# Project Report

## Vulnerability Assessment & Reverse Engineering

# Automated Vulnerability Scanner with AI-Powered Classification and Remediation

**Team Members:** Muhammad Umer Farooq, Jalil Ahmad, Aleena Fatima

**Roll Numbers:** I221661, I221635, I222353

**Class:** CY-D

# Contents

# 1. Introduction

## 1.1 Background

As cyber threats grow in complexity and frequency, automated systems that not only detect but also assess and remediate vulnerabilities are critical. Traditional vulnerability scanners like OpenVAS (Greenbone Vulnerability Management) are effective at detection but lack automated reasoning and response generation capabilities.

## 1.2 Objective

This project aims to build a comprehensive vulnerability scanning framework that integrates:

- OpenVAS for detection
- Machine learning for severity classification
- Natural language generation (NLG) for remediation recommendations

The system streamlines the vulnerability management process, providing actionable insights in addition to detection.

# 2. System Overview

## 2.1 Core Modules

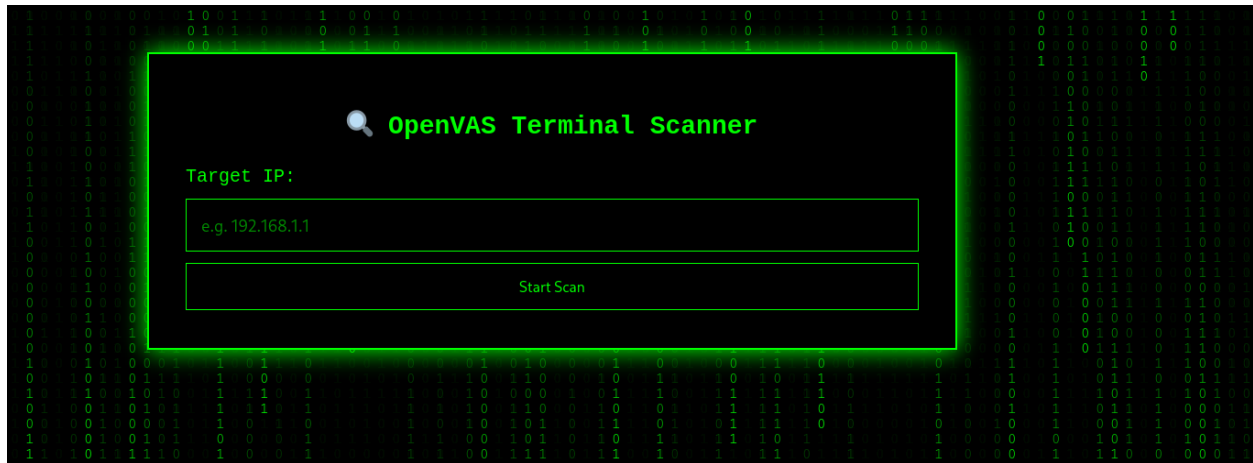The system is composed of three core modules:

| Module | Function |
| --- | --- |
| Vulnerability Scanning | Conducts system scans using OpenVAS |
| Severity Classification | Predicts vulnerability severity using ML models |
| Remediation Recommendation | Generates remediation steps via a fine-tuned T5 model |

## 2.2 Pipeline

Target IP → [OpenVAS Scan] → [Cleaned Report] → [Severity Classifier] → [Remediation Generator] → [Final Output]

# 3. Technical Components

## 3.1 Vulnerability Scanning



### 3.1.1 Tool Used: OpenVAS

OpenVAS is an open-source vulnerability scanner that performs a wide range of network checks against known CVEs.
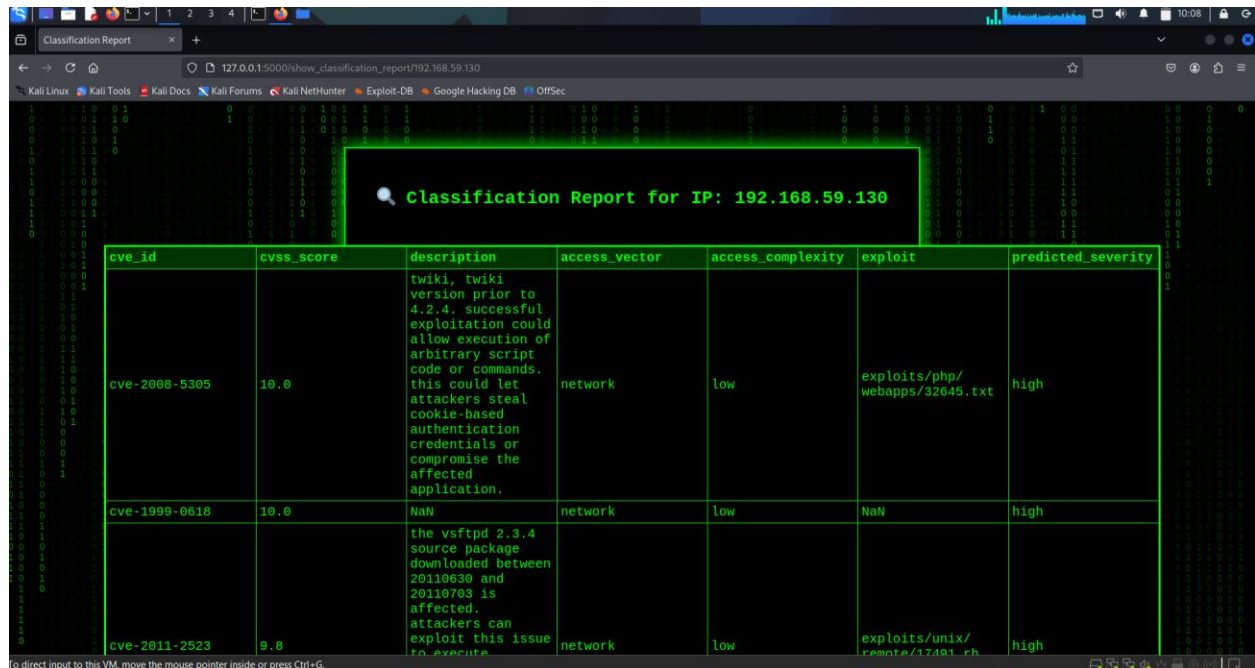
### 3.1.2 Workflow

- Scan tasks are initiated using gvm-cli via Unix socket.

- The scan is configured to use "Full and fast" policy.

- Results are fetched in CSV format and parsed using pandas.

### 3.1.3 Data Enrichment

- Extracted data includes CVE ID, CVSS score, solution, impact, and affected software.

- A secondary lookup (if available) is used to enrich the data with:

  o access_vector

  o access_complexity

  o exploit

## 3.2 Severity Classification



### 3.2.1 Problem Framing

Severity classification is treated as a supervised classification task with multi-class output (e.g., Low, Medium, High).

### 3.2.2 Model Architecture

- **Base Classifier**: Random Forest
- **Text Embedding**: SBERT (all-MiniLM-L6-v2) is used to convert descriptions to numerical vectors.

### 3.2.3 Feature Set

- Encoded categorical fields: access_vector, access_complexity, exploit
- Continuous fields: cvss_score
- Text embeddings: SBERT vectors from the vulnerability description

## 3.3 Remediation Generation



Remediation is not fetched from a database but generated using a **fine-tuned T5 (Text-to-Text Transfer Transformer)** model. The model generates contextual, human-readable remediation steps tailored to each CVE and its description.

# 4. Implementation Details

## 4.1 Backend

- **Framework**: Flask

- **Task Management**: Dictionary-based tracking of running/completed scans

- **Data Handling**: pandas, joblib, torch, transformers

## 4.2 Frontend

Implemented using HTML/CSS/JavaScript with a consistent hacker-themed visual identity (Matrix-style green-on-black). Each stage has a dedicated UI:

- index.html: Scan initiation

- results.html: Scan data

- classification.html: Severity predictions

- remediations.html: Final output with recommendations

# 5. User Workflow

**Step Action**

**1**    User enters target IP on the web interface

**2**    OpenVAS scan begins in background

**3**    Upon completion, results are shown in a table

**4**    User can trigger severity classification

**5**    Final step generates AI-powered remediation suggestions

# 6. Datasets Utilised

The data for training severity classification model is gathered from NVD database and ExploitDB for corresponding exploits. A corpus of 200,000 entries is created. Following is the glimpse of it:



The dataset for remediations is gathered through web scraping and contains remediation steps for over 14000 vulnerabilities. It is as follows:

# 7. Model Evaluation and Testing

This section presents the evaluation results for both core models developed as part of the vulnerability scanning framework: the **Severity Classification Model** and the **Remediation Generation Model**. Each model was rigorously tested to ensure accuracy, generalization, and utility in a real-world cybersecurity context.

## 7.1 Severity Classification Model

The severity classification model is a Random Forest classifier trained on enriched CVE metadata. Input features included the cvss_score, access_vector, access_complexity, and SBERT-encoded descriptions of vulnerabilities. The model was trained using 100 trees and evaluated on a held-out test set.

**Evaluation Metrics**

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Critical | 0.99 | 0.89 | 0.94 | 3489 |
| High | 0.96 | 0.99 | 0.97 | 7878 |
| Low | 1.00 | 0.95 | 0.97 | 3816 |
| Medium | 0.99 | 1.00 | 0.99 | 20834 |
| **Accuracy** | | | **0.98** | 36017 |

| Class | Precision | Recall | F1-Score | Support |
|-------|-----------|--------|----------|---------|
| **Macro Avg** | 0.98 | 0.96 | 0.97 | 36017 |
| **Weighted Avg** | 0.98 | 0.98 | 0.98 | 36017 |



## 7.2 Remediation Generation Model

The remediation model is a fine-tuned T5 transformer trained on a custom dataset containing CVE IDs, descriptions, and corresponding remediation steps. The model is tasked with generating human-readable, context-specific mitigation recommendations from structured inputs.

| Metric | Score |
|--------|-------|
| **ROUGE-1** | 0.597 |
| **ROUGE-2** | 0.566 |
| **ROUGE-L** | 0.618 |
| **ROUGE-Lsum** | 0.618 |

# 8. Limitations and Future Work

## 8.1 Current Limitations

- OpenVAS may miss zero-day vulnerabilities.
- The classifier's accuracy is bounded by the training data and SBERT embeddings.
- The T5 model may generate plausible but non-authoritative remediations.

## 8.2 Proposed Enhancements

- Integrate NVD API for real-time vulnerability enrichment.
- Replace SBERT with a domain-specific LLM fine-tuned on MITRE CVE data.
- Build a dashboard for multi-host scanning and vulnerability trend analytics.

# 9. Conclusion

This project bridges the gap between detection and response by combining the strengths of traditional vulnerability scanners and modern AI models. It provides:

- End-to-end automation

- Human-like remediation suggestions

- An intuitive UI to manage scan workflows

This system is a step toward building **smart cybersecurity assistants** that can reason and recommend—automatically.