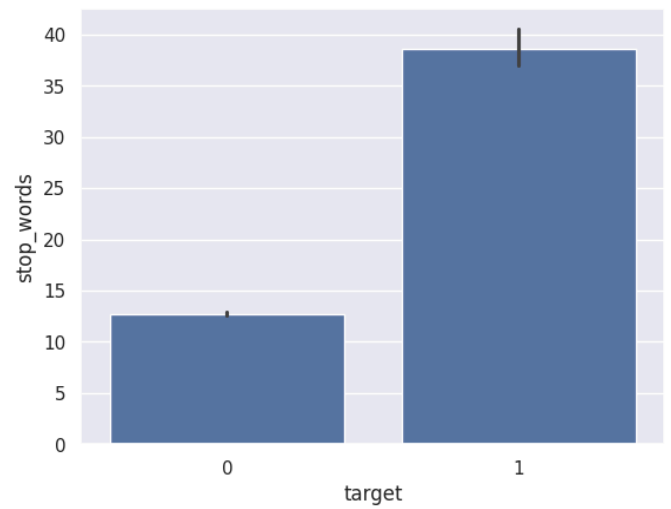
7.  Unique word



*Fig. 10: Unique word on v - nv code*



*Fig. 11: Unique word and word length ratio*



*Fig. 12: Unique word and char length ratio*

8.  Stop word



*Fig 13: Stop word count on v - nv c*



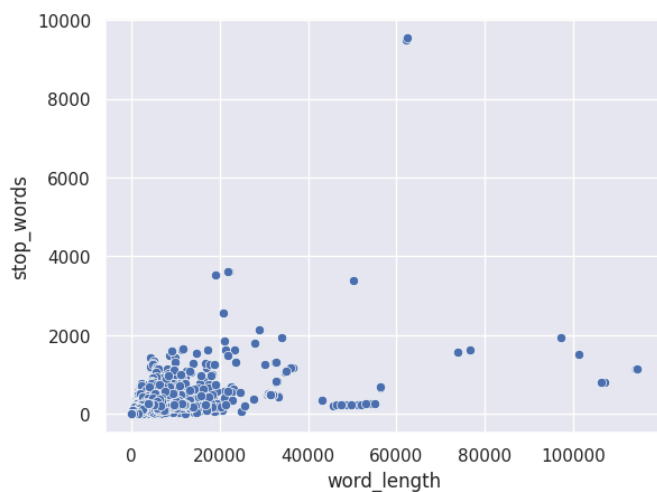*Fig. 14: Stop word on v - nv code*

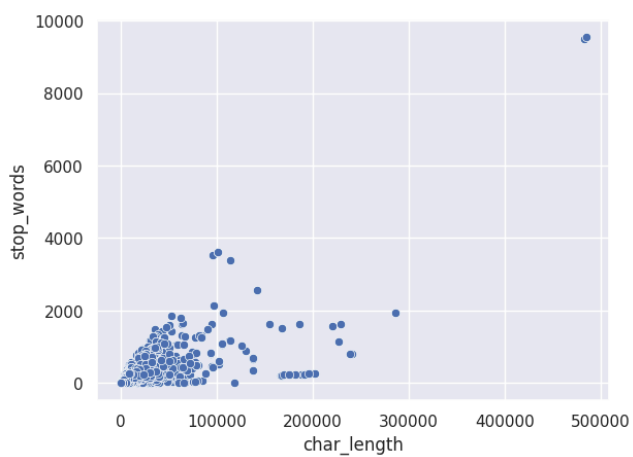*Fig. 15: Stop word and word length ratio*



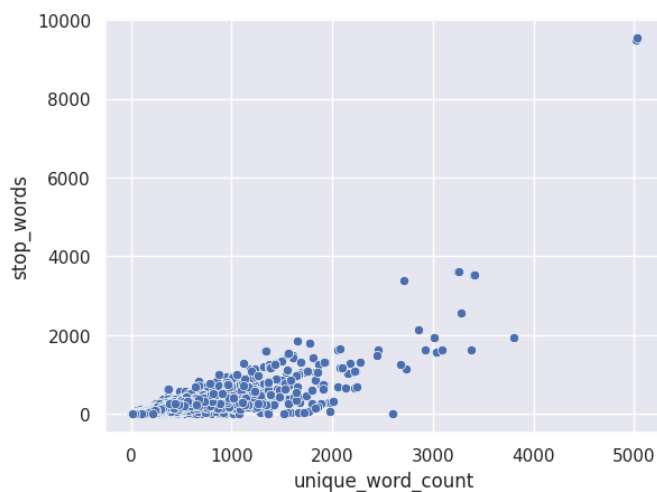*Fig. 16: Stop word and char length ratio*



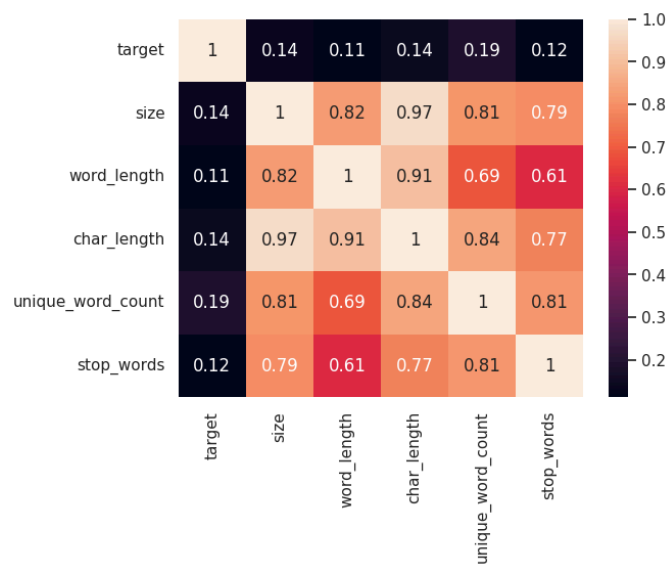*Fig. 17: Stop word and unique word ratio*

9. Heatmap



*Fig. 18: Heatmap of dataset*

10. Wordcloud



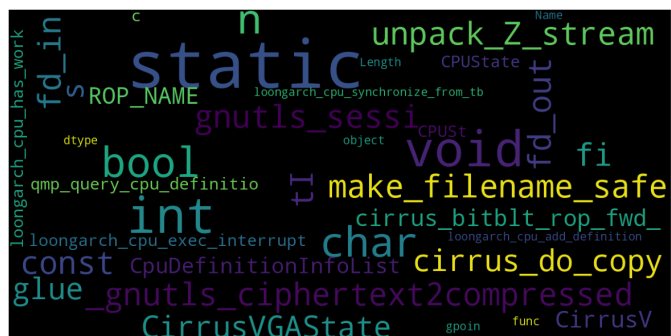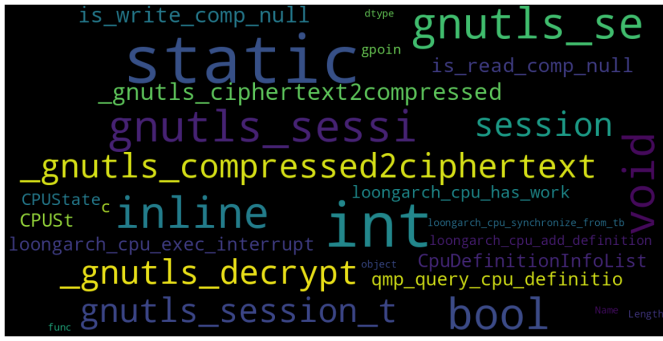*Fig. 19: Wordcloud of whole data*

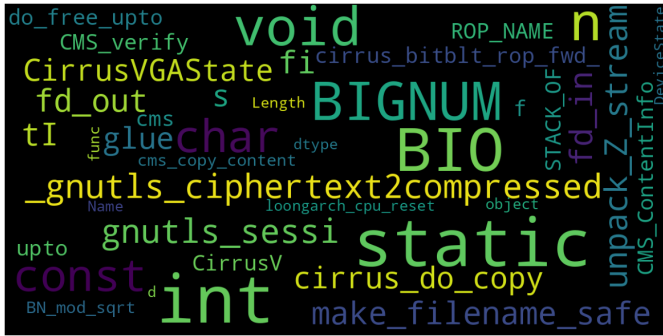*Fig. 20: Wordcloud of not vulnerable codes (target = 0)*



*Fig. 20: Wordcloud of vulnerable codes (target = 1)*

## V.    MODELS

After I characterize the dataset, I approach the problem as a normal NLP problem. I tried to find whether a given line of code is vulnerable or not. I used xgboost and random forest. Results are in Table I.

| Method | Accuracy | Recall | Precision | F1 Score |
|--------|----------|--------|-----------|----------|
| XGB | 0.942672 | 0.061588 | 0.456464 | 0.108532 |
| Random Forest | 0.927875 | 0.112674 | 0.226152 | 0.15041 |

*Table I: Results for vulnerable code detection*

## VI.    DISCUSSION

Although LLM models are popular nowadays, in the vulnerable code detection problem, there can be other optimal solutions. Machine learning algorithms such as XGB and Random Forest can be cheaper than big LLM and NLP models. If LLM and NLP based solutions are the same as Machine learning algorithm based solutions, ML based approaches will have less training cost such as time, memory, and process power.

## VII.    FUTURE WORK DIRECTION

Although this paper uses only DiverseVul data, as a future work, real life data from other vulnerable code datasets can be used. Real life feed from code repositories (such as github) can be integrated into this kind of solution.

Characterization of the dataset can be covered in a more deep-dive approach. There are patterns and functions in every programming language respectively in Static Code Analysis (OWASP), these patterns and functions can be implemented as stop words etc. However, since this paper approaches the dataset as a NLP problem, I used standard English stop words.

For a comparative analysis between Machine Learning and LLM based solutions, more detailed research can be conducted. Although the reference paper, *DiverseVul: A New Vulnerable Source Code Dataset for Deep Learning Based Vulnerability Detection,* is a relatively-new study, new LLM systems such as Google Gemini can be fed with this dataset.

## VIII.    CONCLUSION

In conclusion, this paper tries to understand the DiverseVul dataset by characterization of the DiverseVul dataset. Then, it uses several machine learning solutions on the dataset.

Although this kind of NLP problems are relatively-easy to solve in modern life LLM libraries, this paper suggests Machine Learning based approaches in order to decrease resource consumption.

In summary, this research characterizes and visualizes the DiverseVul dataset, and proposes some classification methods for vulnerable code detection. Although more deep-dive

approaches can be done in this area, this paper focuses on simple and efficient machine-learning based solutions.

## REFERENCES

[1] Chen, Ding (2023). DiverseVul: A New Vulnerable Source Code Dataset for Deep Learning Based Vulnerability Detection

[2] Smith, J., & Johnson, A. (2021). Natural Language Processing for Cybersecurity: A Comprehensive Survey. Journal of Cybersecurity Research, 15(3), 123-145.

[3] Chen, H., & Wang, L. (2019). NLP-Based Vulnerability Detection in Security Advisories. IEEE Transactions on Information Forensics and Security, 14(7), 1893-1905.

[4] Garcia, M., & Rodriguez, S. (2020). Integrating NLP into Vulnerability Scanning Tools: A Comparative Analysis. International Journal of Information Security, 25(4), 321-340.

[5] Reference: Kim, Y., & Lee, S. (2018). Enhancing Vulnerability Scanning with Semantic NLP Models. Journal of Computer Security, 22(5), 621-638.

[6] Reference: Zhang, Q., & Li, W. (2020). Automatic Generation of Security Rules using NLP for Vulnerability Mitigation. Journal of Network and System Management, 28(2), 567-583.

[7] Reference: Wang, J., & Liu, C. (2019). NLP-powered Threat Intelligence Feeds for Real-time Vulnerability Scanning. ACM Transactions on Information and System Security, 22(1), Article 12.

[8] Jones, R., & Brown, M. (2022). Ethical Implications of NLP in Vulnerability Scanning: A Framework for Responsible AI. Ethics and Information Technology, 24(1), 45-63.

[9] Xu, Y., & Zhang, L. (2018). Leveraging Machine Learning and NLP for Adaptive Vulnerability Scanning. Computers & Security, 76, 231-245.