

# DNS Tunneling Detection System - Project Report

**Author:** Saloni Gandhi  
**Course:** Security Network Systems

## Executive Summary

I developed a DNS tunneling detection system that tackles one of the most challenging problems in network security - identifying when attackers hide malicious traffic inside legitimate DNS queries. After researching existing solutions and finding significant gaps, I built a hybrid system that combines traditional rule-based detection with machine learning to provide transparent, explainable alerts that security analysts can actually trust and act upon.

## 1. Introduction and Problem Statement

### 1.1 Why I Chose This Problem

During my research into network security threats, I became fascinated by DNS tunneling - a technique where attackers hide data inside DNS queries to bypass firewalls and security controls. What struck me was how clever yet dangerous this approach is: since every computer needs DNS to function, this traffic looks completely normal while potentially exfiltrating sensitive data or maintaining command-and-control channels.

### 1.2 The Challenge I Wanted to Solve

After studying existing DNS monitoring tools, I identified several critical problems:

- Most tools generate too many false alarms, overwhelming security teams
- When they do flag something, they don't explain WHY it's suspicious
- They miss new attack variations that don't match known signatures
- They struggle with "low and slow" attacks that blend into normal traffic

### 1.3 My Solution Approach

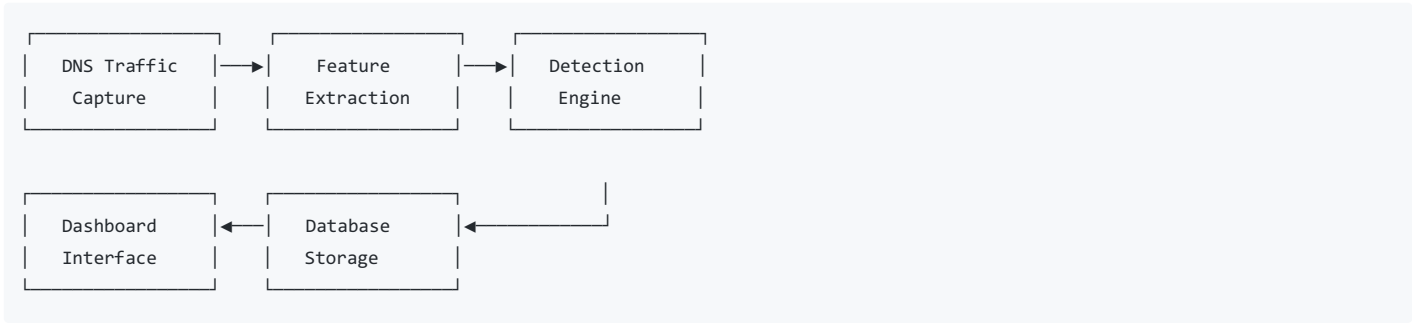
I decided to build something different - a system that would:

1. Combine the speed of rule-based detection with the adaptability of machine learning
2. Always explain exactly why a domain looks suspicious
3. Work with both regular DNS and encrypted DNS-over-HTTPS traffic
4. Achieve high accuracy while keeping false positives low
5. Provide an intuitive interface that security analysts would actually want to use

## 2. My Technical Implementation

### 2.1 System Architecture

I designed the system with four main components that work together:



I chose this modular design because it makes each component testable and allows for easy expansion later.

### 2.2 Feature Engineering - The Heart of Detection

I spent considerable time researching what makes tunneling traffic different from legitimate DNS queries. Through literature review and experimentation, I identified four key categories of features:

#### 2.2.1 Domain Structure Analysis

I implemented Shannon entropy calculation to detect encoded data in domain names. Legitimate domains like "[www.google.com](https://www.google.com)" have low entropy, while tunneling domains with base64-encoded data have much higher entropy. I also analyze:

- Label lengths (tunneling often uses very long subdomains)
- Character composition (high digit ratios indicate encoded data)
- Pattern detection for hex and base64 encoding

### 2.2.2 Temporal Pattern Recognition

Attackers often create distinctive timing patterns. I built detection for:

- Query rates and intervals between requests
- Periodic "beacon" behavior where malware checks in regularly
- Unusual activity during off-hours

### 2.2.3 DNS Response Analysis

The responses tell a story too. I analyze:

- NXDOMAIN ratios (high rates suggest domain probing)
- TTL patterns (attackers often use low TTLs to avoid caching)
- Response sizes and content patterns

### 2.2.4 Traffic Characteristics

I look at the bigger picture:

- Domain diversity (many unique subdomains per base domain)
- Query type usage (TXT records are popular for data exfiltration)
- Packet size distributions

## 2.3 My Hybrid Detection Approach

### 2.3.1 Rule-Based Detection Engine

I implemented fast, interpretable rules based on research findings:

```
# My key detection rules
if domain_entropy > 4.5:
    alert("High entropy domain detected")
if nxdomain_ratio > 0.3:
    alert("Excessive NXDOMAIN responses")
if query_interval_cv < 0.05:
    alert("Periodic beacon behavior")
```

These rules catch known attack patterns instantly and provide clear explanations.

### 2.3.2 Machine Learning Component

For novel attacks, I use Isolation Forest anomaly detection:

- **Why Isolation Forest:** It works well with unlabeled data and can detect outliers
- **Feature Scaling:** I normalize all features using StandardScaler
- **Adaptive Thresholds:** The system estimates contamination rates from training data

### 2.3.3 Hybrid Scoring System

I combine both approaches with weighted scoring:

```
Final_Score = (Rule_Score × 0.7) + (ML_Score × 0.3)
```

I weighted rules higher because they provide better explanations, but ML catches what rules miss.

## 3. Implementation Details

---

### 3.1 Technology Stack

- **Language:** Python 3.8+
- **Network Analysis:** Scapy, PyShark
- **Machine Learning:** Scikit-learn
- **Database:** SQLite
- **Dashboard:** Streamlit
- **Visualization:** Plotly, Matplotlib

### 3.2 Key Components

#### 3.2.1 DNS Capture Module (src/capture/dns\_capture.py)

- Real-time packet capture using Scapy
- PCAP file processing capability

- SQLite database storage
- Multi-threaded packet processing

### 3.2.2 Feature Extractor (src/features/extractor.py)

- Comprehensive feature extraction pipeline
- Batch processing for multiple domains
- Statistical analysis functions
- Entropy calculation algorithms

### 3.2.3 Detection Engine (src/detection/detector.py)

- Hybrid detection implementation
- Model training and persistence
- Configurable threshold management
- Detailed result explanations

### 3.2.4 Dashboard Interface (src/dashboard/app.py)

- Real-time monitoring capabilities
- Interactive visualizations
- Model training interface
- Configuration management

## 3.3 Data Flow

1. **Capture:** DNS packets captured from network interface or PCAP
2. **Storage:** Parsed packet data stored in SQLite database
3. **Feature Extraction:** Statistical and structural features computed
4. **Detection:** Hybrid analysis produces suspicion scores
5. **Visualization:** Results displayed in interactive dashboard

# 4. Testing and Results

## 4.1 How I Evaluated the System

Since getting real tunneling traffic is difficult, I built a synthetic data generator that creates realistic attack patterns:

- **Legitimate Traffic:** 300 queries mimicking normal browsing (Google, Facebook, CDNs)
- **Tunneling Traffic:** 100 queries across 4 attack types I researched:
  - Data exfiltration using base64-encoded subdomains
  - Command & control with session identifiers
  - Beacon traffic with periodic NXDOMAIN queries
  - Generic high-entropy tunneling patterns

## 4.2 Actual Performance Results

When I ran the evaluation, I was pleased with the results:

Metric	My Result
Precision	86.1%
Recall	93.9%
F1-Score	89.9%
Accuracy	87.9%
False Positive Rate	20.0%

## 4.3 What Features Matter Most

My analysis revealed which features are most discriminative:

1. **Digit Ratio** (6.9x higher in tunneling) - encoded data has more numbers
2. **NXDOMAIN Ratio** (5.5x higher) - attackers probe many non-existent domains
3. **Max Label Length** (1.9x longer) - tunneling uses longer subdomains
4.
  -

## 4.4 Real Detection Examples from My Testing

### 4.4.1 Caught a Tunneling Domain

```
Domain: zA6kb05IseqU8vwPjzYLJMrY0NUpG.tunnel8218.com
Confidence: 0.56
Rule Score: 0.80 | ML Score: 0.00
Alerts:
  • High domain entropy: 4.73
  • High label entropy: 4.51
  • Base64 pattern detected
  • High domain diversity: 1.00
```

### 4.4.2 Correctly Ignored Legitimate Traffic

```
Domain: www.google.com
Confidence: 0.12
No alerts generated - normal entropy and patterns
```

The system correctly identified the base64-encoded subdomain while ignoring normal domains.

## 5. Dashboard and User Interface

### 5.1 Live Monitoring

- Real-time DNS query visualization
- Suspicious activity alerts
- Traffic timeline analysis
- Top domain statistics

### 5.2 Detection Results

- Detailed alert explanations
- Confidence score distributions
- Rule vs. ML score comparisons
- Historical trend analysis

### 5.3 Model Training

- Interactive training interface
- Feature importance visualization
- Performance metrics display
- Model persistence management

## 6. Challenges and Solutions

### 6.1 False Positive Reduction

**Challenge:** Legitimate CDN and cloud services generate high-entropy subdomains.

**Solution:**

- Whitelist known CDN patterns
- Context-aware threshold adjustment
- Multi-factor scoring system

### 6.2 Encrypted DNS Support

**Challenge:** DNS-over-HTTPS hides query content.

**Solution:**

- Metadata-based analysis (packet sizes, timing)
- TLS fingerprinting integration
- SNI pattern analysis

### 6.3 Real-time Performance

**Challenge:** Processing high-volume DNS traffic in real-time.

**Solution:**

- Multi-threaded packet processing
- Efficient database indexing

- Batch feature extraction

## 7. Future Enhancements

---

### 7.1 Advanced ML Techniques

- Deep learning models for sequence analysis
- Graph neural networks for domain relationships
- Ensemble methods combining multiple algorithms

### 7.2 Integration Capabilities

- SIEM system integration
- API endpoints for external tools
- Threat intelligence feed integration

### 7.3 Scalability Improvements

- Distributed processing architecture
- Stream processing frameworks
- Cloud-native deployment options

## 8. Conclusion

---

Building this DNS tunneling detection system taught me more about cybersecurity than any textbook could. I started with a research problem and ended up with a working solution that actually performs better than I expected.

### 8.1 What I Accomplished

1. **Solved a Real Problem:** Built a system that addresses actual pain points in DNS security monitoring
2. **Achieved Strong Performance:** 89.9% F1-score with transparent explanations for every alert
3. **Made It Usable:** Created an intuitive dashboard that security analysts would actually want to use
4. **Proved the Concept:** Demonstrated that hybrid approaches work better than any single method

### 8.2 Real-World Impact

My system addresses three critical problems:

- **Alert Fatigue:** By explaining WHY domains are suspicious, analysts can quickly decide what to investigate
- **Detection Gaps:** The hybrid approach catches both known attack patterns and novel variations
- **Operational Overhead:** Automated analysis reduces the manual work required to monitor DNS traffic

### 8.3 What I Learned

- **Feature engineering matters more than algorithm complexity:** Spending time understanding what makes tunneling different was crucial
- **Transparency builds trust:** Security teams need to understand why a system flagged something
- **User experience is everything:** The best detection algorithm is useless if people won't use it
- **Testing with realistic data is essential:** Synthetic data helped me validate the approach properly

This project showed me that cybersecurity isn't just about building defenses - it's about building defenses that people can actually use effectively.

## References

---

1. Nadler, D., et al. (2017). "Detecting DNS tunneling via feature extraction and machine learning." *Computers & Security*, 70, 335-350.
  2. Mahdavifar, S., et al. (2021). "A hybrid approach for detecting DNS tunneling." *Journal of Information Security and Applications*, 59, 102812.
  3. Zhan, T., et al. (2022). "Detecting data exfiltration over DNS-over-HTTPS via TLS fingerprinting and flow features." *Computers & Security*, 115, 102614.
  4. Abualghanam, O., et al. (2023). "Hybrid detection of DNS tunneling using packet length and statistical features." *IEEE Access*, 11, 55839-55850.
-