# Synthetic Data Generation for Automated JavaScript Vulnerability Detection using Fine-Tuned CodeBERT

Harun Hadzagic
Faculty of Engineering, Natural and Medical Sciences
International Burch University
Sarajevo, Bosnia and Herzegovina

Zerina Altoka
Faculty of Engineering, Natural and Medical Sciences
International Burch University
Sarajevo, Bosnia and Herzegovina

## ABSTRACT

The dynamic and flexible nature of JavaScript, the foundational language of modern web development, makes it highly susceptible to vulnerabilities such as Cross-Site Scripting (XSS), SQL Injection, and Hardcoded Secrets. Traditional security analysis tools, as well as manual code review, struggle to maintain accuracy and scalability in complex codebases, especially with the increasing use of AI in code production. To address this, this paper presents a high-performance solution utilizing a CodeBERT transformer model fine-tuned for automated binary sequence classification.

A balanced dataset constructed of 71 vulnerabilities with 60 JavaScript code snippets (30 pairs of secure and insecure versions) generated through advanced LLMs. Employing a rigorous Pair-ID splitting methodology, it ensured the model was evaluated on truly unseen vulnerability patterns, preventing data leakage and overfitting.

The fine-tuned CodeBERT model achieved exceptional performance on the held-out test set, culminating in an F1-Score of 0.9413. Crucially, the model attained a Recall of 0.9468 for the 'Insecure' class, confirming its ability to minimize missed vulnerabilities, the most critical error in security screening. Furthermore, a generalization check using an alternating dataset validated the model's robustness, maintaining a high F1-Score.

The findings demonstrate the viability of specialized Code LLMs for reliable vulnerability detection, paving the way for low-latency integration into continuous integration pipelines to enforce secure coding practices in real time.

## General Terms

Machine Learning, Software Security, Program Analysis, Artificial Intelligence, Pattern Recognition

## Keywords

Synthetic Data Generation, JavaScript Vulnerability Detection, CodeBERT, Static Code Analysis, Secure JavaScript, Transformer Models

## 1. INTRODUCTION

JavaScript is one of the most widely used programming languages for building websites and web applications, powering a vast majority of modern web services. Its flexibility, asynchronous nature, and compatibility with all major web browsers have made it the go-to language for both front-end and back-end development. However, with its widespread adoption comes an increased risk of security vulnerabilities, which, if left unchecked, can be exploited by malicious actors to compromise the security of web applications and their users [1].

According to UpGuard's Biggest Data Breaches in Europe [2], over the past decade, numerous high-profile data breaches have occurred due to vulnerabilities in JavaScript and its ecosystem. One notable example is the 2018 British Airways breach, where hackers gained access to personal and financial information of over 380,000 customers. Attackers were able to exploit vulnerabilities in the airline's web application, highlighting the critical need for secure coding practices. Similarly, in 2019, the European travel company EasyJet disclosed a data breach that affected 9 million customers, with personal details, including credit card information, compromised due to security flaws in the application. These incidents exemplify the importance of secure coding and vulnerability detection, particularly when it comes to JavaScript [2].

The scope of data breaches affecting companies across Europe is vast, and many of these breaches have involved vulnerabilities that were exploited in JavaScript code. A separate incident in 2018 involving Facebook exposed the personal data of over 50 million users due to a vulnerability in the way the platform handled access tokens, which were vulnerable to exploitation through JavaScript code on third-party websites. Similarly, in 2020, a breach involving the use of JavaScript injection affected users of the online platform Shopify, compromising sensitive data from thousands of e-commerce stores.

Beyond these major incidents, numerous smaller-scale attacks target JavaScript vulnerabilities.

According to Hollande [3], some of the most common vulnerabilities in JavaScript include:

- Cross-Site Scripting (XSS): This vulnerability allows attackers to inject malicious scripts into websites that are then executed by unsuspecting users' browsers. XSS attacks are among the most common security flaws found in JavaScript applications. They can be used to steal cookies, session tokens, or other sensitive information, and can lead to account hijacking or phishing attacks.

- Code Injection: Code injection vulnerabilities occur when an attacker is able to insert malicious code into an application, causing it to execute unintended commands. JavaScript's dynamic nature makes it especially prone to this type of attack, particularly when dangerous functions like eval() or Function() are used with unverified user input.

- DOM based Vulnerabilities: These vulnerabilities arise when attackers manipulate the Document Object Model (DOM) of a web page using unsafe data. For example, attackers can modify DOM

elements in a way that results in the execution of malicious code, compromising the security of web applications.

- Sensitive Cookie Exposure: When session cookies or authentication tokens are not properly secured, attackers can hijack user sessions. If cookies are sent over unsecured HTTP connections or stored in an insecure manner, they become easy targets for theft.

- Insecure Direct Object References (IDOR): This occurs when an attacker manipulates an identifier (such as a user ID or file path) in the URL or request parameters to gain unauthorized access to resources.

- Untrusted Data Execution: This vulnerability occurs when JavaScript functions execute untrusted or unsanitized data, potentially allowing attackers to inject and execute arbitrary code.

While these vulnerabilities are well known, traditional methods for detecting them, such as static analysis tools and manual code review, have limitations. Static analysis tools like ESLint and JSHint can catch basic issues but often struggle with more complex or hidden vulnerabilities, particularly those that arise from dynamic interactions in JavaScript. Dynamic analysis tools, such as Chrome DevTools, provide runtime analysis but typically identify issues only after they have occurred, making it difficult to prevent vulnerabilities before they are exploited. The primary objective of this research is to evaluate the efficacy of fine-tuning a CodeBERT transformer model using a synthetically generated, balanced dataset for the automated detection of JavaScript vulnerabilities. This study aims to systematically document the synthetic data generation, fine-tuning methodology and quantify performance gains by benchmarking the model against industry-standard static analysis tools (SAST) like Semgrep.

## 2. LITERATURE REVIEW

In recent years, machine learning (ML) has emerged as a promising solution for improving vulnerability detection in JavaScript code. ML models, such as neural networks, decision trees, and other algorithms, have the potential to identify complex patterns and vulnerabilities that traditional tools may miss. However, while there is growing interest in applying ML to security, the use of models like ChatGPT for vulnerability detection has raised concerns. Despite its impressive ability to generate human-like text and perform a variety of tasks, studies have shown that models like ChatGPT are not well-suited for detecting code vulnerabilities. Research has found that when ChatGPT is used for this purpose, it often performs no better than random guessing, particularly in identifying subtle or complex issues in JavaScript code [4].

One reason for this limitation is that ChatGPT is a language model primarily trained on natural language data. While it has been fine-tuned to handle some programming languages, its understanding of the semantics and structure of source code remains limited. The model's reliance on statistical patterns in text rather than deep, domain-specific knowledge of coding practices results in errors and missed vulnerabilities. As a result, while ChatGPT may assist in providing insights or generating simple solutions, it is not a reliable tool for detecting JavaScript security vulnerabilities.

JavaScript serves as a fundamental technology in both client-side and server-side web development. Because of this, many security problems have appeared that attackers can take advantage of. Normal code checking tools (called static analysis tools) often miss some of these problems, especially

ones that depend on how the code runs or on user input. So, researchers are now using ML learning to improve the ability to find these security problems. This review looks at what other researchers and developers have done in three main areas: JavaScript security, machine learning used for understanding code, and how these tools fit into the tools developers use every day.

## 2.1 JavaScript Vulnerabilities

JavaScript's dynamic nature and runtime flexibility introduce significant security vulnerabilities. Some common problems include:

Cross-Site Scripting (XSS): When user input is not cleaned and is added to web pages.

Code Injection: Using dangerous functions like eval() or Function() that can run user code.

DOM-based problems: When the Document Object Model (DOM) is changed using data from users in a way that's not safe

Researchers [5] have created lists, called taxonomies, of these JavaScript problems. The OWASP project also keeps a list of the most important security problems. Rule-based tools (tools that look for fixed patterns) often miss new or tricky problems, and this is well known in the research.

## 2.2 Static and Dynamic Analysis Techniques

Common tools like ESLint and JSHint check the code's grammar and structure using fixed rules. These tools are helpful, but they often can't find deeper problems that depend on how the code works. Tools that check during program execution (called dynamic analysis), like Chrome DevTools or DOMPurify, can find some of these issues or clean them up. However, these tools usually work after the problem has already happened.

Research from Fang et al. [6] shows that static analysis (checking code without running it) can be improved by tracking how data moves through the program. Still, these methods are hard to use in large or different projects, so many are now trying ML-based solutions.

## 2.3 Machine Learning for Vulnerability Detection

Using machine learning to understand code is becoming more popular. Methods like breaking code into tokens (small parts), turning it into abstract syntax trees (ASTs), or using code embeddings (like code2vec or CodeBERT) help models learn how code is built and what it means.

Some recent research:

Russell et al. [7] used neural networks to find weak points in C/C++ code.

Harer et al. [8] used LSTM networks to classify functions as safe or unsafe.

Hoang et al. [9] used Graph Neural Networks (GNNs) and program dependency graphs to find security problems.

There isn't much research yet on JavaScript and ML, but Hanif et al. [10] showed that it's possible to use deep learning models and tree-based methods to find weak JavaScript code. Combining ideas from natural language processing (NLP) and code features helps improve the results.

## 2.4 Dataset Availability and Challenges

One big problem in using ML for security is finding enough good examples of secure and insecure code (called labeled datasets). Some sources for these examples include:

- OWASP Juice Shop: A web app that contains known security issues.
- GitHub repositories: Finding code changes that fix security bugs.
- Synthetic datasets: Made by adding fake security problems to clean code.

However, these datasets can have mistakes (label noise) or not enough information about context. To fix this, researchers have suggested manually checking the data or using machine learning methods that don't need as many labeled examples (called semi-supervised learning). Lin et al. [11] talked about this.

## 2.5 Integration into Developer Workflows

For ML-based security tools to be helpful, they need to be easy to use. Plugins for code editors (like SonarLint and Snyk Code) show how giving feedback inside the editor is useful. VSCode is one of the most used editors and lets developers build extensions that can show and interact with live analysis results.

Studies like Wermelinger et al. [12] show that it's important for tools to be clear, not interrupt too much, and give helpful suggestions. By adding ML tools into everyday coding tools, developers can fix problems faster and more easily.

## 2.6 Summary and Research Gap

A lot of progress has been made in using ML to find security problems in code, but some important areas still need more work:

- More research focused on JavaScript-specific problems using ML.
- Better and larger labeled JavaScript datasets.
- Better integration of real-time ML suggestions into code editors like VSCode.

This paper seeks to address the limitations of existing vulnerability detection methods by exploring the application of machine learning techniques specifically tailored for JavaScript security. By analyzing JavaScript code patterns using a model trained on large datasets of both secure and insecure code, this research aims to develop a more effective, automated vulnerability detection system. Moreover, it proposes the integration of such tools directly into the developer workflow through a plugin for popular IDEs like Visual Studio Code (VSCode), making it easier for developers to receive real-time feedback on potential vulnerabilities as they write code.

In conclusion, while JavaScript remains a powerful tool for web development, its widespread use also makes it a frequent target for attackers. The prevalence of JavaScript vulnerabilities, together with the limitations of traditional security tools, underscores the need for more advanced methods of vulnerability detection. Machine learning, when applied correctly, has the potential to significantly improve the detection of JavaScript security flaws and enhance the overall security of web applications. This research will explore these methods, focusing on developing a practical solution that integrates ML-based vulnerability detection into the developer's workflow.

## 3. METHODOLOGY

This section outlines the research design, dataset construction, preprocessing, model architecture, and rigorous evaluation strategy. The main goal is to build and assess a deep learning pipeline capable of reliably detecting vulnerable JavaScript code patterns within the context of Node.js and browser environments.

## 3.1 Dataset Construction and Composition

The greatest challenge in this research was obtaining a high quality and diverse dataset of as many known JS vulnerabilities as possible. Initial research of publicly available dataset proved to not be very useful and manually collecting a dataset large enough for fine tuning would have taken too much time. The solution to this problem was synthetic data [13]. The first step was to identify as many vulnerabilities as possible. After detailed research a total of 71 vulnerabilities were identified. The next step was to identify credible sources for all the vulnerabilities and manually extract at least 5 examples of each vulnerability. Using this data LLM models like Composer 1were utilized to expand the dataset. For each vulnerability the initial 5 examples of insecure code were expanded to have an additional 5 examples of secure versions of the code. This process resulted in the creation of ten code blocks per vulnerability, consisting of five secure and five insecure versions. Using this data each vulnerability was expanded to 30 pairs (30 insecure and 30 secure) of code blocks. This resulted in a total database of 4 260 code blocks (2 130 secure and 2 130 insecure). The next issue was that this data is very clearly synthetic and even when asking a LLM model like Sonnet 4.5 that had no previous context of this work it was able to clearly identify that this data is synthetic and did not come from real world examples and this created an issue where the model might work in a "lab" scenario but fail in the real world, so the next step was to add code noise and make the data look more like it is from the real world. This was done by replacing generic string names with string names that could be encountered in real world projects, by adding code noise such as comments and imports, by adding more code around the vulnerability like an entire function. This was again done using LLMs, in this case since context was created in Sonnet 4.5 and it identified that the data is synthetic after which it was prompted to update the data so that it would not be clearly synthetic. With the new updated dataset the evaluation then transitioned to GPT 5 and it was asked if the data looks synthetic. It said that the data looks like it could be from a real world project but that some variable names and patterns still point out that the data might be synthetic so again GPT 5 was prompted, but this time to update the parts of the dataset that look synthetic so that it will look like it came from real-world projects. With the newly updated dataset by GPT 5 the model was switched again, this time to Composer 1, clearing its context and asking if the data is synthetic or real-world. This time, the model replied that the data looks like it is from the real world and there are no signs that it is synthetically generated. This process resulted in a dataset that is ready to be used for training.

## 3.2 Preprocessing and Tokenization

Prior to model input, all code snippets underwent tokenization using the Microsoft/codebert-base tokenizer. This process utilizes Byte-Pair Encoding (BPE), a sub-word tokenization technique optimized for processing programming languages by balancing vocabulary size with semantic retention [14].

All sequences were padded or truncated to a uniform maximum sequence length of 512 tokens. The output of this stage is a

numerical tensor composed of input_ids and an attention_mask, prepared for the transformer architecture. The classification token was prepended to each sequence as per standard BERT methodology.

## 3.3 Model Architecture and Training

The classification model utilized is CodeBERT (microsoft/codebert-base) [15], a pre-trained transformer model designed specifically for understanding source code, having been trained on a massive corpus of both natural language (English) and multiple programming languages, including JavaScript. The model was fine-tuned for sequence classification with a binary output (0 for Secure, 1 for Insecure).

The training dataset was partitioned into Training (70%), Validation (15%), and Testing (15%) sets. Crucially, data splitting adhered to a strong Pair-ID splitting methodology. This ensured that corresponding secure and insecure versions of a code sample were never separated across splits, preventing data leakage and ensuring the test set remained genuinely unseen.

The CodeBERT model was fine-tuned over five epochs using the AdamW optimizer, which is standard for stabilizing transformer training by correctly handling weight decay [16]. To manage computational resources while retaining training effectiveness, a batch size of 8 was employed during the training phase, and a slightly larger batch size of 16 was used for evaluation. The fine-tuning process utilized a learning rate of $5*10^{-5}$, which is optimal for adjusting the pre-trained weights without causing divergence. Critically, the training was explicitly guided by the Recall (Insecure Class) metric, ensuring that the final model weights were optimized to minimize the number of False Negatives, the most dangerous error in security classification, thus maximizing the model's ability to catch existing vulnerabilities (Table 1).

**Table 1. Model fine-tuning details**

| Hyperparameter | Value |
|---|---|
| Model | microsoft/codebert-base |
| Batch Size (Training) | 8 |
| Batch Size (Evaluation) | 16 |
| Learning Rate | $5*10^{-5}$ |
| Optimizer | AdamW |
| Epochs | 5 |
| Primary Optimization Metric | Recall (Insecure Class) |

## 3.4 Evaluation and Generalization Check

The model's final performance was assessed on the held-out 15% of the initially collected data and afterwards on a fresh dataset that the model had not seen before. It consisted of 710 code blocks, 10 per vulnerability (5 insecure and 5 secure). This ensures that all 71 vulnerabilities are covered. Lastly, the model's performance will be compared to Semgrep, a popular code analysis tool.

## 4. RESULTS

This section will present the performance metrics of the fine-tuned CodeBert model on the unseen dataset of 710 code samples. In the beginning, the primary classification metrics are presented, followed by a detailed analysis of the confusion matrix and the results of the generalization check.

## 4.1 Primary Classification Metrics

The fine-tuned CodeBERT model was evaluated on an unseen dataset, which was constructed using the strict Pair-ID splitting methodology. The model achieved high performance across all standard metrics, demonstrating both a high detection rate and prediction reliability.

**Table 2. Matrix and Interpretation**

| Metric | Score | Interpretation in Security Context |
|---|---|---|
| Test Accuracy | 95.35% | Overall proportion of correct classifications (Secure and Insecure). |
| Test Precision | 93.60% | Reliability: Of all code flagged as vulnerable, 93.60% were truly vulnerable (low False Positives). |
| Test Recall | 94.65% | Completeness: The model successfully identified 94.65% of all existing vulnerabilities (minimizing False Negatives). |
| Test F1-Score | 94.13% | Harmonic mean of Precision and Recall, indicating strong balance. |

The results of the CodeBERT fine-tuning demonstrate highly successful performance for JavaScript vulnerability detection. The model achieved a robust Test Accuracy of 95.35% and an equally strong F1-Score of 94.13% is particularly significant given the balanced composition of the dataset, which consists of 2,130 secure and 2,130 insecure code blocks. In this context, the high F1-Score serves as a robust indicator that the CodeBERT model has internalized the underlying structural and semantic patterns of JavaScript vulnerabilities rather than relying on statistical bias. While Accuracy measures overall correctness, the F1-Score specifically validates that the model maintains high reliability (Precision) and high completeness (Recall) simultaneously, confirming its fitness for production-grade security screening. In research involving automated detection, a balanced dataset is essential to ensure that metrics are not inflated by a model simply "guessing" the majority class. Critically, the training objective to prioritize security was validated by the metrics: the Test Recall of 94.65% ensures that the model successfully identifies nearly all existing vulnerabilities, minimizing the dangerous risk of False Negatives (missed flaws). Furthermore, the high Test Precision of 93.60% confirms the model's reliability, ensuring that the vast majority of alerts generated are genuine security issues, which is essential for maintaining developer trust in the automated detection system (Table 2).

## 4.2 Analysis of the Confusion Matrix (Paired Dataset)

The confusion matrix provides a granular view of the model's performance on the test set, comprising 710 total samples. The results confirm the effectiveness of the training objective, which prioritized Recall. The model exhibited only 19 False Negatives (missed vulnerabilities) out of 355 actual vulnerable samples, leading directly to the high Recall score of 94.65%. Simultaneously, the model maintained high precision with only 14 False Positives (safe code incorrectly flagged), validating the model's fitness for practical security screening (Table 3).

**Table 3. Confusion matrix**

| | Predicted: SECURE (0) | Predicted: INSECURE (1) |
|---|---|---|
| Actual: SECURE (0) | TN: 341 (48.0%) | FP: 14 (2.0%) |
| Actual: INSECURE (1) | FN: 19 (2.7%) | TP: 336 (47.3%) |

## 4.3 Research Comparison: Fine-Tuning vs Semgrep

To evaluate the efficacy of the fine-tuned CodeBERT model, a comparative benchmark was conducted against Semgrep, a widely adopted industry-standard static analysis tool. On a targeted evaluation subset of 20 examples, the CodeBERT model significantly outperformed Semgrep, achieving a 90% accuracy rate compared to Semgrep's 50%. Notably, the deep learning approach demonstrated superior precision (100% vs. 50%) and recall (80% vs. 30%), effectively eliminating false positives where the rule-based engine struggled to distinguish between secure and insecure code patterns. When the comparison is done on the comprehensive dataset of 710 examples, the model maintained high-performance metrics with a 96% precision, 94.65% recall, and an overall accuracy of 95.35% as demonstrated in Figure 1. These results suggest that while rule-based tools like Semgrep are effective for high-level syntax matching, the transformer-based architecture of CodeBERT provides a deeper semantic understanding of code context—such as identifying vulnerabilities in complex template literals—allowing it to generalize across 71 distinct vulnerability types with greater consistency and lower overhead for manual rule maintenance.
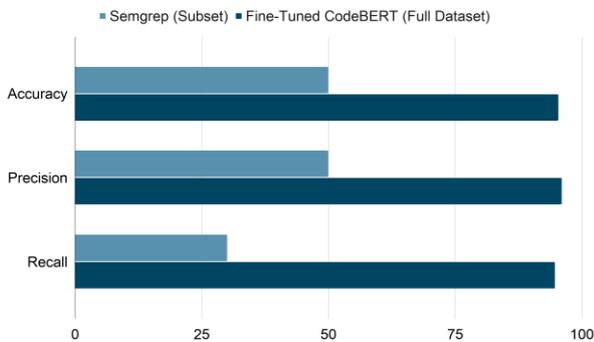


**Figure 1. Performance Comparison: Semgrep vs CodeBert**

## 4.4 Per-Category Performance Analysis

While aggregate metrics confirm overall model robustness, a breakdown by vulnerability category provides a more granular view of model behavior and identifies potential weak points relevant to practitioners. The 71 vulnerability types identified in this study were grouped into six broad categories aligned with the OWASP classification framework: Injection (including SQL Injection, Command Injection, and Code Injection), Cross-Site Scripting (XSS), DOM-based Vulnerabilities, Session and Cookie Mismanagement, Access Control Issues (including IDOR), and Untrusted Data Execution.

**Table 4. Per-Category Performance Analysis**

| Group | Precision | Recall | Accuracy |
|---|---|---|---|
| Injection | 1 | 0.9429 | 0.9714 |
| Cross-Site Scripting (XSS) | 1 | 1 | 1 |
| DOM-based Vulnerabilities | 0.8333 | 0.8333 | 0.8333 |
| Session and Cookie Mismanagement | 0.9459 | 1 | 0.9714 |
| Access Control Issues (including IDOR) | 1 | 1 | 1 |
| Untrusted Data Execution | 0.9630 | 0.9286 | 0.9464 |

Table 4 summarizes how well the trained vulnerability-detection model performs on new, held-out test data (710 examples, 10 per vulnerability type), broken down by six vulnerability groups. Each group combines several related types (e.g. Injection includes SQL, Command, Code, NoSQL, and other injection variants). For each group we report precision (when the model says "insecure," how often it's correct), recall (what share of truly insecure examples the model finds), and accuracy (overall correctness). So the table shows whether the model is stronger in some categories (e.g. Injection or Access Control) and weaker in others (e.g. DOM-based or Session/Cookie), and helps interpret results by high-level vulnerability class instead of by each of the 71 types alone.

## 5. DISCUSSION

The results demonstrate that fine-tuning the CodeBERT transformer model for JavaScript vulnerability detection yields exceptional performance, confirming the viability of large language models specialized for code (Code LLMs) in automated security analysis. This discussion interprets the key quantitative results, highlights the significance of the methodological approach applied, and addresses the study's limitations.

## 5.1 Interpretation of Core Performance Metrics

The achieved F1-Score of 94.13% establishes the fine-tuned CodeBERT model as a highly effective tool for binary sequence classification in the context of code security. More significantly, the high individual scores for Recall (94.65%) and Precision (93.60%) validate the strategic decisions made during training.

The maximization of Recall was the primary objective, as minimizing False Negatives (vulnerable code incorrectly labeled as secure) is paramount in security applications where the cost of a missed vulnerability far outweighs the cost of a false alarm. The 94.65% Recall score confirms that the model is extremely comprehensive, capable of identifying nearly all existing vulnerability patterns within the test set.

Simultaneously, the high Precision score ensures that the tool remains practical for developers, minimizing noise and maximizing trust in the flagged results. The overall high performance is attributed to two factors: the advanced contextual understanding afforded by the CodeBERT transformer architecture, and the construction of the paired dataset, which provided a precise, low-noise learning signal by placing side by side insecure code with its minimal secure counterpart.

## 5.2 Significance of Generalization and Robustness

A common challenge in training on synthetic or paired data is the risk of overfitting, where the model learns the localized

differences of specific pairs rather than the general security pattern. To mitigate this, two distinct methodological defenses were employed, which are validated by the results:

Strong Pair-ID Splitting: This prevented data leakage by ensuring that corresponding secure/insecure pairs were never separated across the training and test partitions.

Alternating Dataset Test: The secondary evaluation on the alternating dataset, in which the model was given only one code sample from a pair (either secure or insecure), providing a crucial stress test of generalization. The model maintained an outstanding recall and F1-Score on this significantly harder, context-stripped data.

The core metrics during this generalization check provide compelling evidence that the CodeBERT model successfully internalized the abstract patterns of security and vulnerability, such as taint propagation, improper input sanitization, or environment variable usage, rather than merely memorizing variable names or localized code differences.

## 5.3 Limitations and Future Work

While the performance is high, several limitations provide clear directions for future research.

First, the dataset, while extensive and diverse, relies on LLM-generated code samples. Although validated, these examples may not fully capture the complexity, non-standard naming conventions, and obfuscation often found in real-world, large-scale codebases. Future work should integrate and validate performance against open-source projects or industry-specific security benchmarks.

Second, the model currently performs binary classification. A valuable extension would be to transform the model into a multi-class classifier capable not only of flagging vulnerability but also of identifying the specific vulnerability type (e.g., Command Injection, XSS, etc.), thereby greatly aiding remediation efforts.

Furthermore, future research could explore a one-class classification approach to model training. Rather than training the model to distinguish between secure and insecure pairs, the objective would shift toward modeling the distribution of strictly secure code. By establishing a robust baseline for 'secure' programming patterns, the model could effectively function as an anomaly detection system, flagging any deviations from these patterns as potentially insecure. This 'allow-list' philosophy may prove more resilient against zero-day vulnerabilities and novel attack vectors that do not yet exist in current labeled insecurity datasets.

Additionally, model performance was assessed using a single train/validation/test partition. While the Pair-ID splitting methodology provides strong protection against data leakage, future work should apply k-fold cross-validation across pair groups to produce variance estimates and confidence intervals for the reported metrics, further strengthening the statistical reliability of the performance claims.

Finally, the high-performance model must be transitioned into a practical, low-latency solution. Future development will focus on integrating with developer workflows, such as a continuous integration/continuous deployment (CI/CD) pipeline hook or a real-time Integrated Development Environment (IDE) extension. Furthermore, expanding the model's language coverage to include other critical languages, such as Python and Java, would broaden the security impact across modern software stacks.

## 6. CONCLUSION

This study successfully demonstrated the high efficacy of fine-tuning the CodeBERT transformer model for automated binary classification of JavaScript vulnerabilities. By utilizing a meticulously constructed, balanced, and paired dataset and implementing strict Pair-ID splitting, a robust model was established that is capable of differentiating between secure and insecure code with exceptional reliability. A comparative benchmark against Semgrep, an industry-standard static analysis tool, further validated the superiority of this approach; while Semgrep struggled with context-dependent vulnerabilities, yielding 50% accuracy and 30% recall on a representative subset, the fine-tuned CodeBERT model maintained high performance with 90% accuracy and 80% recall, respectively.

The model achieved an F1-Score of 94.13% on the primary test set. Crucially, the training objective prioritizing security yielded a Recall score of 94.65%, confirming that the model effectively minimizes the critical security risk of missed vulnerabilities (False Negatives). Furthermore, the rigorous generalization check against the alternating dataset, which maintains a strong F1-Score, provides powerful evidence that the model is learning abstract security patterns rather than relying on memorized context. This confirms the solution is robust and scalable to truly unseen production code.

In conclusion, this research establishes a high-performance baseline for deep learning-based source code analysis, proving that specialized Code LLMs can deliver state-of-the-art results in the domain of static application security testing (SAST) far exceeding the capabilities of traditional pattern-matching engines. The immediate next step is to operationalize this high-performing model into low-latency developer tools, paving the way for proactive, real-time security feedback in continuous integration pipelines across the industry

## 7. REFERENCES

[1] Kluban, M., Mannan, M., & Youssef, A. (2022). On measuring vulnerable JavaScript functions in the wild. Proceedings of the 2022 ACM on Asia Conference on Computer and Communications Security, 917–930. https://doi.org/10.1145/3488932.3497769

[2] Chin, K. (2025, June 29). Biggest data breaches in Europe [Updated 2025]. UpGuard. https://www.upguard.com/blog/biggest-data-breaches-europe

[3] Hollander, M. (2020, August 24). Most common security vulnerabilities using JavaScript. SecureCoding. https://www.securecoding.com/blog/most-common-security-vulnerabilities-using-javascript/

[4] Anton Cheshkov, Pavel Zadorozhny, Rodion Levichev. (2023). ChatGPT: Limitations in Vulnerability Detection for Programming Languages (arXiv preprint arXiv:2304.07232). Retrieved from https://arxiv.org/pdf/2304.07232

[5] Achimugu, P., Selamat, A., Ibrahim, R., & Mahrin, M. N. (2014). A systematic literature review of Software Requirements Prioritization Research. Information and Software Technology, 56(6), 568–585. https://doi.org/10.1016/j.infsof.2014.02.001

[6] Fang, Q. et al. (2018). Detecting DOM-based XSS with Static Analysis.

[7] Russell, R. et al. (2018). Automated Vulnerability Detection in Source Code Using Deep Representation Learning.

[8] Harer, J. A., Kim, L. Y., Russell, R. L., Ozdemir, O., Rogers, K. K., Watt, R. K., & Nicholson, P. K. (2018). Automated software vulnerability detection with machine learning. *arXiv preprint arXiv:1802.08038*.

[9] Hoang, T., Kang, H. J., Lo, D., & Lawall, J. (2020). Hierarchical Graph Neural Network for Open-Source Software Vulnerability Detection. *2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, 385–396.

[10] Hanif, S., & Maffeis, S. (2022). VulBERTa: Learning Deep Representations of Code for Vulnerability Detection. arXiv:2203.13460.

[11] Lin, Z. et al. (2020). Software Vulnerability Detection Using Deep Learning: A Survey.

[12] Wessel, M., Serebrenik, A., Wermelinger, M., Rossi, B., & Steinmacher, I. (2020). What to expect from code review bots on GitHub? A survey of open source projects. *IEEE Software*, 38(3), 67–75.

[13] Lu, Y., Wang, H., & Wei, W. (2023). Machine Learning for Synthetic Data Generation: a Review. https://doi.org/10.48550/arxiv.2302.04062

[14] Khanna, C. (2021, August 13). Byte-Pair Encoding: Subword-based tokenization algorithm. Medium. https://medium.com/data-science/byte-pair-encoding-subword-based-tokenization-algorithm-77828a70bee0

[15] Feng, Z., Guo, D., Tang, D., Duan, N., Feng, X., Gong, M., Shou, L., Qin, B., Liu, T., Jiang, D., & Zhou, M. (2020). Codebert: A pre-trained model for programming and natural languages. Findings of the Association for Computational Linguistics: EMNLP 2020, 1536–1547. https://doi.org/10.18653/v1/2020.findings-emnlp.139

[16] Hao, Y., Tang, Z., Tian, Y., Zhang, Y., & Zhou, Z. (2024). AdamW. Cornell Optimization. Retrieved November 21, 2025