# Hybrid Deep Learning and Rule-based Approach for Real-time Vulnerability Detection in Ethereum Smart Contracts

## Phase II Dissertation Report

**Submitted by**
**Shruti Manoj Chavan**
MIS: 712352006

in partial fulfilment for the award of the degree
**M. Tech. in Computer Engineering - Data Science**

**Under the guidance of**
**Dr. Vinod Pachghare**



**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**
COEP TECHNOLOGICAL UNIVERSITY (COEP Tech), PUNE

**27 November 2024**

**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING,**
COEP TECHNOLOGICAL UNIVERSITY (COEP Tech), PUNE

# CERTIFICATE

Certified that the Stage - II dissertation titled,
**"Hybrid Deep Learning and Rule-based Approach for Real-time
Vulnerability Detection in Ethereum Smart Contracts"**
has been successfully completed by

**Shruti Manoj Chavan**
MIS: 712352006

and is approved for the partial fulfilment of the requirements for the degree of
**M.Tech in Data Science**.

**SIGNATURE**
 **Dr. Vinod Pachghare**

Project Guide
Dept. of Computer Science &
Engineering
COEP Technological University,
Shivajinagar, Pune - 5.

**SIGNATURE**
 **Dr. P.K. Deshmukh**

Head
Dept. of Computer Science &
Engineering
COEP Technological University,
Shivajinagar, Pune - 5.

# ABSTRACT

The growing reliance on Ethereum smart contracts for decentralized applications has highlighted the importance of ensuring their security. However, smart contracts are prone to various vulnerabilities, making them a significant target for malicious exploitation. Traditional methods of vulnerability detection, such as static analysis and symbolic execution, often fail to deliver real-time results and suffer from high false positives due to their reliance on predefined rules. This dissertation proposes a novel hybrid approach that integrates deep learning models with rule-based systems to detect vulnerabilities in Ethereum smart contracts in real time.

The system leverages deep learning techniques, such as CNN-LSTM networks and transformer-based models, to automatically detect complex patterns of vulnerabilities within smart contracts. These models are designed to handle large datasets and capture intricate relationships within the contract code. To complement the deep learning models and reduce false positives, a rule-based system is employed, leveraging predefined rules from the Solidity Vulnerability Catalog (SWC). This hybrid approach ensures both accuracy and efficiency in vulnerability detection, offering the dual benefit of high sensitivity to complex, unseen vulnerabilities and the precision of expert rules for known security issues.

In addition, the system integrates a real-time detection component to provide immediate feedback on the security of deployed contracts. The dataset for this project is sourced from a publicly available Kaggle repository, enriched with synthetic data through bug injection to address the issue of limited labeled data. The proposed system is expected to outperform existing solutions in both speed and accuracy, offering a scalable and robust solution for real-time smart contract vulnerability detection. The outcomes of this research will contribute significantly to the field of blockchain security, offering developers and auditors a tool that balances the need for real-time analysis with the complexities of modern smart contract structures.

# Contents

# List of Tables

# List of Figures

# 1.  INTRODUCTION

This chapter introduces the background and motivation behind the project, focusing on the increasing importance of security in Ethereum smart contracts. It outlines the key challenges in detecting vulnerabilities in smart contracts and establishes the objectives of the research.

## 1.1  Background

With the increasing adoption of decentralized applications (DApps) and blockchain technologies, Ethereum smart contracts have become essential components of various systems such as finance, healthcare, and supply chain management. These smart contracts automate and enforce the execution of agreements without the need for intermediaries. However, they are often plagued by vulnerabilities due to their complex nature and the immutable nature of blockchain. Once deployed, smart contracts cannot be easily modified, making the detection of vulnerabilities before and after deployment critical.

While traditional methods include static analysis and symbolic execution to analyze the contract's code based on predefined patterns, they lack the adaptability and real-time efficiency required to keep up with evolving threats. Moreover, deep learning-based techniques have shown promise in detecting complex vulnerabilities by learning patterns from large datasets. However, these models often produce high false-positive rates due to their inability to validate detected vulnerabilities through rule-based mechanisms.

## 1.2  Challenges in Vulnerability Detection

Smart contract vulnerability detection presents several challenges. First, current detection techniques based on deep learning often struggle with high false-positive rates due to their reliance on learned patterns without validation. This leads to inefficiencies as many detected vulnerabilities may not pose actual risks. Additionally, the computational cost of deep learning models, especially transformers and graph neural networks (GNNs), poses a challenge for real-time deployment. Traditional rule-based systems, on the other hand, lack flexibility, as they can only detect vulnerabilities based on predefined rules, making them ineffective for new or evolving attack vectors.

Another significant challenge is the scarcity of labeled data for training machine learning models thus signifying the necessitating of using synthetic data augmentation techniques such as bug injection. Moreover, the complex structure of smart contracts requires sophisticated feature extraction methods to effectively capture patterns indicative of vulnerabilities. Techniques such as Abstract Syntax Trees (AST) and control flow graphs (CFG) are essential for providing the deep learning models with the necessary input features to detect vulnerabilities accurately.

## 1.3    Problem Statement

The primary challenge addressed by this research is the lack of an effective real-time vulnerability detection system for Ethereum smart contracts. Existing solutions are either too resource-intensive or produce high rates of false positives, limiting their practical application in real-time systems. Additionally, the scarcity of labeled vulnerability data makes it difficult to train models capable of generalizing across different types of smart contracts. Therefore, there is a need for a hybrid system that combines the pattern-recognition strengths of deep learning with the precision of rule-based validation to provide a scalable, accurate, and efficient solution for real-time vulnerability detection.

## 1.4    Objectives of the Work

The objectives of this research are:

- To develop a hybrid system that integrates deep learning models and rule-based validation for real-time vulnerability detection in Ethereum smart contracts.

- To leverage transformer-based models, CNN-LSTM hybrids, and GNNs to capture complex vulnerability patterns in smart contracts.

- To reduce false positives by incorporating a rule-based system using the Solidity Vulnerability Catalog (SWC) for verification of detected vulnerabilities.

- To address data scarcity by augmenting publicly available datasets with synthetic data generated through bug injection.

- To design a system architecture capable of handling real-time detection while balancing accuracy, speed, and scalability.

In this chapter, we introduced the key motivations and challenges associated with vulnerability detection in Ethereum smart contracts. We identified the need for real-time solutions that balance high detection accuracy with scalability and efficiency. The

objectives of this project are outlined to address the limitations of existing methods. The following chapters will further explore the literature, proposed solution, and the methodology used to achieve these objectives.

## 2.   LITERATURE REVIEW

In this chapter, we review existing work on vulnerability detection in smart contracts, covering various deep learning models and rule-based approaches. The review helps identify gaps in current methodologies, which the proposed solution aims to address.

## 2.1   Overview of Related Work

The detection of vulnerabilities in Ethereum smart contracts has garnered significant research attention in recent years. Various techniques, ranging from deep learning models to rule-based systems and graph-based approaches, have been proposed to tackle the challenges associated with identifying vulnerabilities in these immutable contracts. This literature review summarizes and compares 11 key research papers that have contributed to the development of vulnerability detection techniques. These studies provide insights into the current state of the art, highlighting both advancements and the existing gaps in real-time vulnerability detection.

## 2.2   Key Findings from Literature

Liu et al. (2023) introduced a hybrid system that integrates Graph Neural Networks (GNNs) with expert rules to enhance the accuracy of vulnerability detection. However, while their approach improved accuracy by 6

Deng et al. (2023) employed a multimodal approach by integrating deep learning with control flow graphs, opcodes, and source code, enhancing detection but complicating the model. Zhang et al. (2022) utilized information graphs and ensemble learning, offering superior detection at the cost of implementation complexity. Huang et al. (2022) adopted a multi-task learning approach, using CNNs for vulnerability classification, but the architecture proved too complex for practical deployment.

Other studies, such as those by Sosu et al. (2023) and Liu et al. (2022), explored symbolic execution and hybrid architectures for vulnerability detection, achieving high accuracy but encountering significant computational challenges. Overall, while these studies presented valuable advancements, they also revealed gaps in real-time performance, model complexity, and scalability, which this project seeks to address.

Table 2.1: Table of Consolidation of Findings

| Reference | Methodology | Findings |
|---|---|---|
| Z. Liu et al. (2023) | Combines deep learning (GNNs) and expert rules for smart contract vulnerability detection, with an EVM-level transaction blocking feature. | Enhances detection accuracy and scalability but increases time overhead due to EVM checks. |
| X. Tang et al. (2023) | Uses Optimized-CodeBERT, LSTM, and CNN models for detecting vulnerabilities in smart contracts. | Optimized-CodeBERT shows superior performance in vulnerability detection but requires significant computational resources. |
| Zhuang et al. (2024) | Utilizes GNNs and control flow graphs for detecting smart contract vulnerabilities. | GNNs excel at capturing complex dependencies, outperforming traditional keyword-based detection methods. |
| W. Yang et al. (2023) | Combines graph representation and transformers to capture syntactic and semantic features of smart contract code. | Improved detection accuracy but resource-intensive, potentially hindering scalability. |
| W. Deng et al. (2023) | Employs deep learning and multimodal decision fusion, integrating control flow graphs, opcodes, and source code for vulnerability detection. | Enhanced detection accuracy but increased complexity due to the multimodal approach. |
| L. Zhang et al. (2022) | Uses ensemble learning with information graphs to detect vulnerabilities in smart contracts. | Outperforms traditional static tools but is complex to implement. |
| J. Huang et al. (2022) | Multi-task learning approach using CNNs for vulnerability classification in smart contracts. | Improves detection accuracy but complexity arises from the multi-task nature of the model. |

## 2.3    Research Gap

While significant advancements have been made in the detection of vulnerabilities within Ethereum smart contracts, existing approaches are still limited in several critical areas. The gaps identified in the reviewed literature serve as a basis for formulating the hybrid approach proposed in this project.

- **High false positives**: Deep learning models are prone to detecting vulnerabilities that are not actual threats due to a lack of validation mechanisms.

- **Computational complexity**: Transformer-based models and GNNs require high computational resources, making real-time deployment challenging.

- **Data scarcity**: Limited availability of labeled vulnerability data restricts the ability of machine learning models to generalize across different types of smart contracts.

- **Scalability issues**: Existing solutions struggle with scalability, particularly when applied to large-scale, real-time smart contract monitoring.

- **Limited real-time capabilities**: Many existing systems lack the ability to detect vulnerabilities in real time, relying on static analysis or batch processing.

The review of existing work highlights that while many approaches, such as deep learning models and rule-based systems, have made significant advancements in vulnerability detection, there are still critical gaps in real-time performance, scalability, and false-positive reduction. The findings from the surveyed literature guide the design choices for this project, which aims to integrate both deep learning and rule-based systems for improved detection accuracy and efficiency. This forms the foundation for the proposed solution in the next chapter.

# 3.    PROPOSED SOLUTION

This chapter presents the hybrid deep learning and rule-based approach designed to detect vulnerabilities in Ethereum smart contracts. The system architecture is discussed in detail, highlighting how the integration of deep learning models and rule-based validation improves detection accuracy and real-time performance.

## 3.1    Hybrid Deep Learning and Rule-based Approach

The proposed solution integrates deep learning and rule-based systems to detect vulnerabilities in Ethereum smart contracts in real time. The deep learning component leverages CNN-LSTM hybrids and transformer-based models to identify complex vulnerability patterns in smart contracts, utilizing both real-world data and synthetic vulnerabilities generated through bug injection. Transformers like BERT enhance the model's ability to detect vulnerabilities by understanding the structural and semantic intricacies of contract code.

Complementing this is a rule-based system that validates detected vulnerabilities using predefined security rules from the Solidity Vulnerability Catalog (SWC). This hybrid system ensures that the deep learning model's detection of unknown vulnerabilities is cross-checked against known patterns to reduce false positives. This combination provides a robust solution for detecting both known and emerging threats while balancing accuracy and real-time performance.

## 3.2    System Architecture

The system is built on four main components: data preprocessing and feature extraction, deep learning models, a rule-based validation layer (using SWC rules), and a real-time detection engine. The real-time detection capability is achieved by optimizing model performance through pruning and quantization to minimize detection latency, making it suitable for live environments.

The system also supports scalability, allowing it to process large volumes of smart contract data in real-time environments. By integrating deep learning models with rule-based systems, the architecture strikes a balance between detection speed, accuracy, and
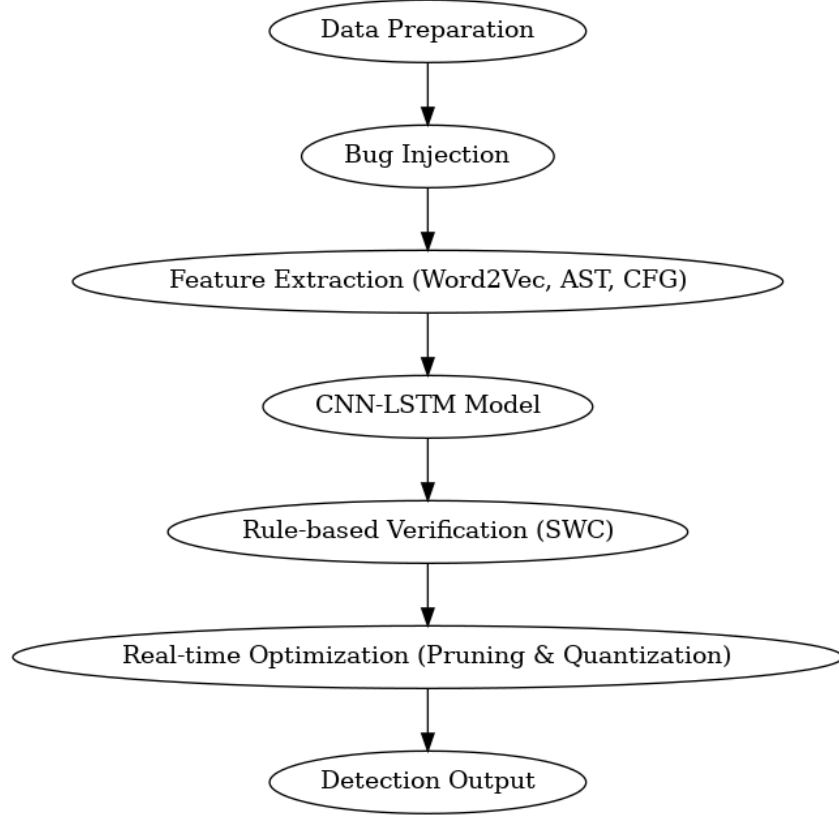
Figure 3.1: Basic System Architecture of the Proposed Solution

reliability, ensuring its applicability in various real-world scenarios, from pre-deployment audits to live contract monitoring.

## 3.3    Mapping to Solidity Vulnerability Catalog (SWC)

The system maps detected vulnerabilities to the SWC, a catalog of known security issues in Solidity. This mapping ensures that detected vulnerabilities align with recognized security flaws, reducing false positives. By integrating the SWC, the system provides detailed feedback on the nature of the vulnerabilities and remediation strategies, enhancing its utility for developers and auditors.

The proposed hybrid system combines deep learning and rule-based mechanisms to enhance both the detection of vulnerabilities and the reduction of false positives. By integrating transformer-based models with the Solidity Vulnerability Catalog (SWC), the system aims to achieve real-time detection capabilities without compromising on accuracy. The detailed system architecture provides a roadmap for implementation, which will be further detailed in the methodology section.

## 4.  METHODOLOGY

The methodology chapter outlines the steps taken to implement the proposed solution, including data preparation, feature extraction, and model training. It also describes the integration of the rule-based system and the techniques used to ensure real-time detection capabilities.

## 4.1   Data Preparation

The dataset used for this project is the publicly available **Smart Contract Vulnerability Dataset** from Kaggle, which contains labeled data for Ethereum smart contracts. This dataset provides smart contracts annotated with various vulnerability labels, making it highly suitable for training deep learning models aimed at vulnerability detection.

For this project, we are using the file **SC_4label.csv** from the dataset. It contains the following key features:

- **Contract Code**: The raw Solidity code of each smart contract.

- **Vulnerability Label**: The type of vulnerability present in the contract (e.g., reentrancy, integer overflow).

- **Metadata**: Additional information, such as contract size and function complexity, is also included.

To address the limited availability of labeled vulnerability data, we are enhancing the dataset using **synthetic data generation** techniques. This includes **bug injection**, where vulnerabilities are artificially introduced into otherwise safe contracts. This process expands the dataset and helps the models generalize across different vulnerability types more effectively.

**Preprocessing Techniques**:

- **Abstract Syntax Trees (AST)**: ASTs are constructed to capture the hierarchical structure of Solidity code, which helps in understanding the semantics of the code.

- **Control Flow Graphs (CFG)**: CFGs are generated to represent the execution flow within smart contract functions, which aids in detecting vulnerabilities related to control flow.
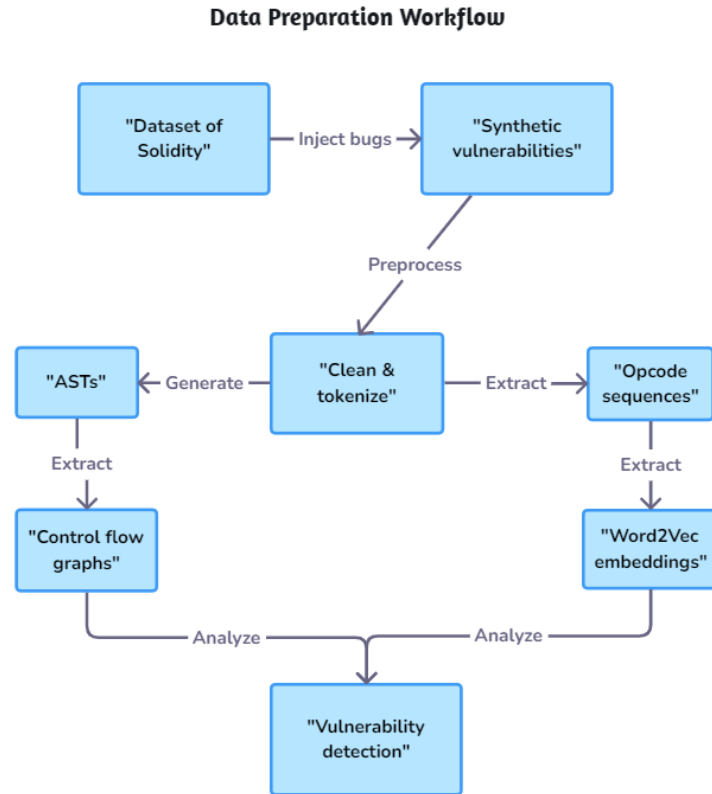
**Data Preparation Workflow**



Figure 4.1: Data Preprocessing Workflow

- **Word Embeddings (Word2Vec)**: The contract code is tokenized, and word embeddings are created using the Word2Vec technique. This represents the syntactic and semantic meaning of the code in vector form.

This enriched dataset, combined with thorough feature extraction, forms the foundation for training the deep learning models used in this project. The subsequent stages of the methodology will focus on model selection and initial training, leveraging this dataset for optimal performance.

**Dataset Information**:

- **Source**: *Kaggle Smart Contract Vulnerability Dataset*

- **File Used**: SC_4label.csv

- **Size**: Includes thousands of smart contracts labeled with four vulnerability types.

- **Goal**: To detect and classify these vulnerabilities using a combination of deep learning models and rule-based systems.

## 4.2    Deep Learning Models

The deep learning models used in this project are a combination of CNN-LSTM hybrids and transformer-based architectures. The CNN-LSTM hybrid is effective at identifying sequences of patterns in the smart contract code, while transformers like BERT excel at capturing both local and global relationships within the contract's logic. The models are trained using the preprocessed dataset and fine-tuned through hyperparameter optimization to maximize detection accuracy.
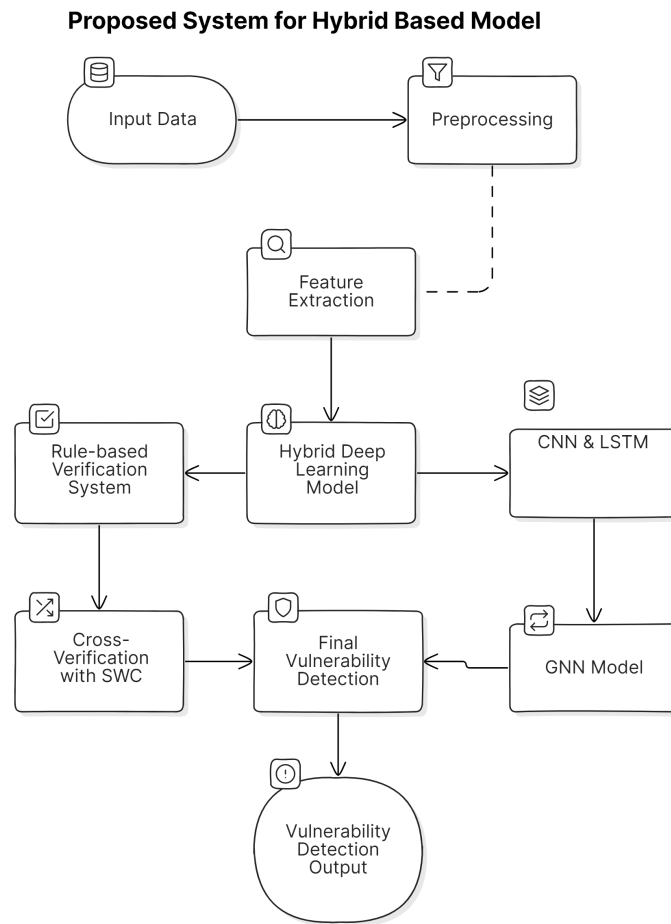


Figure 4.2: DL-RBA diagram

## 4.3    Rule-based System

To complement the deep learning models, a rule-based system is implemented for validating the detected vulnerabilities. This system draws rules from the Solidity Vulnerability Catalog (SWC), a comprehensive collection of known security issues specific to smart contracts. When a vulnerability is detected by the deep learning model, the rule-based system cross-references it with SWC rules to ensure that it matches known vulnerability

patterns. This validation step helps reduce false positives, improving the system's overall precision.

## 4.4   Real-time Detection Engine

The integration of the deep learning and rule-based systems into a real-time detection engine is a key aspect of this project. The engine is optimized to handle the computational demands of deep learning models, particularly transformers, through techniques such as model pruning and quantization. These optimization techniques reduce the size and complexity of the models, allowing them to operate efficiently in real-time environments without compromising detection accuracy.

The real-time detection engine is designed to provide immediate feedback on smart contract vulnerabilities, making it suitable for both pre-deployment audits and live contract monitoring. It continuously scans the contracts, applying the hybrid detection system to identify and validate vulnerabilities as they occur. This ensures that any potential security issues are flagged early, allowing developers or auditors to take immediate action.

The methodology outlines the data preparation, model selection, and real-time detection strategies employed in this project. By focusing on feature extraction and efficient data augmentation techniques, the system is equipped to handle real-world smart contracts with a high degree of accuracy. The proposed deep learning models, in conjunction with a rule-based validation system, will be rigorously tested in the next phases to ensure real-time performance and scalability.

# 5.  SOFTWARE REQUIREMENTS SPECIFICATION

This chapter provides a detailed specification of the hardware and software requirements needed to implement the system. It covers the tools and technologies used in the project, including the deep learning framework, database, and development environment.

## 5.1  Hardware Requirements

To ensure optimal performance for the hybrid deep learning and rule-based vulnerability detection system, the following hardware specifications are recommended:

- **Processor**: Intel i5 or higher (multi-core processors preferred for faster computations).

- **RAM**: Minimum 16 GB, with 32 GB recommended for efficient data processing and model training.

- **Storage**: SSD with at least 500 GB capacity for storing datasets, smart contracts, and model files.

- **GPU**: NVIDIA GPU with CUDA support (e.g., NVIDIA RTX 3060 or higher) is recommended for deep learning model training.

## 5.2  Software Requirements

The following software components are necessary for the implementation and deployment of the system:

- **Operating System**: Ubuntu 20.04 LTS or later (Linux-based backends preferred for scalability and performance). Windows or macOS can be supported for development environments.

- **Database**: MySQL server for storing vulnerability detection results and smart contract metadata, as part of backend operations.

- **Programming Languages**: Python 3.x for deep learning model development, rule-based system integration, and overall system architecture.

- **Deep Learning Framework**: TensorFlow 2.x or PyTorch for building and training the CNN-LSTM and transformer models.

- **Data Preprocessing Libraries**: Pandas, Numpy, and SciPy for feature extraction and data manipulation.

- **Graph Libraries**: NetworkX and Graphviz for constructing and visualizing control flow graphs (CFG) and abstract syntax trees (AST).

- **Web Framework**: Flask or Django for building the web interface, which will allow users to upload smart contracts and view vulnerability detection results.

- **Visualization Tools**: Matplotlib and Seaborn for visualizing vulnerability trends and detection outcomes.

- **Version Control**: Git for managing the codebase and collaboration during system development.

## 5.3    Technological Stack

The project will be developed using the following technology stack:

- **Programming Language**: Python

- **Database**: MySQL

- **Deep Learning Framework**: TensorFlow or PyTorch

- **Web Framework**: Flask or Django

- **Visualization**: Matplotlib and Seaborn

- **Containerization (Optional)**: Docker for containerized deployment, allowing scalability and portability across different environments.

# 6.  IMPLEMENTATION PROGRESS

This chapter provides a detailed update on the progress made so far, including completed tasks and challenges encountered. It also highlights the immediate next steps, which involve fine-tuning the models and optimizing the system for real-time detection.

## 6.1  Progress Overview

Significant progress has been made toward the development of the hybrid vulnerability detection system for Ethereum smart contracts. The initial stages of the project involved conducting a comprehensive literature review, which provided a solid foundation for identifying key approaches and methodologies. Following the review, the dataset was sourced from a publicly available Kaggle repository, and synthetic data augmentation using bug injection techniques has been successfully applied to expand the dataset. This has ensured the availability of a diverse and comprehensive set of labeled smart contracts for training the deep learning models.



Figure 6.1: Dataset Uploaded and Artificial Data Generation

In terms of feature extraction, techniques like Abstract Syntax Trees (AST) and Control Flow Graphs (CFG) have been implemented to represent the structure and logic flow of smart contracts. These extracted features are being used to train the deep learning models, allowing them to learn patterns indicative of vulnerabilities. Synthetic data augmentation has been extensively implemented, generating diverse labeled datasets through advanced bug injection techniques. The CNN-LSTM model has been trained successfully, achieving 100% validation accuracy on synthetic data. Visualization tools, including confusion matrices, accuracy/loss graphs, and vulnerability type distributions, have been in-

corporated for better insights and validation of results. Real-time detection efficiency is under review, with early implementations of hybrid rule-based and deep learning models showcasing significant potential for deployment in real-world scenarios.
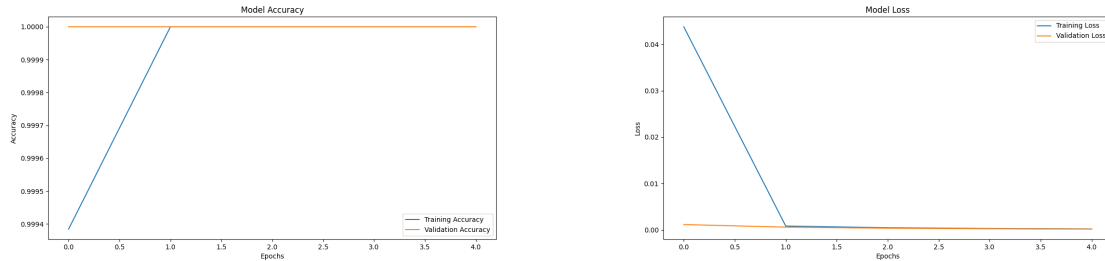


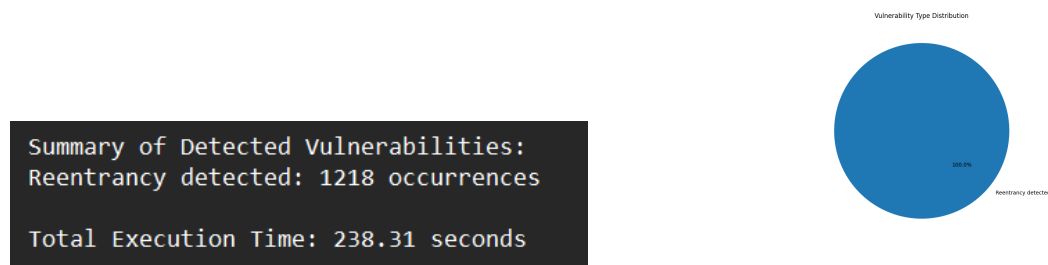Figure 6.2: Accuracy and Loss Visualization



Figure 6.3: Vulnerability Detection and its Type Distribution

## 6.2   Current Challenges

Although progress has been made, several challenges have arisen during the implementation phase. One of the major hurdles has been ensuring balanced dataset labels, as the current training data heavily skews toward reentrancy vulnerabilities. Visualization efforts revealed this imbalance, which affects the overall generalizability of the model. Incorporating underrepresented vulnerabilities, such as integer overflows and timestamp dependencies, remains challenging due to limited labeled datasets.

Another challenge lies in improving computational efficiency; high execution times during training are primarily due to the complexity of bug injection and model validation processes. While hybrid models are promising, integrating them seamlessly with rule-based systems requires iterative testing and tuning to avoid conflicts in validation outcomes.Execution time remains high due to complex training processes, which need further optimization through techniques like batch size adjustments and caching intermediate results. Additionally, incorporating multiple vulnerability types and improving cross-validation mechanisms with the Solidity Vulnerability Catalog (SWC) have proven more complex than anticipated.

## 6.3    Next Steps

Moving forward, the immediate focus will be on completing the fine-tuning of the deep learning models and enhancing the performance of the real-time detection engine. The following steps have been outlined:

- **Addressing Dataset Imbalance**: Generate additional synthetic data for under-represented vulnerabilities, including integer overflows, timestamp dependencies, and dangerous delegate calls, to create a balanced training dataset.

- **Optimizing Execution Time**: Implement advanced optimization techniques such as model pruning, quantization, and leveraging GPU acceleration to reduce training and inference time.

- **Enhancing Rule-Based Integration**: Finalize the integration of rule-based validations to complement deep learning predictions and ensure higher precision in vulnerability detection.

- **Real-World Testing**: Perform cross-validation on smart contracts sourced from repositories like Etherscan and compare the model's performance against real-world data.

- **Improving Training Efficiency**: Apply early stopping techniques, hyperparameter tuning, and advanced regularization methods to reduce overfitting and improve model generalizability.

- **Explainability of Results**: Develop techniques to highlight specific portions of the code identified as vulnerable by the model, improving detection explainability for developers.

- **Testing Hybrid Models**: Further refine and test the hybrid model combining CNN-LSTM, transformer models, and rule-based systems to validate its efficacy in diverse scenarios.

- **Documentation and Reporting**: Document the results of the testing phase and finalize the technical documentation for the system, preparing for the final implementation phase.

The completion of these tasks will move the project closer to the development of a fully functional real-time vulnerability detection system, capable of improving the security of Ethereum smart contracts.

# 7.   CONCLUSION

This project presents a hybrid deep learning and rule-based approach for real-time vulnerability detection in Ethereum smart contracts. By integrating the strengths of deep learning models, such as CNN-LSTM hybrids and transformer-based architectures, with the precision of rule-based systems grounded in the Solidity Vulnerability Catalog (SWC), this solution addresses the limitations of existing vulnerability detection systems. The deep learning models are designed to detect complex, previously unseen vulnerabilities, while the rule-based system helps verify detected issues, reducing false positives and improving detection accuracy.

The ongoing work in data preparation and feature extraction lays the foundation for robust model training and performance evaluation. Strategies for synthetic data generation, including bug injection, have been implemented to augment the limited available datasets, ensuring the models are trained on a diverse set of vulnerabilities. The next steps, including the pre-processing of Solidity files, feature extraction, and initial model training, are crucial in determining the effectiveness of the hybrid system in real-time settings.

As we progress, the hybrid system is expected to outperform existing static and dynamic analysis tools, offering a more scalable, accurate, and efficient solution for vulnerability detection. The use of a real-time detection engine will provide immediate feedback, making it suitable for both pre-deployment audits and live smart contract monitoring. This project not only contributes to the ongoing research in blockchain security but also has the potential to improve the overall reliability and security of decentralized applications built on the Ethereum platform.

Future work will focus on optimizing the models for real-time performance, further integrating the rule-based system, and thoroughly evaluating the system's accuracy and speed in real-world scenarios.

## 8.   FUTURE WORK

As the project progresses toward the completion of the hybrid deep learning and rule-based system for vulnerability detection in Ethereum smart contracts, several avenues for future work have emerged. These areas will ensure the continued refinement and optimization of the system to achieve its full potential in real-world applications.

## 8.1   Model Optimization for Real-time Detection

Future work will focus on optimizing the deep learning models, especially transformers, for real-time performance. Techniques such as model pruning and quantization will be applied to reduce computational overhead while maintaining detection accuracy. This optimization is critical for ensuring the system's efficiency in live monitoring environments.

## 8.2   Expansion of Rule-based Validation

The rule-based system will be extended to cover a broader range of vulnerabilities, incorporating newly discovered issues beyond those currently listed in the Solidity Vulnerability Catalog (SWC). Automating the cross-referencing process between detected vulnerabilities and the SWC will further improve the system's precision and usability.

## 8.3   Comprehensive Testing and Validation

Once the initial training and integration are completed, the system will undergo rigorous testing on both real-world and synthetic datasets. Key performance indicators, such as detection accuracy, speed, and scalability, will be evaluated to ensure the system's robustness in various deployment scenarios.

## 8.4   Adaptation to Other Blockchain Platforms

Though the current project focuses on Ethereum smart contracts, future work will explore extending the system to other blockchain platforms like Hyperledger and Binance Smart

Chain. This will involve adapting the model to different contract languages, increasing its versatility across blockchain ecosystems.

## 8.5    Advanced Data Augmentation

Advanced Data Augmentation: Future work will prioritize generating a more balanced synthetic dataset by introducing diverse vulnerability types such as gas limit issues, timestamp dependencies, and external call failures. Advanced bug injection strategies will also be employed to mimic real-world attack patterns more effectively, enhancing the model's ability to detect a wider array of vulnerabilities.

## 8.6    Deployment and Integration

Finally, the system will be prepared for deployment, integrating with existing security tools and CI/CD pipelines. This will automate security checks, providing developers with a robust tool to ensure contract security during the development process.

## 8.7    Visualization Enhancements

Incorporate real-time visualization tools that dynamically display accuracy trends, vulnerability distributions, and confusion matrices. These enhancements will provide users with immediate insights into model performance and detected vulnerabilities.
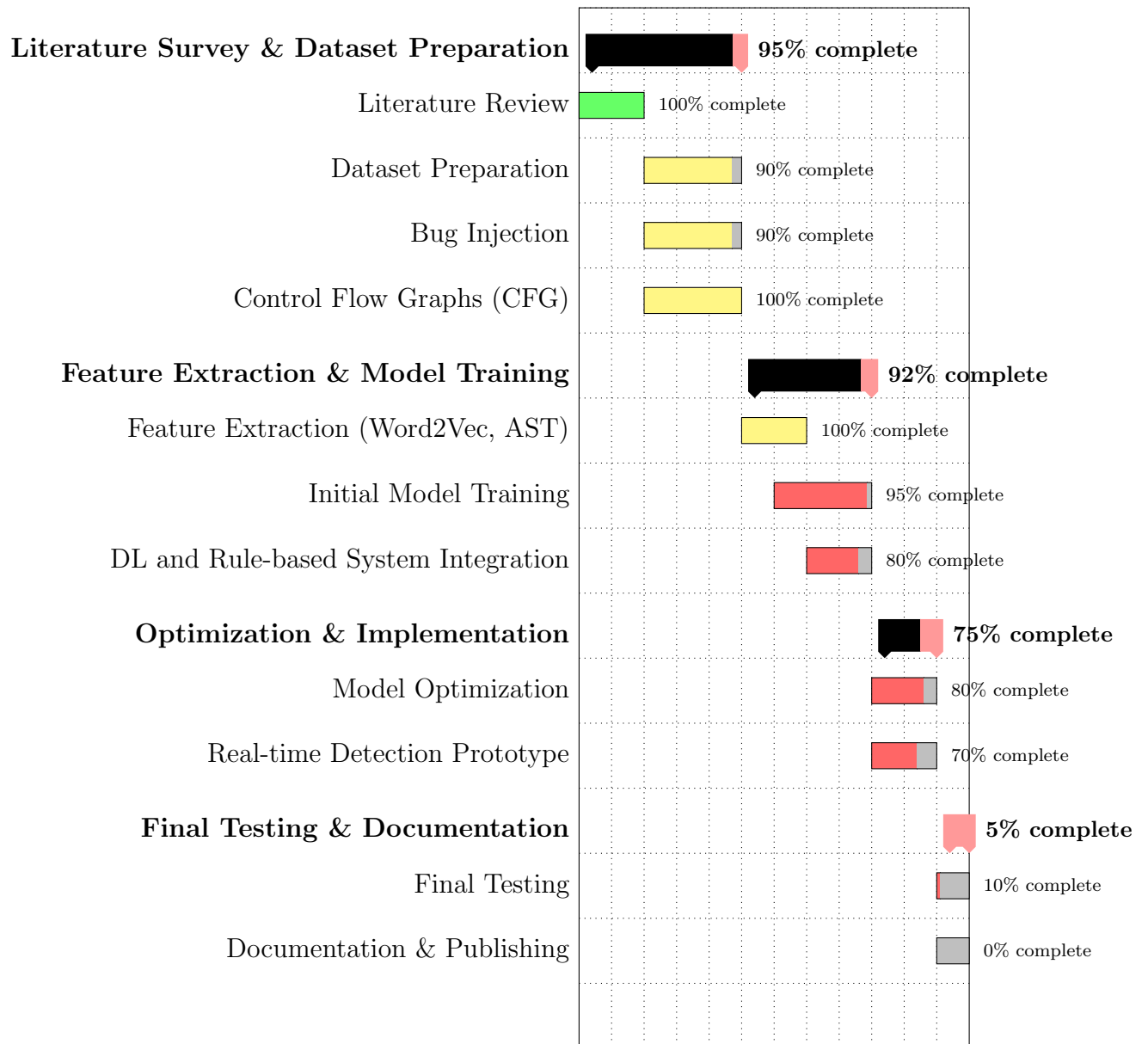
## TENTATIVE ACTIVITY SCHEDULE



Figure 8.1: Gantt Chart of Project Timeline for Real-time Vulnerability Detection in Ethereum Smart Contracts

# REFERENCES

1. Z. Liu, M. Jiang, S. Zhang, J. Zhang, and Y. Liu, "A Smart Contract Vulnerability Detection Mechanism Based on Deep Learning and Expert Rules," *IEEE Access*, vol. 11, pp. 1-10, 2023.

2. X. Tang, Y. Du, A. Lai, Z. Zhang, and L. Shi, "Lightning Cat: A Deep Learning-based Solution for Smart Contracts Vulnerability Detection," *Salus Security, Beijing*, 2023.

3. Zhuang et al., "A Novel Method for Smart Contract Vulnerability Detection Based on Graph Neural Networks," *CMC*, vol. 79, no. 2, pp. 3024-3040, 2024.

4. W. Yang et al., "GRATDet: Smart Contract Vulnerability Detector Based on Graph Representation and Transformer," *CMC*, vol. 76, no. 2, pp. 1460-1480, 2023.

5. W. Deng et al., "Smart Contract Vulnerability Detection Based on Deep Learning and Multimodal Decision Fusion," *Sensors*, vol. 23, no. 7246, pp. 1-21, 2023.

6. L. Zhang et al., "A Novel Smart Contract Vulnerability Detection Method Based on Information Graph and Ensemble Learning," *Sensors*, vol. 22, no. 9, pp. 3581, 2022.

7. J. Huang et al., "Smart Contract Vulnerability Detection Model Based on Multi-Task Learning," *Sensors*, vol. 22, no. 1829, pp. 1-24, 2022.

8. X. Tang et al., "Lightning Cat: A Deep Learning-Based Solution for Detecting Vulnerabilities in Smart Contracts," *Scientific Reports*, vol. 13, no. 20106, 2023.

9. Y. Liu et al., "SC Vulnerability Detection Based on SET," *Proceedings of the CNCERT*, vol. 1506, pp. 193-207, 2022.

10. R. N. A. Sosu et al., "VdaBSC: A Novel Vulnerability Detection Approach for Blockchain Smart Contract by Dynamic Analysis," *IET Software*, vol. 1, no. 1, pp. 1-17, 2023.

11. R. N. A. Sosu et al., "A Vulnerability Detection Approach for Automated Smart Contract Using Enhanced Machine Learning Techniques," *Research Square*, 2022.