# TECHNICAL REPORT

## CYBER-SHIELD INTELLIGENT ASSISTANT

> **Project Type:** Hybrid Cybersecurity Assistant (LLM + Deep Learning)
> **Competition:** Furssah AI Competition - Tunisia
> **Repository:** GitHub Repository

**Authors:** Muhammed Ehab & Ghassen Sellami

May 30, 2025

# Contents

**Abstract**

This project presents an innovative hybrid assistant designed for cybersecurity tasks, combining the power of fine-tuned large language models (LLMs) with deep learning-based threat detection systems. The assistant enables contextual understanding, real-time anomaly detection, and actionable threat analysis in modern networked environments. The system integrates Qwen-2.5B Instruct for natural language processing with machine learning models including Decision Trees, Random Forest, SVM, and LSTM networks for comprehensive threat classification.

# 1   Introduction

In the rapidly evolving landscape of cybersecurity, traditional rule-based systems are increasingly inadequate for addressing sophisticated threats. This project introduces CYBER-SHIELD, a hybrid intelligent assistant that leverages both large language models and deep learning techniques to provide comprehensive cybersecurity analysis and support.

The system combines:

- **Contextual AI Chat:** Fine-tuned Qwen-2.5B model for expert cybersecurity guidance

- **Threat Detection:** ML/DL pipeline for real-time network anomaly detection

- **Web Interface:** User-friendly dashboard for interaction and visualization

- **Modular Architecture:** Scalable design for easy extension and deployment

# 2   Objectives

The primary objectives of this project are:

1. **Develop a contextual chat interface** for cybersecurity using Qwen-2.5B Instruct model

2. **Integrate ML/DL models** (Decision Tree, SVM, Random Forest, RNN/LSTM) for comprehensive threat classification

3. **Design a modular system** with intuitive web UI using Flask backend and HTML/CSS/JS frontend

4. **Enable log ingestion** and produce actionable risk scoring with detailed recommendations

# 3   System Architecture

## 3.1   System Modules

The CYBER-SHIELD system consists of four main modules:

> **System Modules**
>
> - **LLM Module:** Fine-tuned Qwen-2.5B with LLaMA Factory, Hugging Face, and Weights & Biases integration
>
> - **Threat Detection Module:** Hybrid ML + DL models trained on labeled network flow datasets
>
> - **Frontend Interface:** Flask backend with custom HTML/JavaScript dashboard
>
> - **API Layer:** RESTful endpoints for LLM queries and flow data submission

## 3.2   System Workflow

The system operates through the following workflow:

1. User submits either a chat query or threat log data via the web interface

2. Backend routing system directs input to appropriate module (LLM or Classifier)

3. Processing occurs using specialized models for the input type

4. Results are formatted and returned as either conversational responses or threat insights

5. Output is displayed through the web dashboard with visualizations and recommendations

```
[Architecture Diagram Placeholder]
System Architecture Overview
```

Figure 1: CYBER-SHIELD System Architecture

# 4   LLM Component

## 4.1   Model Specification

The language model component utilizes:

- **Base Model:** Qwen-2.5B Instruct

- **Fine-Tuning Framework:** LLaMA Factory with Hugging Face Trainer

- **Monitoring:** Weights & Biases for experiment tracking

- **Deployment:** Local inference with API integration capability

## 4.2   Use Cases

The LLM component provides expert assistance in:

> **LLM Capabilities**
>
> - **CVE Analysis:** Detailed explanation of Common Vulnerabilities and Exposures
> - **Penetration Testing:** Step-by-step guidance on ethical hacking techniques
> - **Malware Analysis:** Behavioral insights and detection strategies
> - **Incident Response:** SOC workflow guidance and IR procedures
> - **Security Best Practices:** Recommendations for system hardening

# 5   Threat Classification Component

## 5.1   Input Processing

The threat detection system processes various data formats:

- Network flow logs (NetFlow format)
- Zeek/Bro JSON outputs
- Custom log formats with standardized schema
- Real-time stream data

## 5.2   Machine Learning Models

The system employs a hybrid approach combining traditional ML and deep learning:

| Model Type | Algorithm | Use Case |
|---|---|---|
| Traditional ML | Decision Tree | Fast classification with interpretability |
| | Random Forest | Ensemble method for improved accuracy |
| | SVM | High-dimensional data classification |
| | KNN | Similarity-based threat detection |
| Deep Learning | RNN/LSTM | Sequential pattern recognition |
| | Neural Networks | Complex feature learning |

Table 1: Machine Learning Models Overview

## 5.3   Output Classification

The system generates comprehensive threat assessments including:

- **Threat Type:** Classification (e.g., DDoS, Port Scan, Malware)

- **Risk Score:** Numerical assessment (0-100 scale) with confidence interval

- **Anomaly Detection:** Identification of unusual patterns

- **Actionable Recommendations:** Specific mitigation strategies

# 6 Example Output

The following JSON structure demonstrates the system's analytical output:

```json
{
  "type": "suspicious",
  "origin": "laptop",
  "protocol": "TCP",
  "risk_score": "75/100",
  "confidence": "85%",
  "threats": ["port scan", "potential DDoS"],
  "anomalies": [
    "Abnormally high packet rate (500K pps)",
    "Unusual unidirectional flow",
    "Non-standard destination port"
  ],
  "recommendation": "Investigate further: Block source IP temporarily
   and analyze traffic patterns for 24 hours"
}
```

Listing 1: Threat Analysis Output Example

# 7 Web Interface

## 7.1 Backend Architecture

The backend utilizes Flask framework with the following components:

- **API Routes:** RESTful endpoints for data processing

- **Authentication:** Session management and user validation

- **Data Processing:** Real-time log parsing and model inference

- **Response Formatting:** JSON serialization for frontend consumption

## 7.2 Frontend Implementation

The user interface provides:

> **Frontend Features**
>
> - **File Upload:** Drag-and-drop JSON log submission
>
> - **Real-time Chat:** Interactive LLM conversation interface
>
> - **Threat Dashboard:** Visualization of classification results
>
> - **Historical Analysis:** Previous scan results and trends
>
> - **Export Functionality:** Report generation in multiple formats

# 8 Technologies Used

## 8.1 Core Dependencies

| Category | Technologies |
| --- | --- |
| Web Framework | Flask, FastAPI |
| LLM Processing | Transformers, LLaMA Factory, Hugging Face |
| Machine Learning | Scikit-learn, TensorFlow, Keras |
| Data Processing | Pandas, NumPy, Matplotlib |
| Monitoring | Weights & Biases |
| Utilities | dotenv, aiohttp, uvicorn |

Table 2: Technology Stack

# 9 Project Structure

The project follows a modular architecture with clear separation of concerns:

```
cybershield-intelligent/
        app/
                api/
                        chat_handler.py          # LLM chat routes
                        threat_analyzer.py       # Threat analysis endpoints
                        main.py                  # Application entry point
                llm/
                        qwen_chat.py             # Qwen model inference
    wrapper
                models/
                        predictor_rnn.py         # LSTM inference logic
                        preprocess.py            # Data normalization
                        model_weights/           # Trained model artifacts
                frontend/
                        streamlit_ui.py          # User interface components
                utils/
                    helpers.py                   # Utility functions
        data/
                samples/                         # Sample datasets
        assets/
                images/                          # Documentation assets
```

```
21              demo.gif                    # Demonstration materials
22          requirements.txt                # Python dependencies
23          README.md                       # Project documentation
24          LICENSE                         # MIT License
```

Listing 2: Project Directory Structure

# 10    Implementation Details

## 10.1    Model Training Process

The threat detection models were trained using:

1. **Data Collection:** Curated datasets from network security repositories

2. **Preprocessing:** Feature engineering and normalization pipelines

3. **Model Selection:** Cross-validation for optimal algorithm selection

4. **Hyperparameter Tuning:** Grid search and Bayesian optimization

5. **Validation:** K-fold cross-validation with stratified sampling

## 10.2    Performance Metrics

The system achieves the following performance benchmarks:

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Random Forest | 94.2% | 93.8% | 94.1% | 93.9% |
| SVM | 91.7% | 90.5% | 92.3% | 91.4% |
| LSTM | 96.1% | 95.7% | 96.3% | 96.0% |

Table 3: Model Performance Comparison

# 11    Conclusion

This project successfully demonstrates a comprehensive approach to cybersecurity assistance through the integration of large language models and machine learning techniques. The CYBER-SHIELD system provides:

- **Intelligent Conversation:** Expert-level cybersecurity guidance through fine-tuned LLM

- **Accurate Threat Detection:** High-performance classification with 96% accuracy

- **User-Friendly Interface:** Intuitive web dashboard for seamless interaction

- **Modular Design:** Extensible architecture for future enhancements

The system is production-ready and suitable for deployment in educational environments, small to medium enterprises, and as a foundation for more complex cybersecurity platforms.

# 12  Future Work

Several enhancements are planned for future iterations:

## 12.1  Data Processing Improvements

- **PCAP Support:** Direct packet capture file analysis

- **Suricata Integration:** Real-time IDS log processing

- **Stream Processing:** Apache Kafka/Flink integration for live data

## 12.2  Interface Enhancements

- **React Frontend:** Modern, responsive user interface

- **Mobile Application:** iOS/Android companion apps

- **API Documentation:** Comprehensive Swagger/OpenAPI specs

## 12.3  Security Features

- **Authentication System:** OAuth2/SAML integration

- **Role-Based Access:** Multi-tenant user management

- **Audit Logging:** Comprehensive activity tracking

# 13  Acknowledgments

This project was developed for the Furssah AI Competition in Tunisia. Special thanks to the organizers for providing an excellent platform for AI innovation in cybersecurity.

The project utilizes open-source technologies and datasets from the cybersecurity community, contributing to the collective advancement of AI-driven security solutions.

# 14  License

This project is released under the MIT License, promoting open-source collaboration and knowledge sharing in the cybersecurity community.

# 15  References

# References

[1] Qwen Team. (2024). *Qwen2.5: A Large Language Model for Diverse Applications.* Alibaba Cloud.

[2] hiyouga. (2024). *LLaMA-Factory: Unified Efficient Fine-Tuning of 100+ LLMs.* GitHub Repository.

[3] Wolf, T., et al. (2020). *Transformers: State-of-the-Art Natural Language Processing.* EMNLP 2020.

[4] Pedregosa, F., et al. (2011). *Scikit-learn: Machine Learning in Python.* JMLR 12, pp. 2825-2830.

[5] Abadi, M., et al. (2016). *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems.* OSDI.

# Author Information

**Muhammed Ehab**

**GitHub:** github.com/muhammedehab35
**Project Repository:** FURSSAH_AI_COMETITION
**Email:** Contact via GitHub profile

*Developed for Furssah AI Competition - Tunisia 2024*