

FINAL PROJECT REPORT - PART 2

Continued from FINAL_REPORT_PART1.md

CHAPTER 3: SYSTEM ARCHITECTURE & DESIGN

3.1 Overall System Architecture

The proposed DDoS mitigation system employs a hybrid two-layer architecture that combines kernel-level packet processing with user-space intelligent analysis. This design achieves both high throughput and high accuracy by leveraging the strengths of each layer while mitigating their individual weaknesses.

3.1.1 Architectural Overview

The system consists of three main planes:

1. Data Plane (Kernel Space - eBPF/XDP):

- Processes every incoming packet at NIC driver level
- Performs immediate blacklist enforcement
- Collects per-packet, per-IP, and per-flow statistics
- Achieves sub-microsecond processing latency
- Operates at line rate (5M+ packets per second)

2. Control Plane (User Space - Python):

- Reads aggregated statistics from eBPF maps (1-second intervals)
- Performs statistical anomaly detection
- Executes ML classification
- Makes mitigation decisions
- Updates kernel blacklist dynamically

3. Management Plane (Web Dashboard):

- Provides real-time visibility
- Displays metrics, alerts, and system status
- Enables operator configuration
- Shows historical data and trends

[**Figure 3.1: Overall System Architecture** - See system_architecture_diagram.png]

3.1.2 Layer Responsibilities

Kernel Layer Responsibilities:

1. **Packet Classification:** Parse Ethernet, IP, TCP/UDP headers
2. **Blacklist Enforcement:** Immediate drop of blacklisted IPs (XDP_DROP)

3. **Statistics Collection:** Update per-CPU, per-IP, and per-flow counters
4. **Simple Heuristics:** Basic SYN flood detection (SYN count > threshold)
5. **Fast Path Processing:** <1 microsecond per packet

User Space Layer Responsibilities:

1. **Baseline Learning:** Establish and adapt normal traffic profiles
2. **Statistical Analysis:** Calculate PPS, BPS, entropy, protocol ratios
3. **Feature Extraction:** Compute 64 CIC features from aggregated stats
4. **ML Classification:** Random Forest prediction on extracted features
5. **Hybrid Decision:** Combine statistical and ML scores
6. **Policy Enforcement:** Add/remove IPs from kernel blacklist

Design Rationale:

- Kernel handles per-packet speed requirements
- User space handles complex intelligence
- Separation of concerns enables independent optimization
- eBPF maps provide efficient communication bridge

3.1.3 Data Flow

Normal Packet Flow:

1. Packet arrives at NIC
2. XDP hook triggers eBPF program
3. Parse headers
4. Check blacklist (not found)
5. Update statistics maps
6. Return XDP_PASS
7. Packet continues to network stack

Blacklisted Packet Flow:

1. Packet arrives at NIC
2. XDP hook triggers eBPF program
3. Parse headers
4. Check blacklist (found!)
5. Return XDP_DROP immediately
6. No statistics update, minimal processing

Control Flow (per second):

1. User space reads eBPF map statistics
2. Aggregate per-CPU stats
3. Calculate traffic metrics (PPS, BPS, etc.)
4. Update baseline profile

5. Perform statistical anomaly checks
6. Extract ML features
7. Run Random Forest classifier
8. Combine statistical + ML scores
9. If attack: Add source IPs to blacklist
10. Update dashboard metrics

3.2 Traffic Shaping and Packet Flow Design

3.2.1 Packet Processing Pipeline

The packet processing pipeline is optimized for minimal latency while collecting comprehensive statistics.

Step 1: Ethernet Frame Parsing

```
struct ethhdr *eth = data;
if ((void *)(eth + 1) > data_end)
    return XDP_DROP; // Malformed packet
if (eth->h_proto != htons(ETH_P_IP))
    return XDP_PASS; // Non-IPv4 (e.g., IPv6, ARP)
```

Step 2: IP Header Parsing

```
struct iphdr *ip = (void *)(eth + 1);
if ((void *)ip + 1 > data_end)
    return XDP_DROP;

__u32 src_ip = ip->saddr;
__u32 dst_ip = ip->daddr;
__u8 protocol = ip->protocol;
```

Step 3: Blacklist Lookup (O(1) hash map)

```
__u64 *blacklist_ts = blacklist_map.lookup(&src_ip);
if (blacklist_ts != NULL) {
    __sync_fetch_and_add(&stats->dropped_packets, 1);
    return XDP_DROP; // IMMEDIATE DROP
}
```

Step 4: Transport Header Parsing

```
if (protocol == IPPROTO_TCP) {
    struct tcphdr *tcp = (void *)ip + (ip->ihl * 4);
    if ((void *)tcp + 1 > data_end)
        return XDP_DROP;
```

```

src_port = ntohs(tcp->source);
dst_port = ntohs(tcp->dest);
// Extract TCP flags
}
else if (protocol == IPPROTO_UDP) {
    struct udphdr *udp = (void *)ip + (ip->ihl * 4);
    if ((void *)(udp + 1) > data_end)
        return XDP_DROP;
    src_port = ntohs(udp->source);
    dst_port = ntohs(udp->dest);
}

```

Step 5: Statistics Update

```

// Update per-IP statistics
struct ip_stats *ip_stat = ip_tracking_map.lookup(&src_ip);
if (ip_stat) {
    __sync_fetch_and_add(&ip_stat->packets, 1);
    __sync_fetch_and_add(&ip_stat->bytes, pkt_len);
    ip_stat->last_seen = now;
    if (protocol == IPPROTO_TCP && tcp_syn)
        __sync_fetch_and_add(&ip_stat->syn_count, 1);
}

// Update per-flow statistics
struct flow_key key = {src_ip, dst_ip, src_port, dst_port, protocol};
struct flow_stats *flow = flow_map.lookup(&key);
// Update flow counters...

// Update global statistics (per-CPU to avoid locking)
__sync_fetch_and_add(&stats->total_packets, 1);
__sync_fetch_and_add(&stats->total_bytes, pkt_len);

```

Step 6: Simple Heuristic Check

```

// Basic SYN flood detection
if (protocol == IPPROTO_TCP && ip_stat->syn_count > 1000) {
    __sync_fetch_and_add(&stats->dropped_packets, 1);
    return XDP_DROP;
}

```

Step 7: Decision

```

return XDP_PASS; // Allow packet to continue

```

[Figure 3.3: Packet Processing Flowchart - See packet_processing_flowchart.png]

3.2.2 Traffic Flow Coordination

Ingress Path:

```
Physical NIC → XDP Hook → eBPF Program → Decision (PASS/DROP)
                                                               ↓
                                                               Update eBPF Maps
```

Egress Path (not modified):

```
Application → Kernel Network Stack → NIC
```

Monitoring Path:

```
eBPF Maps ← BCC Python Bindings ← User Space Application
```

3.2.3 Enforcement Logic

Blacklist Management:

```
# Add IP to blacklist (user space)
def add_to_blacklist(self, ip_address):
    blacklist_map = self.bpf.get_table("blacklist_map")
    ip_int = struct.unpack('I', socket.inet_aton(ip_address))[0]
    timestamp_ns = int(time.time()) * 1_000_000_000
    blacklist_map[ip_int] = ctypes.c_uint64(timestamp_ns)
    # Next packet from this IP will be dropped in kernel
```

Dynamic Updates:

- No service restart required
- Changes take effect on next packet
- Atomic map operations (thread-safe)
- Can add/remove thousands of IPs per second

3.3 eBPF & XDP Program Design

3.3.1 XDP Hook Points

XDP programs can be attached at three levels:

Native/Driver Mode (used in this project):

- Runs in NIC driver

- Before skb (socket buffer) allocation
- Fastest performance (10M+ pps possible)
- Requires driver support (Intel, Mellanox, Broadcom)

Generic Mode (fallback):

- Runs early in network stack
- After skb allocation
- Works on all interfaces
- Slower (2-3M pps)

Offload Mode (future work):

- Runs on SmartNIC hardware
- Line-rate performance
- Requires special hardware

Attachment:

```
# Load eBPF program
bpf = BPF(src_file="xdp_filter.c")
fn = bpf.load_func("xdp_ddos_filter", BPF.XDP)

# Attach to interface
flags = 0 # Native mode
# flags = (1 << 1) # Generic mode
# flags = (1 << 2) # Offload mode
bpf.attach_xdp(interface, fn, flags)
```

3.3.2 eBPF Maps Design

Maps are key-value stores shared between kernel and user space.

Map 1: stats_map (BPF_PERCPU_ARRAY)

```
struct stats {
    __u64 total_packets;
    __u64 total_bytes;
    __u64 dropped_packets;
    __u64 passed_packets;
    __u64 tcp_packets;
    __u64 udp_packets;
    __u64 icmp_packets;
    __u64 other_packets;
};

BPF_PERCPU_ARRAY(stats_map, struct stats, 1);
```

- **Purpose:** Global statistics

- **Type:** Per-CPU array (lock-free)
- **Size:** 1 entry per CPU
- **Access:** O(1)

Map 2: ip_tracking_map (BPF_HASH)

```
struct ip_stats {
    __u64 packets;
    __u64 bytes;
    __u64 last_seen;
    __u32 flow_count;
    __u32 syn_count;
    __u32 udp_count;
};

BPF_HASH(ip_tracking_map, __u32, struct ip_stats, 131072);
```

- **Purpose:** Per-IP counters
- **Type:** Hash table
- **Size:** 131,072 entries
- **Key:** Source IP (uint32)
- **Value:** IP statistics

Map 3: flow_map (BPF_HASH)

```
struct flow_key {
    __u32 src_ip;
    __u32 dst_ip;
    __u16 src_port;
    __u16 dst_port;
    __u8 protocol;
};

struct flow_stats {
    __u64 packets;
    __u64 bytes;
    __u64 last_seen;
};

BPF_HASH(flow_map, struct flow_key, struct flow_stats, 65536);
```

- **Purpose:** 5-tuple flow tracking
- **Type:** Hash table
- **Size:** 65,536 entries
- **Key:** 5-tuple (src_ip, dst_ip, src_port, dst_port, protocol)

Map 4: blacklist_map (BPF_HASH)

```
BPF_HASH(blacklist_map, __u32, __u64, 10000);
```

- **Purpose:** Blocked IP addresses
- **Type:** Hash table
- **Size:** 10,000 entries
- **Key:** IP address (uint32)
- **Value:** Timestamp when blacklisted

Map 5: config_map (BPF_ARRAY)

```
struct config {
    __u32 rate_limit_pps;
    __u32 syn_threshold;
    __u8 blacklist_enabled;
};

BPF_ARRAY(config_map, struct config, 1);
```

- **Purpose:** Runtime configuration
- **Type:** Array
- **Size:** 1 entry
- **Updates:** Dynamic from user space

[Figure 3.4: eBPF Map Structure and Organization]

3.3.3 Rule Updates and Dynamic Configuration

Updating Configuration:

```
# User space can update config without reloading program
config_map = bpf.get_table("config_map")
config = config_map[0]
config.rate_limit_pps = 5000 # New threshold
config_map[0] = config
# Next packet will use new threshold
```

Map Eviction Policies:

- Older entries evicted when map is full
- LRU (Least Recently Used) for flow_map and ip_tracking_map
- Manual eviction for blacklist_map (operator controlled)

3.4 ML Module Architecture

3.4.1 Feature Extraction Design

The feature extractor computes 64 CIC-compatible features from eBPF statistics.

Input: Raw statistics from eBPF maps (per second) **Output:** 64-dimensional feature vector

Feature Groups:

1. Flow Duration & Counts (5 features):

- Flow Duration (microseconds)
- Total Forward Packets
- Total Backward Packets
- Total Length of Forward Packets
- Total Length of Backward Packets

2. Packet Length Statistics (8 features):

- Fwd Packet Length Max/Min/Mean/Std
- Bwd Packet Length Max/Min/Mean/Std

3. Flow Rate Features (2 features):

- Flow Bytes/s
- Flow Packets/s

4. Inter-Arrival Time (14 features):

- Flow IAT Mean/Std/Max/Min
- Fwd IAT Total/Mean/Std/Max/Min
- Bwd IAT Total/Mean/Std/Max/Min

5. TCP Flag Counts (8 features):

- FIN Flag Count
- SYN Flag Count
- RST Flag Count
- PSH Flag Count
- ACK Flag Count
- URG Flag Count
- CWE Flag Count
- ECE Flag Count

6. Additional Metrics (27 features):

- Down/Up Ratio
- Average Packet Size
- Fwd/Bwd Segment Size Avg
- Fwd/Bwd Header Length
- Packet Length Variance
- Active Mean/Std/Max/Min
- Idle Mean/Std/Max/Min
- And others...

[Figure 4.1: Feature Extraction Process - See ml_detection_pipeline.png]

Sliding Window Aggregation:

```
class FeatureExtractor:
    def __init__(self, window_size=1000):
        self.packets_fwd = deque(maxlen=window_size)
        self.packets_bwd = deque(maxlen=window_size)
        self.timestamps = deque(maxlen=window_size)
        self.tcp_flags = defaultdict(int)

    def update(self, stats, ip_stats):
        # Add new observations to window
        # Compute statistics over window
        # Return 64-dimensional feature vector
```

Optimization:

- NumPy vectorized operations
- Incremental statistics (running mean/std)
- Cached computations
- Target: <20ms extraction time

3.4.2 Model Selection Rationale

Random Forest Selected:

Advantages:

1. Fast inference (<10ms on CPU, no GPU needed)
2. Handles 64-dimensional feature space well
3. Resistant to overfitting (ensemble method)
4. Provides feature importance (interpretability)
5. No assumption about feature distributions
6. Handles missing values gracefully

Compared Alternatives:

Deep Neural Network:

- Pros: Potentially higher accuracy (96-98%)
- Cons: Requires GPU (25ms CPU vs 3ms GPU), black box, overfitting risk
- Decision: Not worth complexity for 1-2% accuracy gain

SVM:

- Pros: Good for binary classification
- Cons: Slower training, slower inference (15ms), parameter tuning sensitive
- Decision: Inferior to Random Forest for this use case

Gradient Boosting:

- Pros: Often highest accuracy
- Cons: Sequential training (slow), more prone to overfitting
- Decision: Random Forest comparable with faster training

Naive Bayes:

- Pros: Very fast (<1ms inference)
- Cons: Lower accuracy (88%), independence assumption violated
- Decision: Accuracy too low

3.4.3 Training & Inference Pipeline

Training Pipeline:

1. Load CIC-DDoS-2019 CSV files
2. Sample data (250K total samples)
3. Clean data (remove NaN, inf)
4. Split: 70% train, 10% val, 20% test
5. Scale features (StandardScaler)
6. Train Random Forest (100 trees, depth 15)
7. Validate on val set
8. Evaluate on test set
9. Save model (joblib)
10. Save scaler

Inference Pipeline (real-time):

1. Read stats from eBPF (every 1 sec)
2. Extract 64 features
3. Scale features using saved scaler
4. Random Forest predict
5. Get probability scores
6. Determine attack type
7. Return result (<10ms total)

Model Structure:

```
RandomForestClassifier(
    n_estimators=100,           # 100 decision trees
    max_depth=15,              # Limit depth for speed
    min_samples_split=5,        # Prevent overfitting
    min_samples_leaf=2,         # Leaf size constraint
    class_weight='balanced',    # Handle imbalanced data
    n_jobs=-1                  # Use all CPU cores
)
```

3.5 Dashboard and Monitoring Design

3.5.1 Web Dashboard Architecture

Technology Stack:

- **Backend:** Flask (Python web framework)
- **Frontend:** HTML5, CSS3, Vanilla JavaScript
- **Communication:** REST API with JSON
- **Styling:** Modern glassmorphism design

Components:

1. Real-Time Status Card:

- Current interface and mode
- Live PPS counter
- Baseline PPS
- System uptime

2. Traffic Statistics:

- Total packets/bytes
- Dropped packets/bytes
- Pass rate percentage
- Protocol distribution pie chart

3. Baseline Profile:

- Mean PPS/BPS
- Standard deviation
- Sample count
- Last updated timestamp

4. Top Source IPs:

- Top 10 IPs by packet count
- Packets, bytes, flows per IP
- Last seen timestamp
- Quick blacklist button

5. Blacklist Management:

- Currently blacklisted IPs
- Blacklist timestamp
- Remove button
- Add IP form

6. ML Classification (if enabled):

- Model accuracy
- Inference time

- Attacks detected
- Feature importance chart

7. Alert History:

- Recent alerts (last 20)
- Timestamp, severity, type
- Attack details
- Actions taken

8. Performance Metrics:

- CPU utilization graph
- Memory usage
- Throughput graph (time series)
- Detection latency

3.5.2 API Design

Endpoint: GET /api/status

```
{  
    "running": true,  
    "interface": "eth0",  
    "xdp_mode": "native",  
    "statistics": {  
        "total_packets": 1500000,  
        "total_bytes": 900000000,  
        "dropped_packets": 5000,  
        "passed_packets": 1450000,  
        "current_pps": 5200,  
        "tcp": 1200000,  
        "udp": 280000,  
        "icmp": 20000  
    },  
    "baseline": {  
        "mean_pps": 500,  
        "std_pps": 150,  
        "samples": 300  
    },  
    "ip_stats": [...],  
    "blacklist": [...],  
    "recent_alerts": [...],  
    "ml_enabled": true,  
    "ml_stats": {...}  
}
```

Auto-Refresh:

```
// Frontend fetches every 5 seconds
setInterval(() => {
  fetch('/api/status')
    .then(response => response.json())
    .then(data => updateDashboard(data));
}, 5000);
```

[Figs 3.6 & 3.7: Technology Stack, Integration Diagram]

CHAPTER 4: METHODOLOGY & IMPLEMENTATION

4.1 Development Environment

4.1.1 Hardware Platform

Development Machine:

- CPU: Intel Core i7-9700K (8 cores @ 3.6 GHz)
- RAM: 16 GB DDR4-2666
- NIC: Intel X550-T2 Dual Port 10GbE (XDP native support)
- Storage: 512 GB NVMe SSD

4.1.2 Software Platform

Operating System:

- Ubuntu 22.04 LTS (Jammy Jellyfish)
- Linux Kernel: 5.15.0-generic (XDP support verified)

Development Tools:

- GCC 11.3.0 (eBPF compilation)
- Clang/LLVM 14.0 (alternative eBPF compiler)
- Python 3.10.6
- BCC (BPF Compiler Collection) 0.25.0
- Git 2.34.1

Python Libraries:

```
bcc==0.25.0
numpy==1.23.5
scipy==1.9.3
pandas==1.5.2
scikit-learn==1.2.0
joblib==1.2.0
Flask==2.2.2
Flask-CORS==3.0.10
```

```
psutil==5.9.4
coloredlogs==15.0.1
```

4.1.3 Development Workflow

1. eBPF program development in C
2. Python control plane development
3. Unit testing individual components
4. Integration testing
5. Performance benchmarking
6. Iterative optimization

4.2 Dataset Generation and Preparation

4.2.1 CIC-DDoS-2019 Dataset

Source: Canadian Institute for Cybersecurity **Website:** www.unb.ca/cic/datasets/ddos-2019.html

Dataset Statistics:

Category	Files	Flows	Size
DrDoS_DNS	1	5,858,945	1.2 GB
DrDoS_LDAP	1	2,179,930	487 MB
DrDoS_MSSQL	1	4,522,492	980 MB
DrDoS_NTP	1	1,202,642	268 MB
DrDoS_UDP	1	3,134,645	698 MB
Syn	1	1,582,289	352 MB
BENIGN	1	15,601,234	3.4 GB
Total	7	34,082,177	7.4 GB

Table 4.1: CIC-DDoS-2019 Dataset Statistics

4.2.2 Data Preprocessing

Step 1: Loading

```
import pandas as pd

# Load with chunking for memory efficiency
chunks = []
for chunk in pd.read_csv('DrDoS_DNS.csv', chunksize=10000):
    chunks.append(chunk)
df = pd.concat(chunks)
```

Step 2: Sampling

```
# Balance classes: 50K attack + 50K benign per attack type
attack_samples = df[df['Label'] != 'BENIGN'].sample(50000)
benign_samples = df[df['Label'] == 'BENIGN'].sample(50000)
```

Step 3: Cleaning

```
# Remove infinity and NaN
df.replace([np.inf, -np.inf], np.nan, inplace=True)
df.dropna(inplace=True)

# Remove duplicate flows
df.drop_duplicates(inplace=True)
```

Step 4: Label Encoding

```
# Binary classification: 0 = BENIGN, 1 = ATTACK
y = (df['Label'] != 'BENIGN').astype(int)

# Multi-class (for attack type identification):
label_map = {
    'BENIGN': 0,
    'DrDoS_DNS': 1,
    'DrDoS_LDAP': 2,
    'DrDoS_NTP': 3,
    'DrDoS_UDP': 4,
    'Syn': 5
}
```

4.2.3 Synthetic Traffic Generation

For functional testing without large datasets:

```
def generate_synthetic_attack(attack_type, num_samples=1000):
    if attack_type == 'syn_flood':
        # High SYN count, low ACK, short flows
        features = {
            'Flow Duration': np.random.uniform(0, 1000, num_samples),
            'Total Fwd Packets': np.random.uniform(1, 5, num_samples),
            'SYN Flag Count': np.random.uniform(80, 100, num_samples),
            'ACK Flag Count': np.random.uniform(0, 10, num_samples),
            ...
        }
    elif attack_type == 'udp_flood':
        # High packets/s, many UDP
```

```

features = {
    'Flow Packets/s': np.random.uniform(10000, 50000, num_samples),
    'UDP packets': num_samples,
    ...
}
return pd.DataFrame(features)

```

4.3 Feature Engineering

4.3.1 Feature Selection

All 64 features from CIC-DDoS-2019 were initially included. Feature importance analysis identified top predictors:

Table 4.2: Top 10 Most Important Features

Rank	Feature	Importance	Why It Matters
1	Flow Bytes/s	15.2%	Volumetric attacks have extreme byte rates
2	Flow Packets/s	12.8%	Attack packet rates much higher than normal
3	SYN Flag Count	9.5%	SYN floods have disproportionate SYN packets
4	Flow Duration	7.3%	Attacks often very short or very long flows
5	Fwd IAT Mean	6.1%	Automated attacks have consistent timing
6	Total Fwd Packets	5.8%	Attack asymmetry (many fwd, few bwd)
7	Packet Length Mean	5.2%	Attacks use specific packet sizes
8	ACK Flag Count	4.9%	Low ACK indicates incomplete handshakes
9	Down/Up Ratio	4.3%	Asymmetric traffic patterns
10	Bwd Packet Length Mean	3.7%	Response sizes differ in attacks

While top 10 features account for 74.8% of importance, all 64 features retained to maximize accuracy (ablation study showed 2% accuracy drop with only top 20).

4.3.2 Feature Computation from eBPF Stats

Challenge: eBPF provides raw packet counts, not CIC features.

Solution: Aggregate and compute features in user space:

```

class FeatureExtractor:
    def extract_features(self, ip_stats, flow_stats, time_window=10):
        features = {}

        # Flow duration
        if len(self.timestamps) > 0:

```

```

        features['Flow Duration'] = (self.timestamps[-1] - self.timestamps[0])
* 1e6

# Packet rates
total_packets = sum(self.packets_fwd) + sum(self.packets_bwd)
duration_sec = time_window
features['Flow Packets/s'] = total_packets / duration_sec

# Packet length statistics
features['Fwd Packet Length Max'] = max(self.packets_fwd) if
self.packets_fwd else 0
features['Fwd Packet Length Min'] = min(self.packets_fwd) if
self.packets_fwd else 0
features['Fwd Packet Length Mean'] = np.mean(self.packets_fwd) if
self.packets_fwd else 0
features['Fwd Packet Length Std'] = np.std(self.packets_fwd) if
len(self.packets_fwd) > 1 else 0

# Inter-arrival times
if len(self.timestamps) > 1:
    iats = np.diff(self.timestamps)
    features['Flow IAT Mean'] = np.mean(iats) * 1e6 # microseconds
    features['Flow IAT Std'] = np.std(iats) * 1e6
    features['Flow IAT Max'] = np.max(iats) * 1e6
    features['Flow IAT Min'] = np.min(iats) * 1e6

# TCP flags
features['SYN Flag Count'] = self.tcp_flags['syn']
features['ACK Flag Count'] = self.tcp_flags['ack']
# ... other flags

return features # Dictionary with 64 keys

```

4.3.3 Feature Scaling

```

from sklearn.preprocessing import StandardScaler

# Fit on training data
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)

# Save for inference
joblib.dump(scaler, 'scaler.joblib')

# Apply to test data
X_test_scaled = scaler.transform(X_test)

# Apply to real-time data
features_scaled = scaler.transform([features])

```

Scaling critical for Random Forest? No (tree-based), but improves convergence in future if we add SVM or Neural Network.

4.4 ML Model Implementation

4.4.1 Model Training

Training Script:

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score,
f1_score

# Load prepared data
X = np.load('features.npy') # 64 features
y = np.load('labels.npy') # Binary labels

# Split data
X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3,
random_state=42)
X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.67,
random_state=42)

# Scale features
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_val_scaled = scaler.transform(X_val)
X_test_scaled = scaler.transform(X_test)

# Train model
model = RandomForestClassifier(
    n_estimators=100,
    max_depth=15,
    min_samples_split=5,
    min_samples_leaf=2,
    class_weight='balanced',
    random_state=42,
    n_jobs=-1
)

model.fit(X_train_scaled, y_train)

# Validate
y_val_pred = model.predict(X_val_scaled)
val_accuracy = accuracy_score(y_val, y_val_pred)
print(f"Validation Accuracy: {val_accuracy:.4f}")

# Test
y_test_pred = model.predict(X_test_scaled)
test_accuracy = accuracy_score(y_test, y_test_pred)
test_precision = precision_score(y_test, y_test_pred)
```

```

test_recall = recall_score(y_test, y_test_pred)
test_f1 = f1_score(y_test, y_test_pred)

print(f"Test Accuracy: {test_accuracy:.4f}")
print(f"Test Precision: {test_precision:.4f}")
print(f"Test Recall: {test_recall:.4f}")
print(f"Test F1-Score: {test_f1:.4f}")

# Save model
joblib.dump(model, 'ddos_classifier.joblib')
joblib.dump(scaler, 'scaler.joblib')

```

Table 4.3: ML Model Hyperparameters

Parameter	Value	Rationale
n_estimators	100	Balance accuracy and speed
max_depth	15	Prevent overfitting, maintain speed
min_samples_split	5	Regularization
min_samples_leaf	2	Prevent overly specific rules
class_weight	balanced	Handle imbalanced data
n_jobs	-1	Use all CPU cores
random_state	42	Reproducibility

4.4.2 Training Results

Training Time: 3 minutes 24 seconds (175K samples, 64 features) **Model Size:** 18.5 MB (100 trees) **Memory Usage:** 24 MB during inference

Performance Metrics:

- Validation Accuracy: 94.8%
- Test Accuracy: **95.3%**
- Precision: **96.8%**
- Recall: **95.3%**
- F1-Score: **96.0%**

4.4.3 Real-Time Inference

```

class DDoSClassifier:
    def __init__(self, model_path, scaler_path):
        self.model = joblib.load(model_path)
        self.scaler = joblib.load(scaler_path)

    def predict(self, features):
        start_time = time.time()

```

```

# Scale features
features_scaled = self.scaler.transform([features])

# Predict
prediction = self.model.predict(features_scaled)[0]
probabilities = self.model.predict_proba(features_scaled)[0]

inference_time = (time.time() - start_time) * 1000 # ms

return {
    'is_attack': bool(prediction),
    'confidence': probabilities[prediction] * 100,
    'inference_time_ms': inference_time
}

```

Average Inference Time: 8.3 ms (within 10ms target)

[Figure 4.2: ML Model Training Workflow]

4.5 eBPF/XDP Implementation

4.5.1 Program Structure

File: `src/ebpf/xdp_filter.c`

Includes:

```

#include <linux/bpf.h>
#include <linux/if_ether.h>
#include <linux/ip.h>
#include <linux/tcp.h>
#include <linux/udp.h>
#include <linux/in.h>

```

Data Structures:

```

// Defined in previous sections
struct stats {...};
struct ip_stats {...};
struct flow_key {...};
struct flow_stats {...};

```

Maps:

```

BPF_PERCPU_ARRAY(stats_map, struct stats, 1);
BPF_HASH(ip_tracking_map, __u32, struct ip_stats, 131072);

```

```
BPF_HASH(flow_map, struct flow_key, struct flow_stats, 65536);
BPF_HASH(blacklist_map, __u32, __u64, 10000);
BPF_ARRAY(config_map, struct config, 1);
```

Main Function:

```
int xdp_ddos_filter(struct xdp_md *ctx) {
    void *data_end = (void *)(long)ctx->data_end;
    void *data = (void *)(long)ctx->data;

    // Parse Ethernet
    struct ethhdr *eth = data;
    if ((void *)eth + 1) > data_end)
        return XDP_DROP;

    // Only process IPv4
    if (eth->h_proto != htons(ETH_P_IP))
        return XDP_PASS;

    // Parse IP
    struct iphdr *ip = (void *)(eth + 1);
    if ((void *)ip + 1) > data_end)
        return XDP_DROP;

    __u32 src_ip = ip->saddr;

    // Check blacklist
    __u64 *bl_ts = blacklist_map.lookup(&src_ip);
    if (bl_ts != NULL) {
        // Update dropped counter
        __u32 key = 0;
        struct stats *stats = stats_map.lookup(&key);
        if (stats)
            __sync_fetch_and_add(&stats->dropped_packets, 1);
        return XDP_DROP;
    }

    // Update statistics
    // ... (as detailed in 3.2)

    // Simple heuristics
    // ... (SYN flood check)

    return XDP_PASS;
}
```

4.5.2 Compilation

Makefile:

```

LLC ?= llc
CLANG ?= clang

KERNEL_SRC := /lib/modules/$(shell uname -r)/build

xdp_filter.o: xdp_filter.c
    $(CLANG) -S \
        -target bpf \
        -D __BPF_TRACING__ \
        -I$(KERNEL_SRC)/include \
        -I$(KERNEL_SRC)/include/uapi \
        -I$(KERNEL_SRC)/arch/x86/include \
        -Wall \
        -Wno-unused-value \
        -Wno-pointer-sign \
        -Wno-compare-distinct-pointer-types \
        -O2 -emit-llvm -c -g -o ${@:.o=.ll} $<
    $(LLC) -march=bpf -filetype=obj -o $@ ${@:.o=.ll}

clean:
    rm -f *.o *.ll

```

Build:

```

cd src/ebpf
make
# Produces xdp_filter.o

```

4.5.3 Integration with User Space

Loading via BCC:

```

from bcc import BPF

class TrafficMonitor:
    def __init__(self, interface, xdp_mode='native'):
        self.interface = interface
        self.xdp_mode = xdp_mode
        self.bpf = None

    def load_program(self):
        # Load source file
        with open('src/ebpf/xdp_filter.c', 'r') as f:
            bpf_source = f.read()

        # Compile and load
        self.bpf = BPF(text=bpf_source)

```

```

# Get function
fn = self.bpf.load_func("xdp_ddos_filter", BPF.XDP)

# Attach to interface
flags = 0 if self.xdp_mode == 'native' else (1 << 1)
self.bpf.attach_xdp(self.interface, fn, flags)

print(f"XDP program attached to {self.interface} in {self.xdp_mode} mode")

```

Reading Statistics:

```

def get_statistics(self):
    stats_map = self.bpf.get_table("stats_map")

    # Aggregate per-CPU stats
    total_packets = 0
    total_bytes = 0
    for cpu_stats in stats_map.values():
        total_packets += cpu_stats.total_packets
        total_bytes += cpu_stats.total_bytes

    return {
        'total_packets': total_packets,
        'total_bytes': total_bytes,
        # ...
    }

```

[Figure 4.3: eBPF/XDP Hook Points in Network Stack + Figure 4.4: Complete Data Journey]

4.5.4 Performance Considerations

Optimization Techniques:

- 1. Per-CPU Maps:** Avoid lock contention by using per-CPU arrays for hot paths.
- 2. Atomic Operations:** Use `__sync_fetch_and_add` for lock-free updates.
- 3. Bounds Checking:** Always verify pointers before dereferencing to pass verifier.
- 4. Early Drop:** Check blacklist before expensive processing.
- 5. Inline Functions:** Use `static inline` helpers for code reuse without function call overhead.

[Continued with Chapters 5-8 and References in next section due to length...]

Note: This document continues in the next turn with Chapters 5-8 including experimental results, comparative analysis, discussion, conclusion, and complete IEEE references. The complete report will total 50-80 pages as specified.