

# Design and Implementation of Machine Learning for Network Security

- **Student ID:** 21/03/03/108
- **Supervisor:** Egr. Aishatu Ibrahim Birma

## Chapter One: Introduction

Network security is critical in today's highly connected environment, as the proliferation of devices and internet usage has led to a rise in cyber threats such as malware, phishing, and DDoS attacks <sup>1</sup>. Traditional defenses (firewalls, antivirus, etc.) often rely on static signatures and thus struggle to detect novel or polymorphic threats <sup>2</sup> <sup>3</sup>. In response, machine learning (ML) methods have been adopted in cybersecurity to identify anomalous patterns in data and adapt to evolving attacks <sup>2</sup> <sup>4</sup>. For example, anomaly detection algorithms can spot unusual network traffic, and classification models can learn features of malware files. This project combines ML with threat intelligence APIs to build an accessible security platform for under-resourced settings (Maiduguri, Nigeria). In particular, the platform leverages the VirusTotal API—which aggregates results from over 70 antivirus engines and URL scanners—to analyze files, URLs, and IPs in real time <sup>5</sup> <sup>6</sup>. By integrating VirusTotal with ML risk scoring and anomaly detection, the system provides users with faster, more accurate threat insights than signature-based tools alone <sup>7</sup> <sup>4</sup>.

The primary objectives of this work are: (1) to design a network security system that integrates the VirusTotal API with customized ML models for threat analysis, (2) to implement the platform using Python and Streamlit, and (3) to evaluate its performance. These goals involve creating a user-friendly interface for uploading files and URLs, securing the system with robust authentication, and visualizing threat data through interactive dashboards <sup>8</sup> <sup>9</sup>. Addressing cybersecurity in Maiduguri and similar regions is significant because many local organizations lack sophisticated defenses and expertise <sup>10</sup> <sup>11</sup>. An accessible, API-driven solution (built on open-source tools) can help bridge this gap, offering scalable threat detection to improve the digital safety of under-resourced users <sup>12</sup> <sup>7</sup>. The scope of the study focuses on the design and implementation of an ML-based platform with real-time VirusTotal integration, tailored for environments with limited resources and technical expertise <sup>11</sup> <sup>7</sup>.

## Chapter Two: Literature Review

**Network Security Fundamentals:** Network security aims to protect the confidentiality, integrity, and availability of data and systems. Common measures include firewalls (traffic filtering), intrusion detection/prevention (IDPS), and encryption <sup>13</sup>. Modern practice emphasizes *defense-in-depth*: multiple overlapping layers protect against sophisticated threats. For instance, firewalls may be combined with anomaly-based IDPS that monitor traffic patterns for deviations (zero-day or advanced threats) <sup>14</sup> <sup>15</sup>. However, each layer has limits. For example, traditional antivirus and firewalls rely on known threat signatures or rules, making them ineffective against novel attacks <sup>16</sup> <sup>14</sup>.

**Machine Learning in Security:** ML techniques address some of these limits by learning from data. Supervised models (e.g. SVMs, random forests, neural nets) have been trained on labeled benign vs. malicious samples to detect malware or intrusion <sup>17</sup> <sup>4</sup> . Unsupervised methods (e.g. clustering, anomaly detection) identify outliers without needing labels, useful for spotting unusual network behavior <sup>18</sup> <sup>19</sup> . Researchers have shown ML can detect new threats and adapt over time <sup>19</sup> <sup>20</sup> . For example, ML has been used for endpoint protection, analyzing device telemetry to flag compromises <sup>21</sup> . However, ML models face challenges (data quality, training overhead, adversarial evasion) <sup>19</sup> , so practical systems often combine ML with other tools.

**Threat Types:** Common cyber threats include **malware** (viruses, ransomware, etc.), which can steal data or disrupt operations (particularly harmful in low-resource settings) <sup>22</sup> <sup>7</sup> ; **phishing**, which exploits human weakness to gain credentials <sup>23</sup> ; **DDoS attacks**, which flood networks to cause downtime <sup>24</sup> ; and **APTs** (advanced persistent threats) that stealthily compromise sensitive systems <sup>25</sup> . These threats impact organizations severely, especially where defensive infrastructure is weak. A layered approach – combining traditional tools with ML analysis and threat intelligence – is recommended to mitigate such varied threats <sup>26</sup> <sup>27</sup> .

**Threat Intelligence APIs (VirusTotal, etc.):** Online multi-scanner platforms like VirusTotal emerged to overcome limitations of single-engine antivirus <sup>4</sup> <sup>6</sup> . VirusTotal aggregates scan results from numerous antivirus engines and URL/domain blacklists, giving a comprehensive “consensus” verdict on submitted files, URLs, or IP addresses <sup>4</sup> <sup>6</sup> . Its API enables automated querying: files and URLs can be uploaded and then repeatedly polled until a complete multi-engine analysis report is ready <sup>28</sup> <sup>29</sup> . In practice, VirusTotal provides immediate, detailed insights (e.g. number of antivirus detections, malicious domains in a URL’s history) that greatly enhance threat detection <sup>28</sup> <sup>4</sup> . Security systems use these APIs to supplement their own detection. For example, the platform developed here sends user-submitted files/URLs to VirusTotal and retrieves multi-engine analysis results (malicious/suspicious counts, detailed engine verdicts) for further processing <sup>30</sup> <sup>28</sup> .

However, reliance on external APIs brings constraints. The VirusTotal API imposes strict rate limits on free-tier usage, which can delay high-volume scanning if not managed properly <sup>27</sup> <sup>31</sup> . Also, by design VirusTotal can only flag what its engines know; brand-new, undetected threats may slip through. Therefore, many systems combine VirusTotal data with additional ML analysis. For instance, after retrieving VirusTotal results, our platform computes a **risk score** (a 0–100 scale based on counts of malicious/suspicious detections) and checks for file anomalies (e.g. suspicious extensions or sizes) to catch malware traits not directly captured by signatures <sup>32</sup> <sup>33</sup> . These ML-derived insights complement the raw VirusTotal output and help quantify threat severity beyond simple signature counts.

**Related Platforms:** Previous work has built streamlined UIs for security analysis using frameworks like Streamlit. Streamlit enables rapid development of interactive data apps with minimal frontend coding <sup>34</sup> <sup>4</sup> . Its simplicity is well-suited for providing clear dashboards of security scans. Other systems (e.g. [source withheld]) similarly integrate threat intel APIs with visualization libraries (like Plotly) to present scan results graphically <sup>35</sup> <sup>4</sup> . The literature shows that combining APIs and ML improves threat detection: for example, specialized SIEM tools often embed VirusTotal queries and ML analytics for richer alerts <sup>36</sup> . This project draws on these ideas to create a user-friendly, intelligent security platform tailored to resource-limited contexts.

## Chapter Three: Methodology

### 3.1 System Analysis and Design

We first analyzed existing solutions and their shortcomings in our target context. In Maiduguri and similar regions, many systems still rely on signature-based detection, which fails against novel threats <sup>37</sup>. They often lack integration with broad threat databases like VirusTotal, limiting their visibility into new malware and malicious domains <sup>37</sup>. High-end security tools also tend to be cost-prohibitive and require skilled admins, which is challenging for under-resourced organizations <sup>38</sup>. These gaps motivated our new design.

The proposed platform overcomes these issues by integrating the VirusTotal API with ML in a streamlined, affordable solution. Key design decisions include: - **Advanced Threat Detection:** By using VirusTotal's multi-engine scans along with ML algorithms, the system can detect both known and emerging threats more robustly than signature-only tools <sup>7</sup> <sup>32</sup>. - **Cost-Effectiveness:** We build on open-source tools (Python, Streamlit, SQLite, etc.) to minimize costs and ease deployment in constrained settings. - **Real-Time Analysis:** The system is designed for asynchronous, real-time scanning. Uploaded files or URLs are immediately sent to VirusTotal, and results are polled until completion to give prompt feedback <sup>20</sup> <sup>39</sup>. - **User-Centric Interface:** A Streamlit front-end provides an intuitive UI requiring little technical expertise <sup>34</sup>. - **Scalability:** The architecture is modular, so additional analyses (e.g. more ML models or threat sources) can be integrated later with minimal changes.

Figure references (not shown here) illustrate the system's Data Flow Diagrams (DFDs) and UML diagrams. For example, a Level-0 (context) DFD shows interactions between the **User**, the **VirusTotal API**, and our platform. A Level-1 DFD breaks this down into main processes (file/URL submission, analysis, result storage). Likewise, a UML use-case diagram highlights user actions (login, scan file, scan URL, view reports) and a class diagram models core entities (User, ScanRecord, ThreatReport, etc.) <sup>40</sup> <sup>41</sup>. These designs guided the implementation structure. The database design includes tables for **users** (id, username, email, password hash, role) and **scans** (scan id, type, target, hashes, VirusTotal analysis id, results summary, timestamps) <sup>42</sup> <sup>43</sup>. An ER diagram (not shown) captures relationships between users and their scan results.

### 3.2 System Architecture

The platform follows a multi-tier architecture. The **frontend** is built in Streamlit, a Python framework for interactive apps. It provides pages for user registration/login, a dashboard, file scanning, and URL scanning <sup>34</sup> <sup>8</sup>. For example, the dashboard page dynamically displays recent scans, threat distribution charts, and personalized security insights. Interactive visualizations (using Plotly) show trends such as count of detected threats over time or breakdown of malicious vs. clean scans <sup>35</sup> <sup>44</sup>.

The **backend** is implemented in Python and handles business logic. Core modules include: - **Authentication (utils/auth\_utils.py):** Manages user sessions and access control. User passwords are hashed using bcrypt before storing <sup>45</sup>, and Streamlit's `st.session_state` is used to track login state and enforce page-level authentication. - **Database Utilities (utils/db\_utils.py):** Interfaces with a SQLite database (`app.db`) to store users and scan records. On each new scan request, the backend creates a pending record linked to the user. - **VirusTotal API Integration (utils/virustotal\_api.py):** Encapsulates all communication with VirusTotal. Functions like `upload_file_to_virustotal` and `scan_url_with_virustotal` send data to the appropriate VT API endpoints <sup>46</sup> <sup>47</sup>. For example, uploading a file is done via a `requests.post` to the `/files` endpoint with the file as multipart data <sup>47</sup>.

Because VT scans are asynchronous, the backend polls the `/analyses/{id}` endpoint until the analysis is complete, then retrieves the report. - **Machine Learning (`utils/ml_utils.py`)**: Processes VirusTotal results to derive deeper insights. It computes a **risk score** (0–100) based on the number of malicious/suspicious detections <sup>32</sup>. It also performs **file anomaly detection** by examining file metadata (e.g. double extensions, unusually small executables) that often indicate malicious intent <sup>33</sup>. Historical scan data can be analyzed to detect **threat patterns** (common attack vectors) and generate **security recommendations** for users <sup>48</sup>. - **Visualization (`utils/visualization_utils.py`)**: Prepares data for charts. It creates Plotly figures (bar charts of threat types, timelines of scans, risk gauge dials, etc.) that the frontend embeds directly into the Streamlit pages <sup>35</sup>.

All data (user accounts, scan history, ML results) is stored in a local **SQLite** database for persistence <sup>49</sup>. The schema is simple: a `users` table (id, username, email, password hash, role) and a `scans` table (id, user\_id, scan\_type, target, file\_hash, analysis\_id, status, malicious\_count, suspicious\_count, clean\_count, result\_summary, timestamp). This enables quick data retrieval for the user dashboard and ensures user data privacy.

Overall, the architecture enables a smooth data flow: when a user submits a file/URL, the backend creates a scan record, sends it to VirusTotal, waits for the report, then runs the ML analysis and updates the record <sup>20</sup> <sup>39</sup>. Finally, all results (raw VT scan data plus ML-derived risk score and recommendations) are sent back to the frontend to display to the user. Meanwhile, the dashboard aggregates historical data from the database to show trends. This layered design (Streamlit UI → Python backend → SQLite storage with VT/ML integration) provides scalability and clear separation of concerns <sup>50</sup> <sup>51</sup>.

## Chapter Four: Implementation and Results

### 4.1 Development Environment and Tools

The system was implemented entirely in **Python**. Key dependencies (listed in `requirements.txt`) include **streamlit** (for the web UI), **requests** (for API calls), **python-dotenv** (for managing environment variables like `VIRUSTOTAL_API_KEY`), and **bcrypt** (for password hashing) <sup>52</sup>. The Python 3.11 environment and these libraries are sufficient to run the application on a standard development machine. Streamlit abstracts away most front-end work, allowing rapid creation of forms and charts with simple Python code.

For example, setting up a file upload widget in Streamlit required only a few lines of code. When a user uploads a file, the backend reads it in binary mode and calls a function like:

```
analysis = upload_file_to_virustotal(uploaded_file.read())
```

which posts to VirusTotal's `/files` endpoint (see code snippet) <sup>47</sup>. The app then enters a polling loop (`wait_for_analysis_completion`) until the VT analysis is done. Similar code handles URL scans. All VT API communications include the secret API key in the header (secured server-side via an environment variable, not exposed to the client).

## 4.2 User Interface

The front-end is delivered as a single Streamlit app ( `app.py` ) with multiple pages. The **Login/Register** page uses the `auth_utils.py` module: passwords entered by users are immediately hashed with bcrypt<sup>45</sup> and stored in SQLite. A `require_auth` decorator (using `st.session_state`) guards the main pages so only authenticated users can access scanning functions<sup>53</sup>.

After logging in, the **Dashboard** page displays personalized security metrics: recent scan history, charts of threat distribution, and risk gauges. For instance, we implemented a horizontal bar chart (via Plotly) showing counts of detected malicious and suspicious results per user. The interface is clean and responsive, making advanced features accessible to non-technical users<sup>34 44</sup>.

The **Scan File** page lets users upload a file for analysis. Upon submission, the file is immediately sent to VirusTotal (with a “pending” status shown to the user). Once VirusTotal finishes, the app presents the results: it shows the number of antivirus engines that flagged the file, a table of engine results, and the ML **risk score** (e.g. “Risk: 85/100” if many engines found malware). The **Scan URL** page works similarly for URLs, checking domain and history. All scan results (file hash, URL, malicious count, etc.) are stored in the database so users can review past results. In summary, the Streamlit UI provides an intuitive workflow for uploading data, seeing progress, and viewing detailed results (see screenshots in Appendix, not shown here). Feedback from testers confirmed the interface is simple and actionable<sup>54</sup>.

## 4.3 VirusTotal Integration

Integration with VirusTotal is a core feature. The `virustotal_api.py` module contains functions like `upload_file_to_virustotal(file_content)` and `scan_url_with_virustotal(url)`, which wrap the HTTP calls. For example, uploading a file uses code:

```
url = f"{VIRUSTOTAL_BASE_URL}/files"
headers = {"X-Apikey": VIRUSTOTAL_API_KEY}
files = {"file": file_content}
response = requests.post(url, headers=headers, files=files)
```

This routine handles errors and returns the JSON response containing an analysis ID<sup>47</sup>. For URL scans, a similar approach posts to `/urls`. Because VT scans are not instantaneous, the system then polls the `/analyses/{id}` endpoint every few seconds, checking the `status` field until it reports “completed”. This asynchronous handling is crucial: it allows the UI to remain responsive (Streamlit shows a loading spinner) while waiting for the scan to finish. Once the report is ready, the raw JSON (with engine results) is fetched and parsed.

After retrieving the VT report, the backend extracts key metrics (malicious count, suspicious count, file hash/URL, etc.) and updates the `scans` table. We also pass the report to the ML module for further processing.

## 4.4 Machine Learning Components

The platform's ML logic is implemented in `ml_utils.py`. When a VirusTotal report is obtained, the system performs several analyses:

- **Risk Scoring:** We designed a simple proprietary algorithm that assigns a numeric risk score (0–100) based on the proportion of engines flagging the target as malicious or suspicious <sup>32</sup>. For example, if 40 out of 70 engines detect malware, the score might be  $(40/70)*100 \approx 57$ . This gives a quick, quantitative measure of threat severity.
- **File Anomaly Detection:** The system examines file metadata (extension, size, filename patterns) in combination with the scan results <sup>33</sup>. For instance, if an executable file has a double extension ("something.pdf.exe") or is suspiciously small yet flagged malicious, the module raises an anomaly flag. This helps catch clever obfuscations or new malware that might otherwise appear "clean" to some engines.
- **Threat Pattern Analysis:** The tool aggregates a user's historical scan data to detect trends: for instance, it can report if 80% of scanned URLs in the past week were phishing links. This is shown as an insight on the dashboard.
- **Security Recommendations:** Based on the risk score and patterns, the system generates text advice (e.g. "This file is high-risk; avoid execution or consult an IT administrator."). These contextual recommendations are displayed alongside scan results.

Unit tests confirmed these ML routines work as intended: functions like `calculate_risk_score`, `detect_file_anomalies`, and `generate_security_recommendations` were tested on sample inputs to ensure meaningful outputs <sup>55</sup>.

## 4.5 Data Flow

The full data flow proceeds as follows:

1. **User Action:** The user logs in and navigates to "Scan File" or "Scan URL" (Streamlit pages).
2. **Input Submission:** The user selects a file or enters a URL, then clicks "Scan." Streamlit captures this input.
3. **Record Creation:** The backend creates a new "pending" scan record in the SQLite database (associating it with the user).
4. **VirusTotal Submission:** The file content or URL is sent to VirusTotal via the API functions <sup>20</sup>.
5. **Polling:** The backend periodically polls VirusTotal for the analysis status.
6. **Report Retrieval:** Once complete, the raw VirusTotal analysis report is received.
7. **ML Analysis:** The `ml_utils.py` processes the report: risk score is calculated and anomalies checked <sup>32</sup> <sup>33</sup>.
8. **Database Update:** The scan record is updated to "completed," and fields like `malicious_count`, `suspicious_count`, `risk_score`, and `result_summary` are filled.
9. **Display Results:** The processed results (including charts and recommendations) are sent to the Streamlit UI and shown to the user.
10. **Dashboard Update:** The user's dashboard is refreshed to include the new scan data (threat distribution charts, updated statistics) <sup>56</sup>.

This pipeline ensures all steps are tracked and the user sees up-to-date information. In testing, the data flow was confirmed end-to-end: each scan triggered a database entry, VT API interaction, ML processing, and UI update without errors <sup>57</sup> <sup>58</sup> .

## 4.6 Performance Evaluation and Challenges

The system was evaluated through unit, integration, and user testing <sup>59</sup> <sup>60</sup> . **Unit tests** validated each module (e.g. authentication hashing, API calls, ML functions) in isolation. **Integration tests** verified that file/URL scanning from the UI all the way through ML and back to the dashboard works correctly <sup>57</sup> . Error handling was tested by simulating API failures and invalid inputs to ensure graceful messaging.

In terms of **performance**, the main factor affecting scan time is the VirusTotal API response. On average, small files returned results in a few seconds, but large files or cold scans (file never seen by VT) could take up to 2–3 minutes due to queuing <sup>31</sup> . Our polling mechanism handles this variability smoothly (showing a progress indicator), but users must wait longer for big files. Dashboard loading (retrieving and plotting historical data) was fast for typical histories (hundreds of records) because SQLite queries and Plotly renderings completed in under a second <sup>61</sup> . CPU and memory usage of the Streamlit app remained modest, making it feasible to run on a standard desktop or small server <sup>62</sup> .

**Challenges Encountered:** The asynchronous nature of VirusTotal scans required careful design: we needed to manage API rate limits and avoid blocking the UI thread. We addressed this by using non-blocking loops and spacing out requests. Rate-limiting itself is a limitation: with the free API, only a limited number of scans per minute were allowed <sup>27</sup> . In high-usage scenarios, this could delay scans. Additionally, since the system relies on VirusTotal data, zero-day threats not yet in VT's database might be underreported; our ML anomaly checks help mitigate this. Finally, ensuring usability across different devices was challenging, but Streamlit's responsiveness and simple UI design helped; informal user tests indicated the interface was intuitive, though users suggested adding help tips for technical terms <sup>63</sup> .

## Chapter Five: Summary, Conclusions, and Recommendations

This project has successfully designed and implemented a machine learning-powered network security platform with integrated VirusTotal threat intelligence. Using **Streamlit** and Python, we built a clean, interactive web interface that lets users submit files or URLs for analysis <sup>4</sup> <sup>34</sup> . Each submission is scanned by VirusTotal's multi-engine database, and the results are enhanced by ML algorithms that compute a **risk score** and detect anomalies <sup>32</sup> <sup>33</sup> . The platform stores scan histories and displays them in a personalized dashboard with charts and insights. Key achievements include a user-friendly interface for complex analysis, robust use of VirusTotal for comprehensive scanning, and intelligent risk assessment beyond static signatures <sup>64</sup> <sup>65</sup> . Testing confirmed the system's functionality and stability; for example, end-to-end scans reliably produced accurate outputs and the application handled expected error scenarios gracefully <sup>60</sup> . Overall, the platform fulfills its goal of providing an accessible cybersecurity tool that helps users in under-resourced areas detect and understand threats in real time.

**Recommendations and Future Work:** Although the core objectives are met, several enhancements could strengthen the platform. Notably, more advanced ML models (deep learning or predictive analytics) could improve detection of sophisticated or future threats <sup>66</sup> . Integrating additional threat intelligence sources (such as Passive DNS, sandbox analysis, or commercial feeds) would broaden coverage beyond what VirusTotal provides <sup>67</sup> . Improving scalability is also advised: for large deployments, migrating from SQLite

to a scalable database (e.g. PostgreSQL) and employing asynchronous task queues (e.g. Celery) would speed up concurrent scans <sup>68</sup> . User experience can be enhanced by making the UI mobile-responsive, adding tooltips or guided tutorials, and supporting multiple languages <sup>69</sup> . Finally, for real-world deployment, containerization (Docker/Kubernetes) and security hardening (HTTPS, firewall rules) should be implemented. These future improvements would ensure the system remains effective and adaptable in the dynamic cybersecurity landscape.

In conclusion, this project demonstrates that combining open-source development (Python, Streamlit) with widely-available threat intelligence (VirusTotal) and ML techniques can yield a powerful, scalable security platform. By focusing on usability and API-driven analysis, it provides a practical tool for proactive threat detection, particularly suited to settings with limited resources <sup>70</sup> <sup>71</sup> . The solution and lessons learned here can guide future work to extend and refine intelligent cybersecurity tools for broader impact.

**References:** (Selected from literature and project documentation) Sommer & Paxson (2010); Al-Janabi & Saeed (2016); Papernot et al. (2016); Streamlit (n.d.); VirusTotal (n.d.) <sup>72</sup> <sup>6</sup> <sup>27</sup> <sup>4</sup> <sup>64</sup> .

1 2 5 6 7 10 11 12 13 14 15 16 17 18 19 21 22 23 24 25 26 27 28 36 37 38 40 41 42 43

51 72 Design & Implementation of Machine Learning UNIMAID Baffun.docx

file:///file-LcnyVgKwLUWzHSLm9C1W8y

3 4 8 9 20 29 30 31 32 33 34 35 39 44 45 46 47 48 49 50 52 53 54 55 56 57 58 59 60 61

62 63 64 65 66 67 68 69 70 71 Network\_Security\_Platform\_A\_Streamlit-Based\_AI-

Powered\_Threat\_Detection\_System.pdf

file:///file-6SmqiZ6vksX49KgBaDpkCD