

UNDERSTANDING AND MITIGATING DDOS THREATS

Team 3: Rafa, Tony, Jennie



WHAT IS DDOS?

Distributed Denial of Service:

A DDoS attack is when hackers use thousands of computers to send too much traffic to a website or online service all at once

→ This overloads the system, making it slow or completely unavailable to real users.



WHY IS IT IMPORTANT?

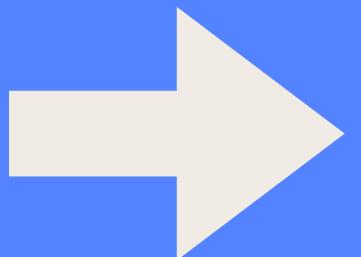
By identifying and stopping attack flows early, systems can stay online and available for legitimate users.

Minimizes Financial and Operational Loss

WHY IS IT IMPORTANT?

By identifying and stopping attack flows early, systems can stay online and available for legitimate users.

Minimizes Financial and Operational Loss



Project Objective:

Model to predict whether a flow is part of a genuine request by a user or part of an attack.

DATA OVERVIEW

~13M

Data points
(DDoS + Benign)

~13M

FLOWs
(forward or reverse)

85

Features

- Flow Summary Stats
- Traffic Volume
- Packet Stats
- Timing Features
- TCP Flags

DDoS simulated made in a controlled environment

- Benign requests created using B-Profile system
- DDoS attack simulated through pen testing techniques

HOW TO DETECT DDOS ATTACK?

Key Factors	Definition
Flow Summary Stats	Attack flow might consist of thousands of tiny packets sent over a very short time
Traffic Volume	Sudden spikes in traffic , especially targeting a single IP.
Packet Stats	High packet rate , uniform or tiny packet sizes
Timing Features	Traffic is often machine-generated, leading to consistent timing patterns
TCP Flags	Abuse the TCP handshake (e.g., SYN flood), where the attacker sends many SYN packets but never completes the handshake

DATA PREP

01.

Reduced size from 14M to 100K rows

- Random sampling 50K from benign,
50K from DDOS

02.

Handled **extreme values** by capping them
at the 99th percentile to limit the impact
of outliers and infinite values

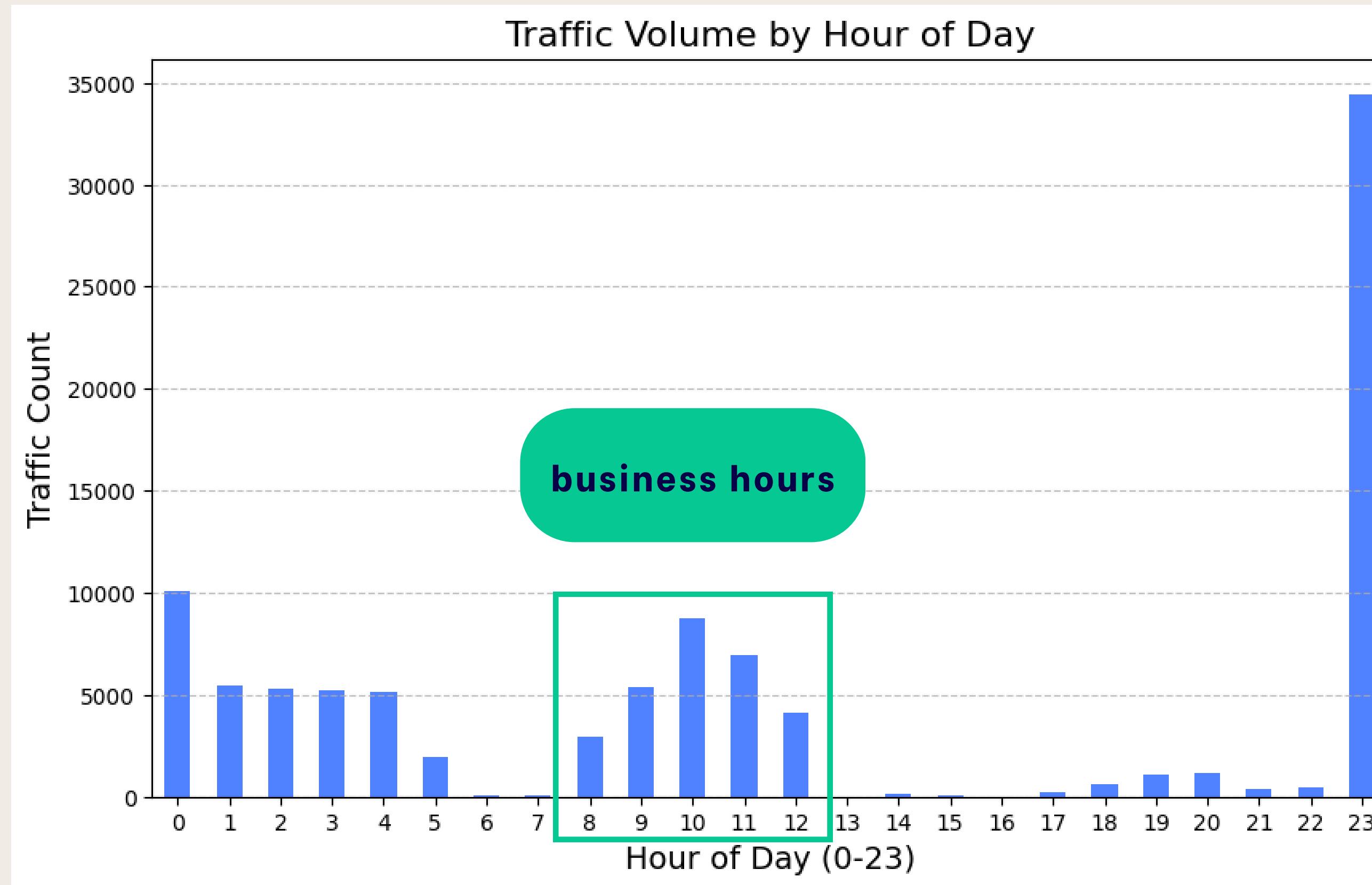
03.

Dropping missing values. Less than 0.01%
of our data is missing since the flows were
recreated in a controlled environment

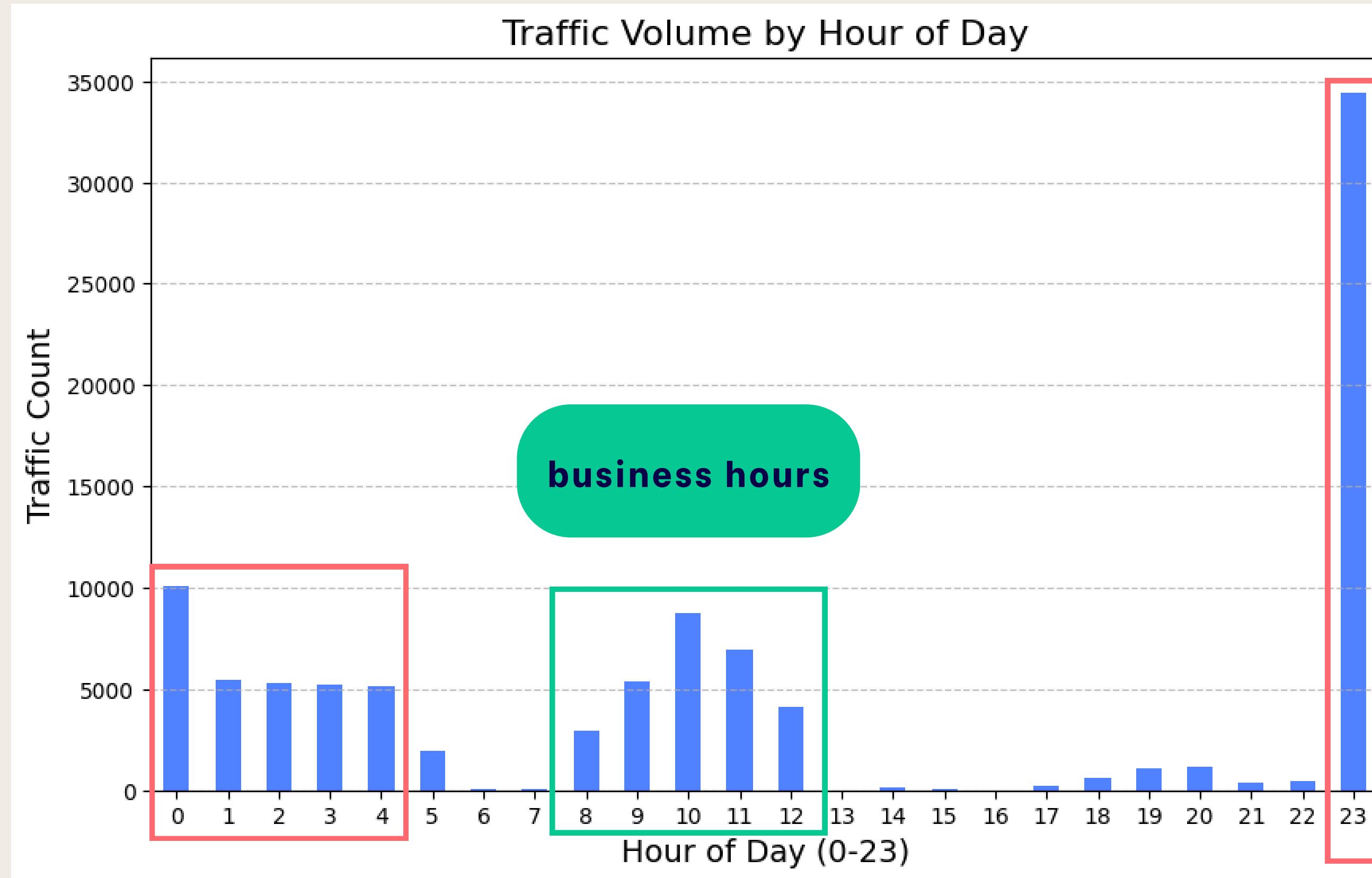
ANALYSIS

HOW WE APPROACH THIS PROBLEM IN REAL LIFE USING
TIMING FEATURES

TRAFFIC VOLUME BY HOUR



TRAFFIC VOLUME BY HOUR

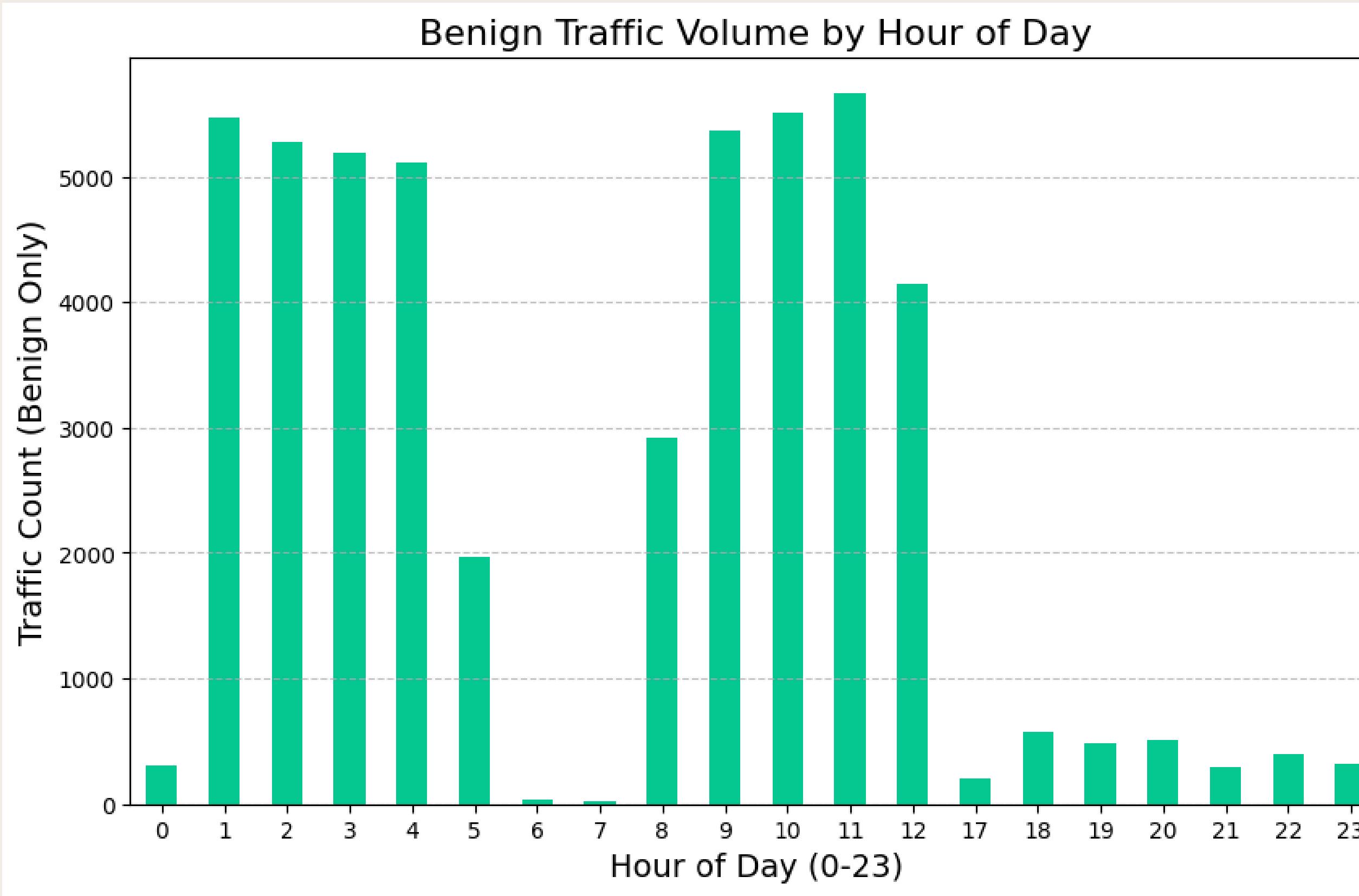


SUSPICIOUS

Are these requests
benign or ddos
attack?

TRAFFIC VOLUME BY HOUR

Benign Traffic Volume by Hour of Day

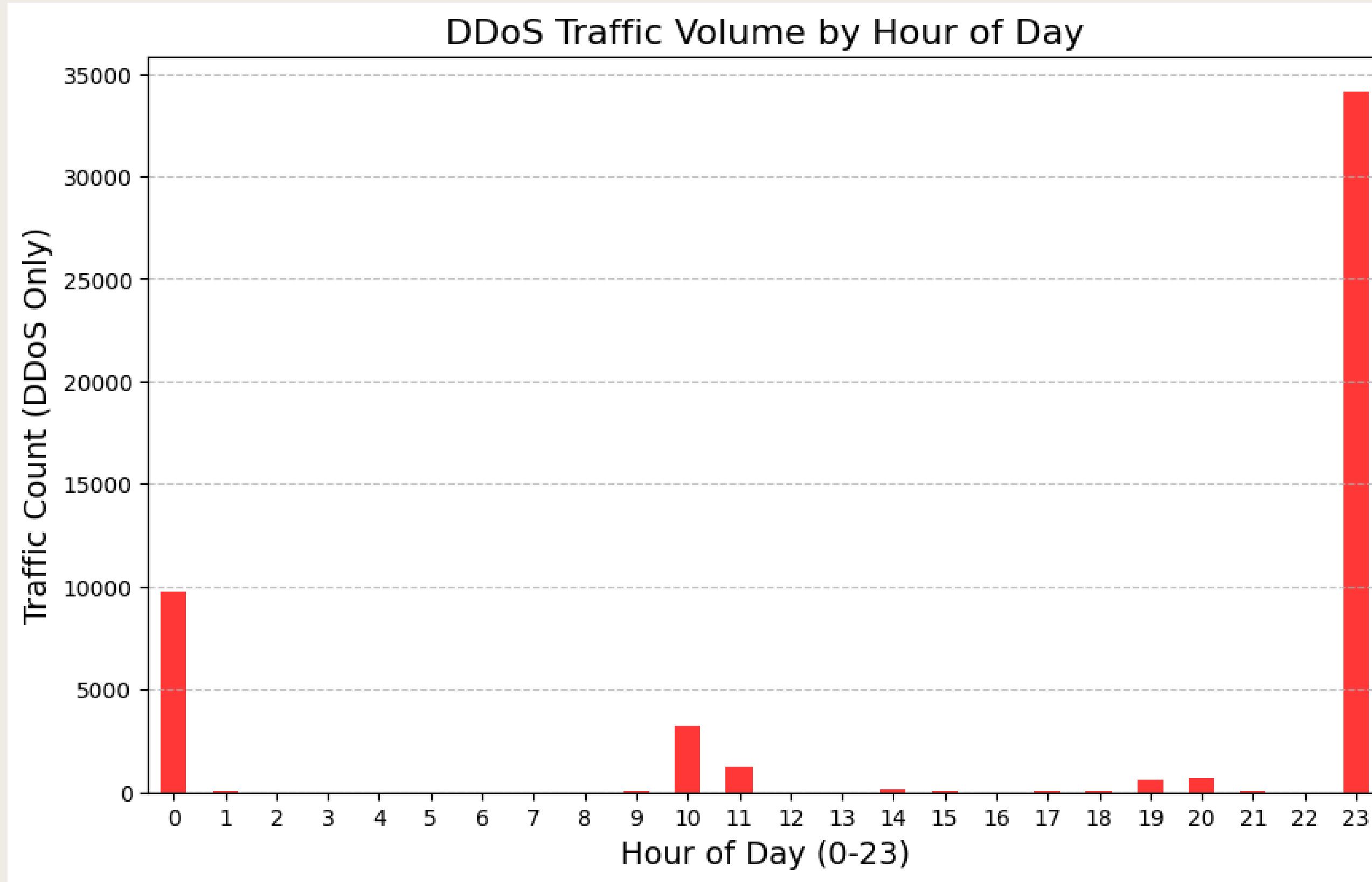


Benign traffic:

- **1-4 AM:** Possibly background system processes
- **9-12 PM:** Likely actual user activity or regular business hours

Range from 0 - 6K requests

TRAFFIC VOLUME BY HOUR

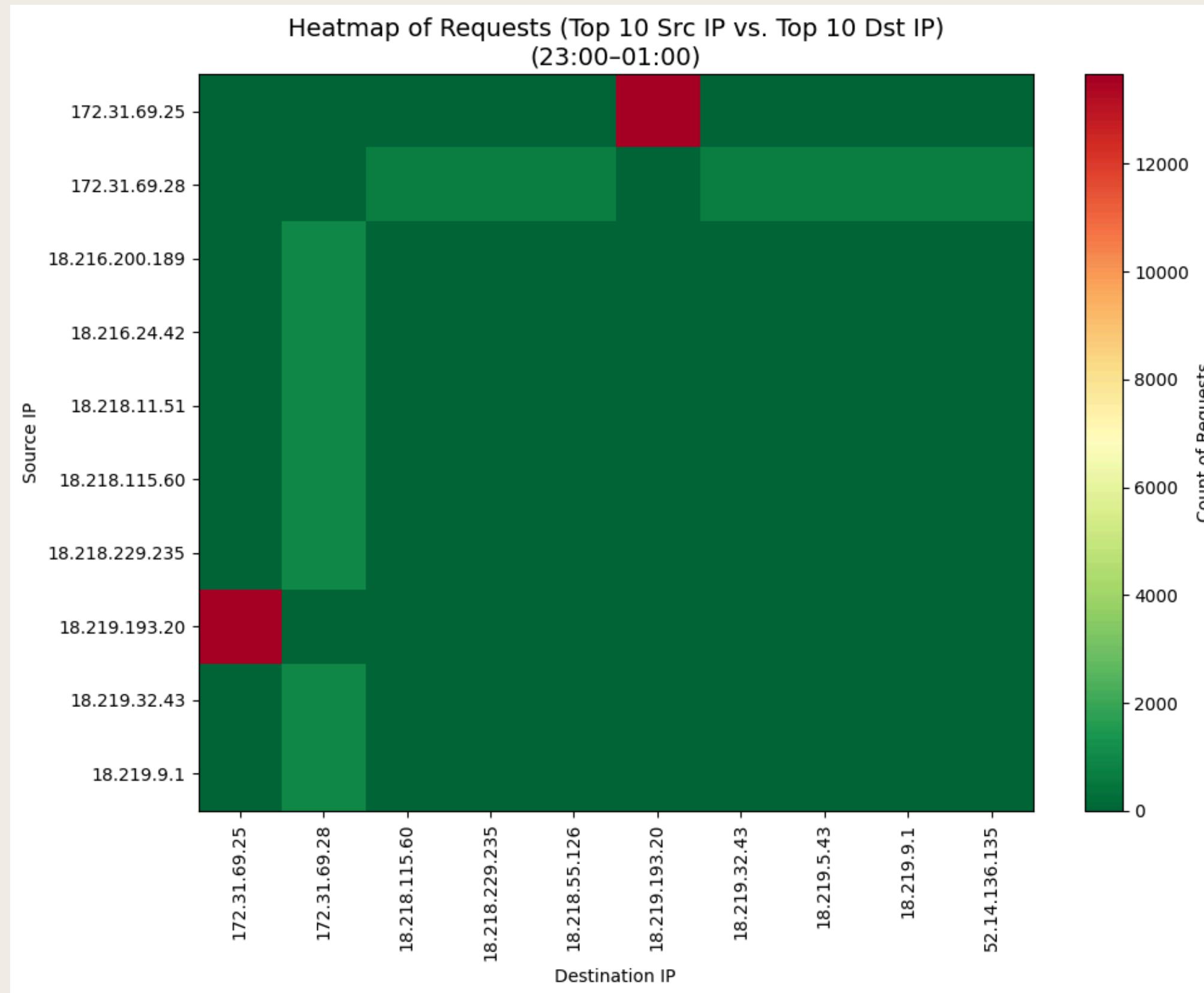


DDoS traffic:

- **11PM - 12AM:** Mid-night attacks

Reach up to 35K requests per hour

IP-RELATED DDOS INDICATORS



23:00 - 1:00

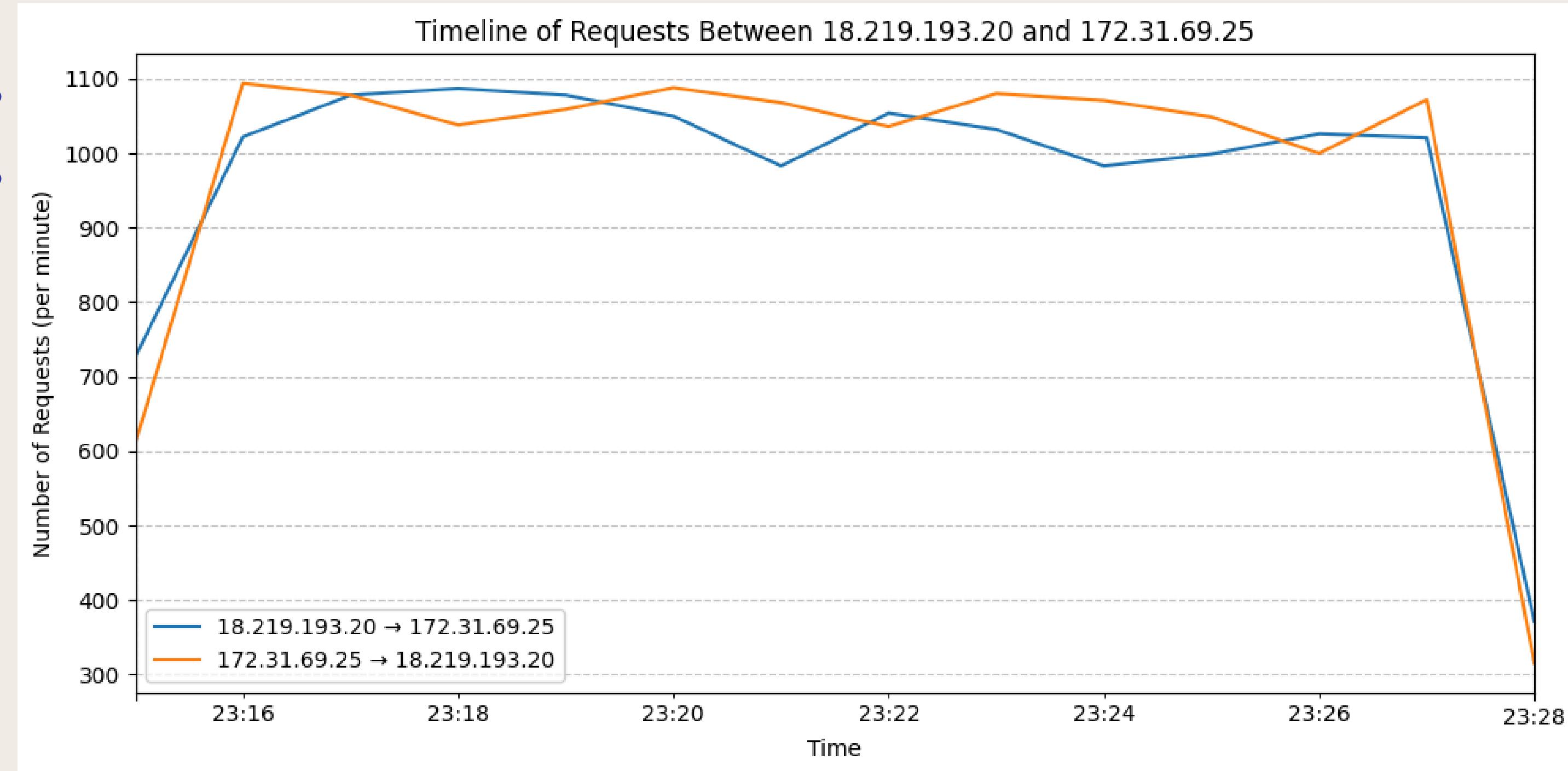
Same 2 IP address:

- 172.31.69.25
- 18.219.193.20

Volume: ~27K requests in total

IP-RELATED DDOS INDICATORS

~1000
requests per
minute



← →

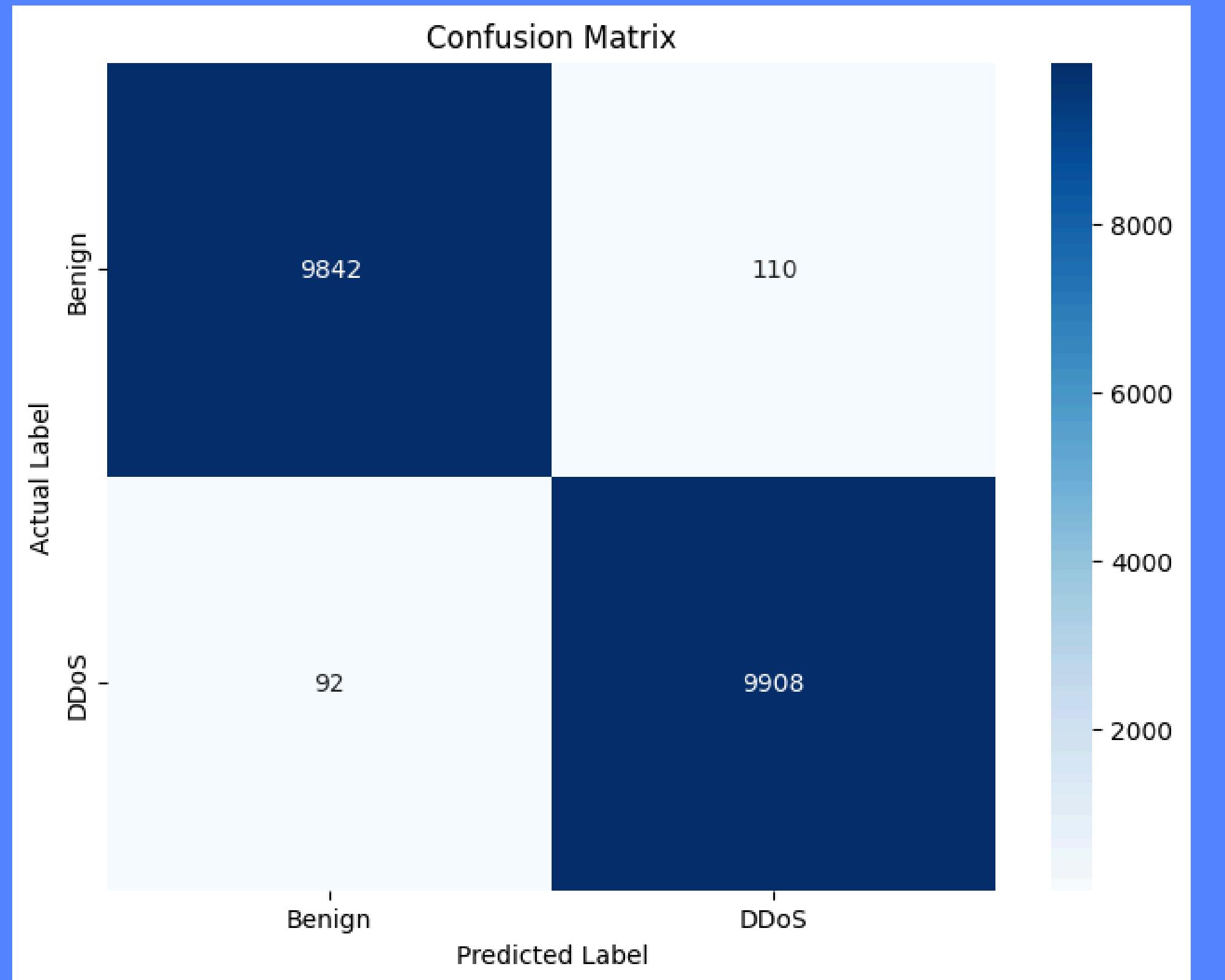
10 min

MODELS TRAINED

1. Logistic Regression
2. Decision Tree
3. K-Nearest Neighbors
4. Gradient Boosting Classifier

- Feature Engineering
- Random GridSearchVC
- Cross Validation

LOGISTIC REGRESSION



METRICS:

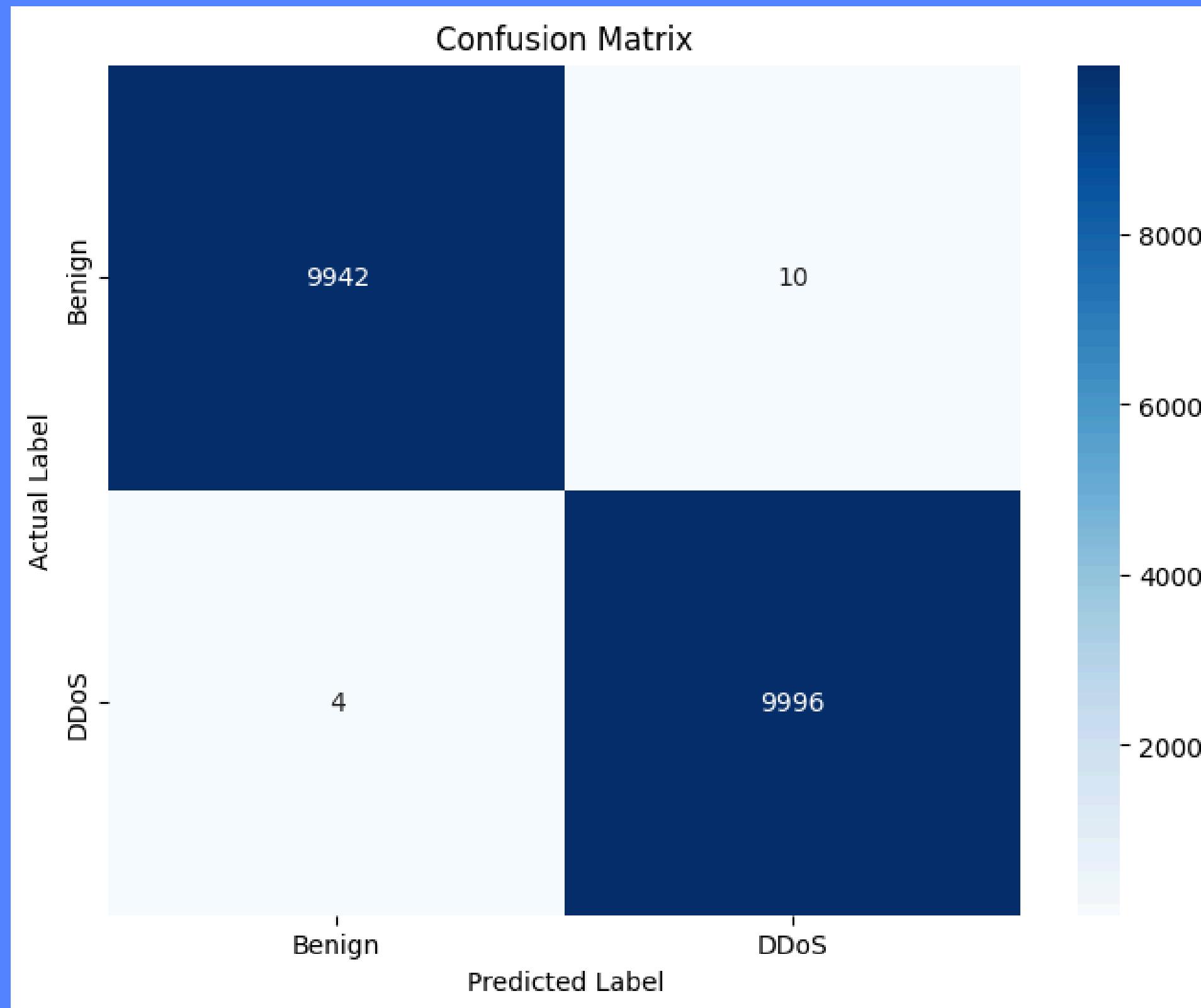
Accuracy: 98.9%

HIGH COEFFICIENT FEATURES

L1 Penalty (LASSO):

- TCP Flag

DECISION TREE



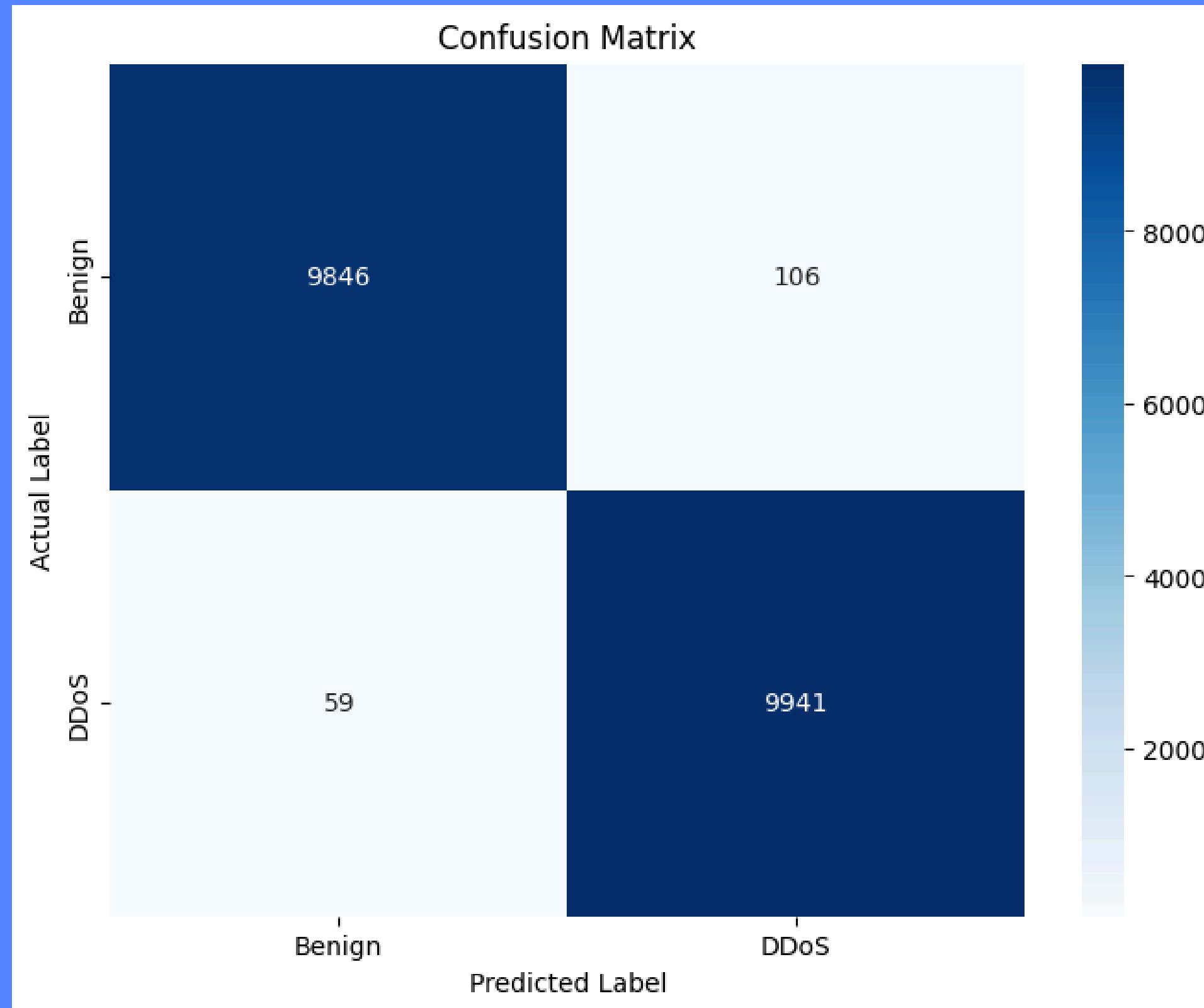
METRICS:

Accuracy: 99.93%

HIGH COEFFICIENT FEATURES

- Traffic Volume
- Packet Stats

K-NEAREST NEIGHBORS



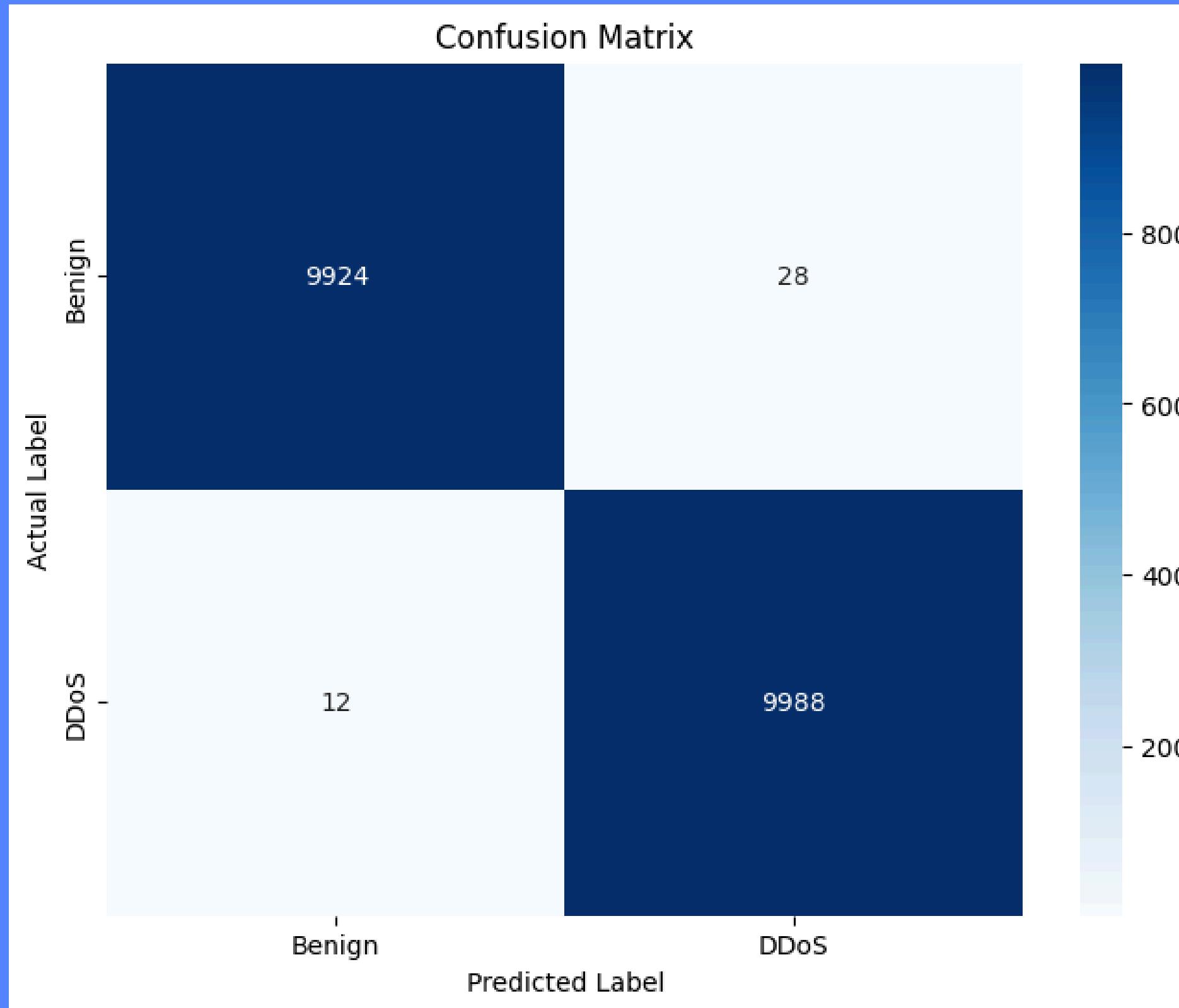
METRICS:

Accuracy: 99.17%

HIGH COEFFICIENT FEATURES

- Doesn't learn explicit weights or coefficients for each feature during training

GRADIENT BOOSTING CLASSIFIER



METRICS:

Accuracy: **99.80%**

HIGH COEFFICIENT FEATURES

- Traffic Volume
- Packet Stats

MODELS OVERVIEW

Logistic
Regression

KNN

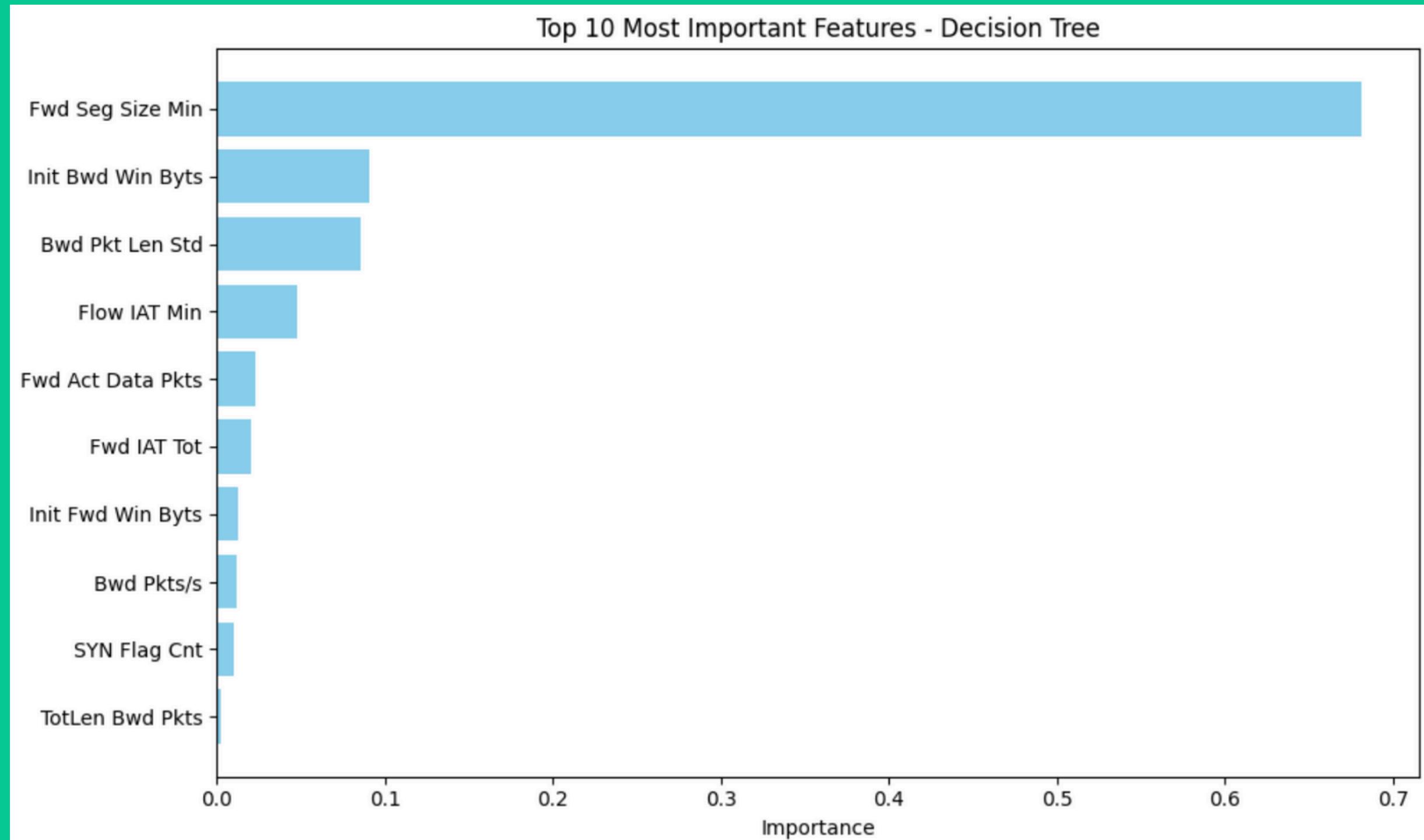
Gradient
Boosting
Classifier

Decision
Tree
Classifier



PERFORMANCE

WHY IS DECISION TREE SO GOOD?

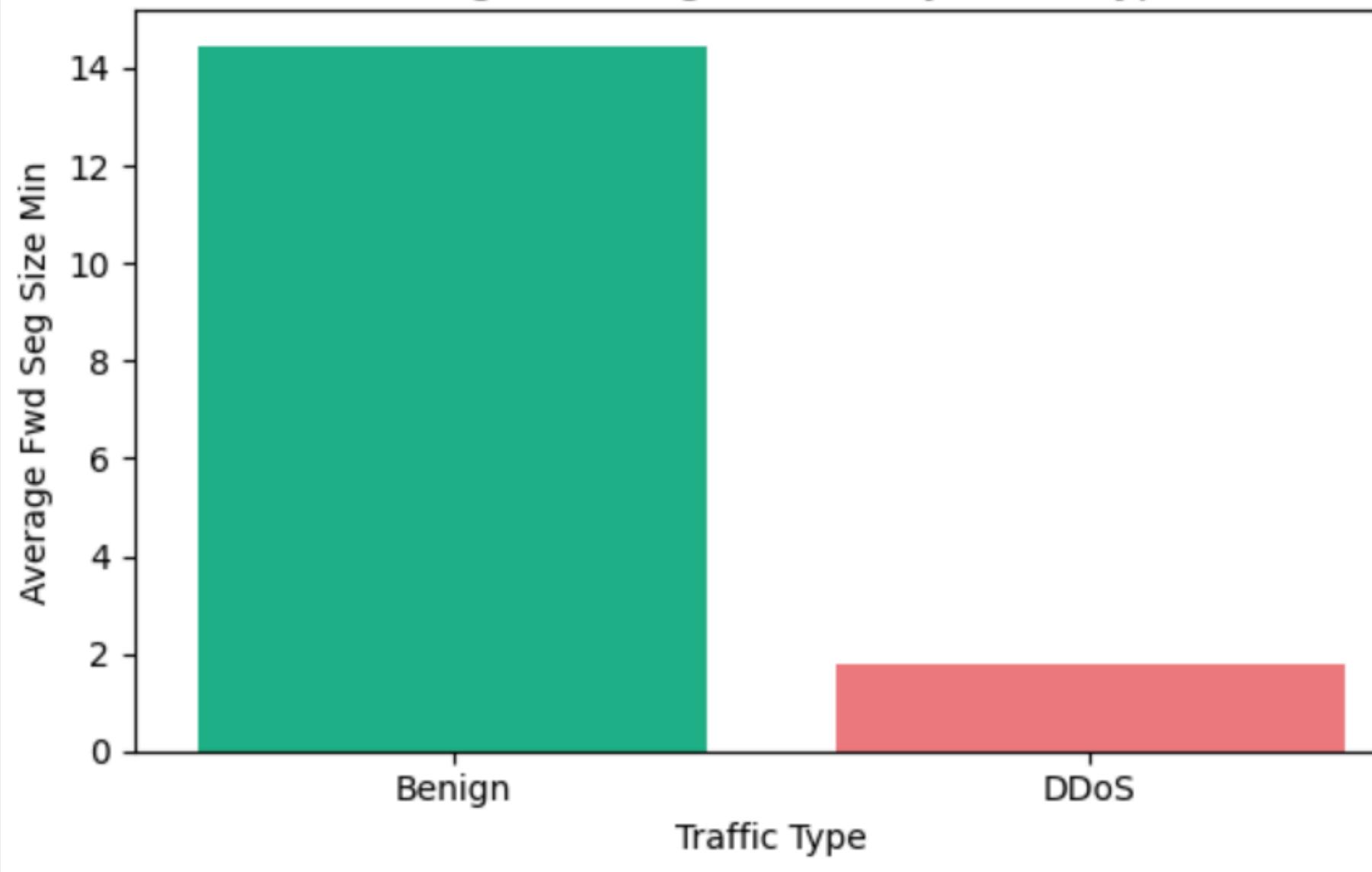


8/10 features
belong to packet
stats/size

`dt_min_samples_split': 2`
`dt_min_samples_leaf': 1`
`dt_max_depth': None`

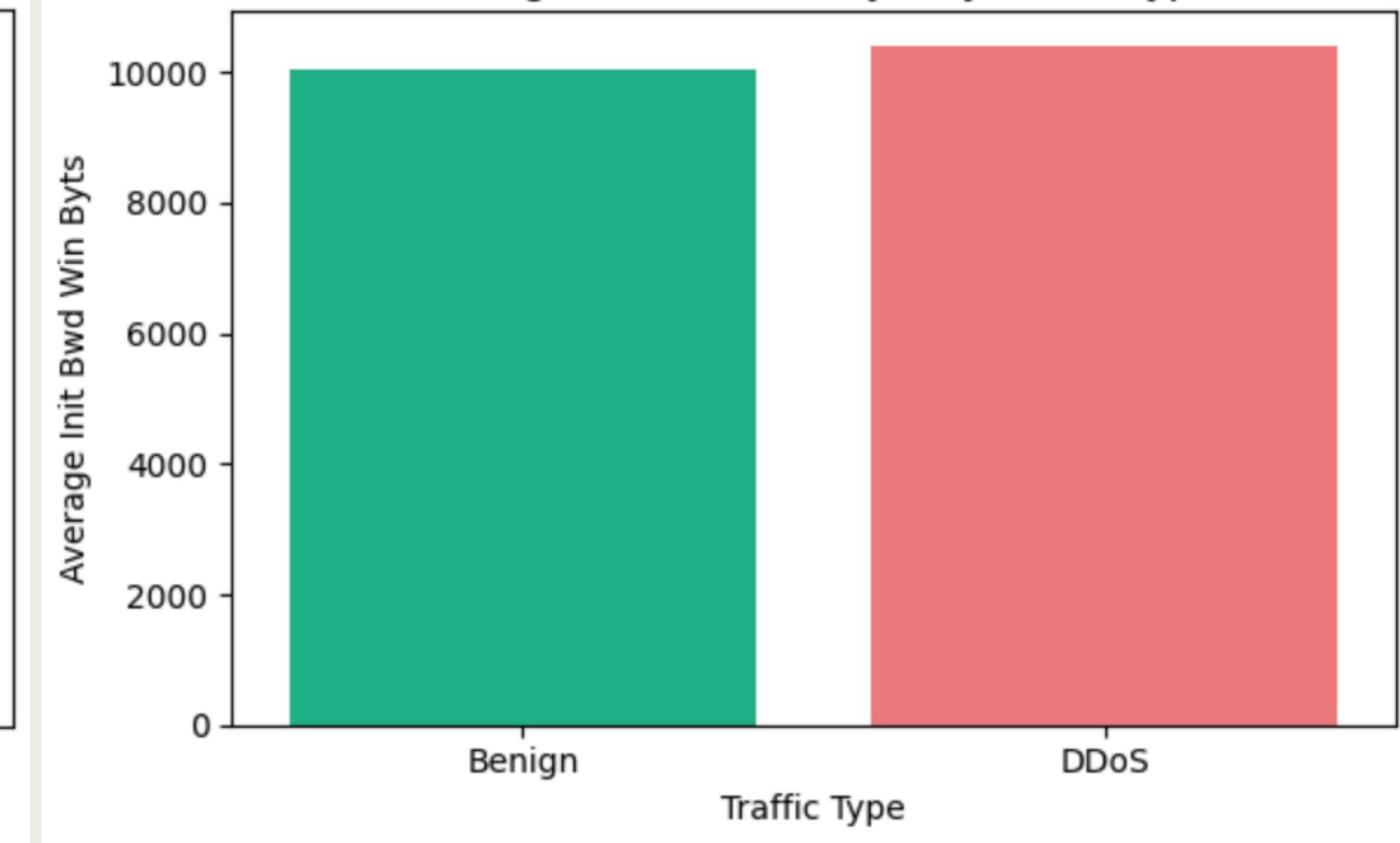
FORWARD SEGMENT SIZE

Average Fwd Seg Size Min by Traffic Type



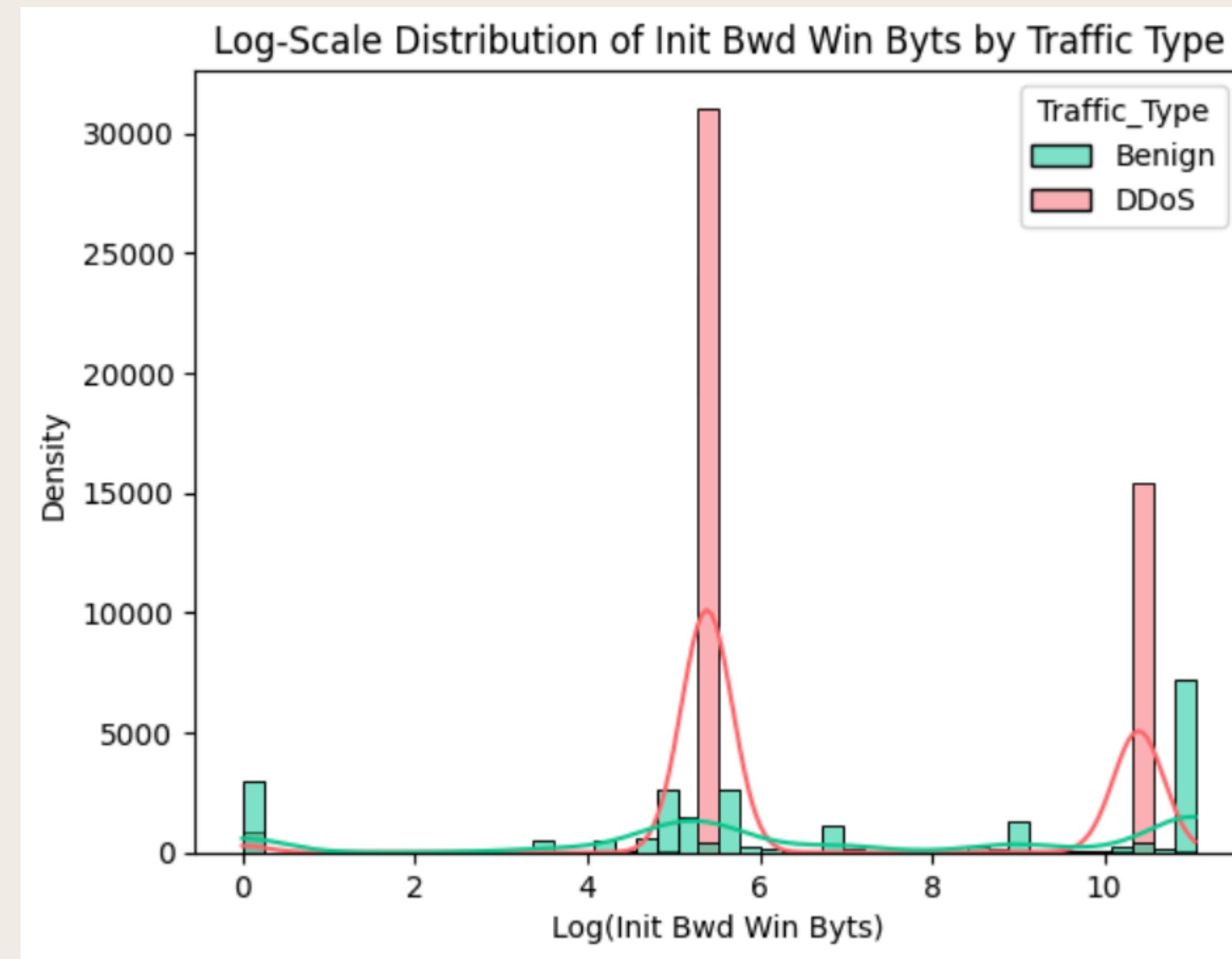
INITIAL FORWARD BYTES

Average Init Bwd Win Byts by Traffic Type

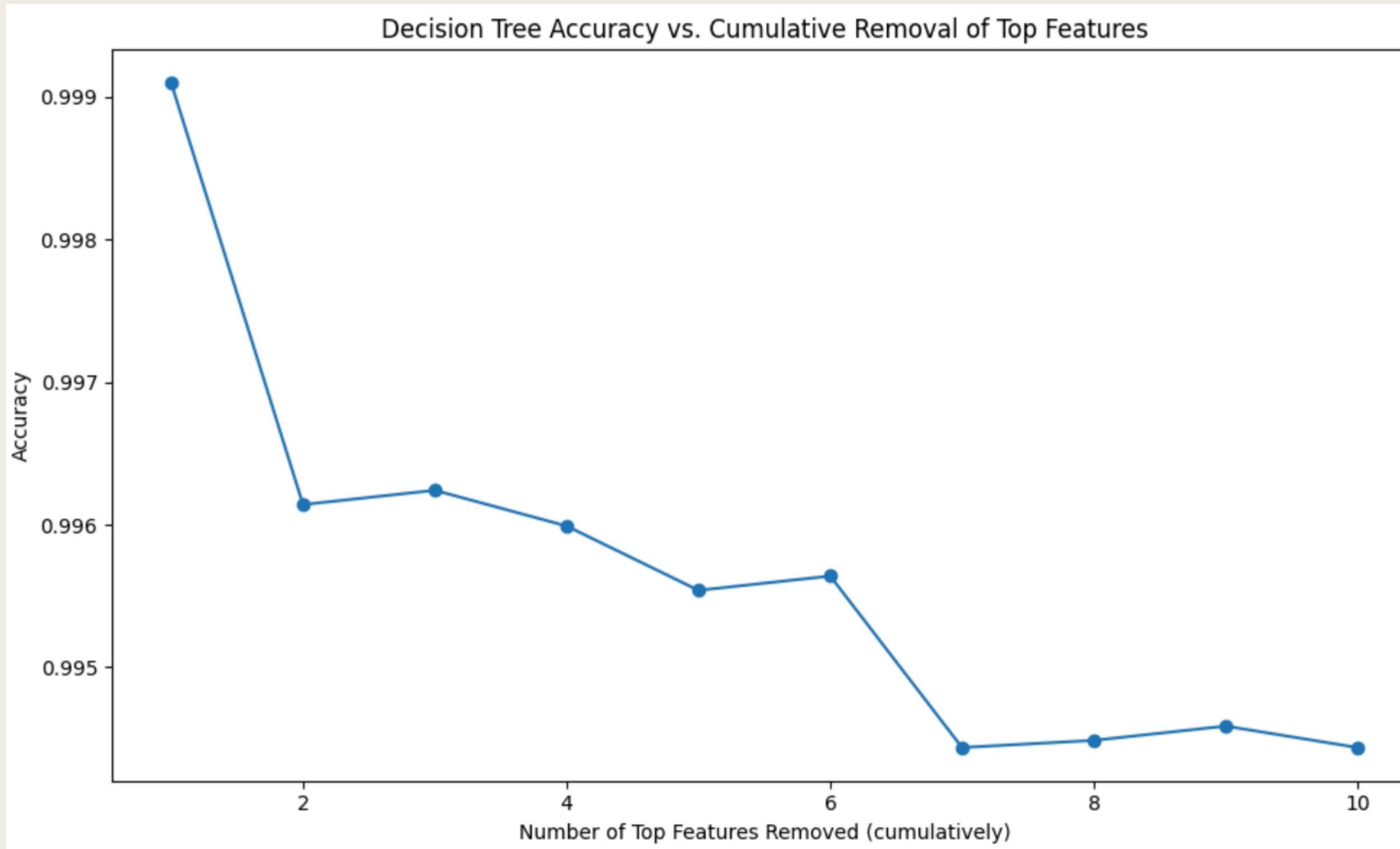


Forward Segment Size clearly separates normal traffic from DDoS attacks, making it a useful feature for detecting DDoS

INITIAL FORWARD BYTES



MOVING FEATURES



Final drop after
8+ removals but
Accuracy only
drops from 0.999
to 0.995

ONE FEATURE TRAINING

Features giving 80%+

- Packet Stats - 8
- Traffic Volume - 6
- Flow Summary - 5

Features giving 90%+

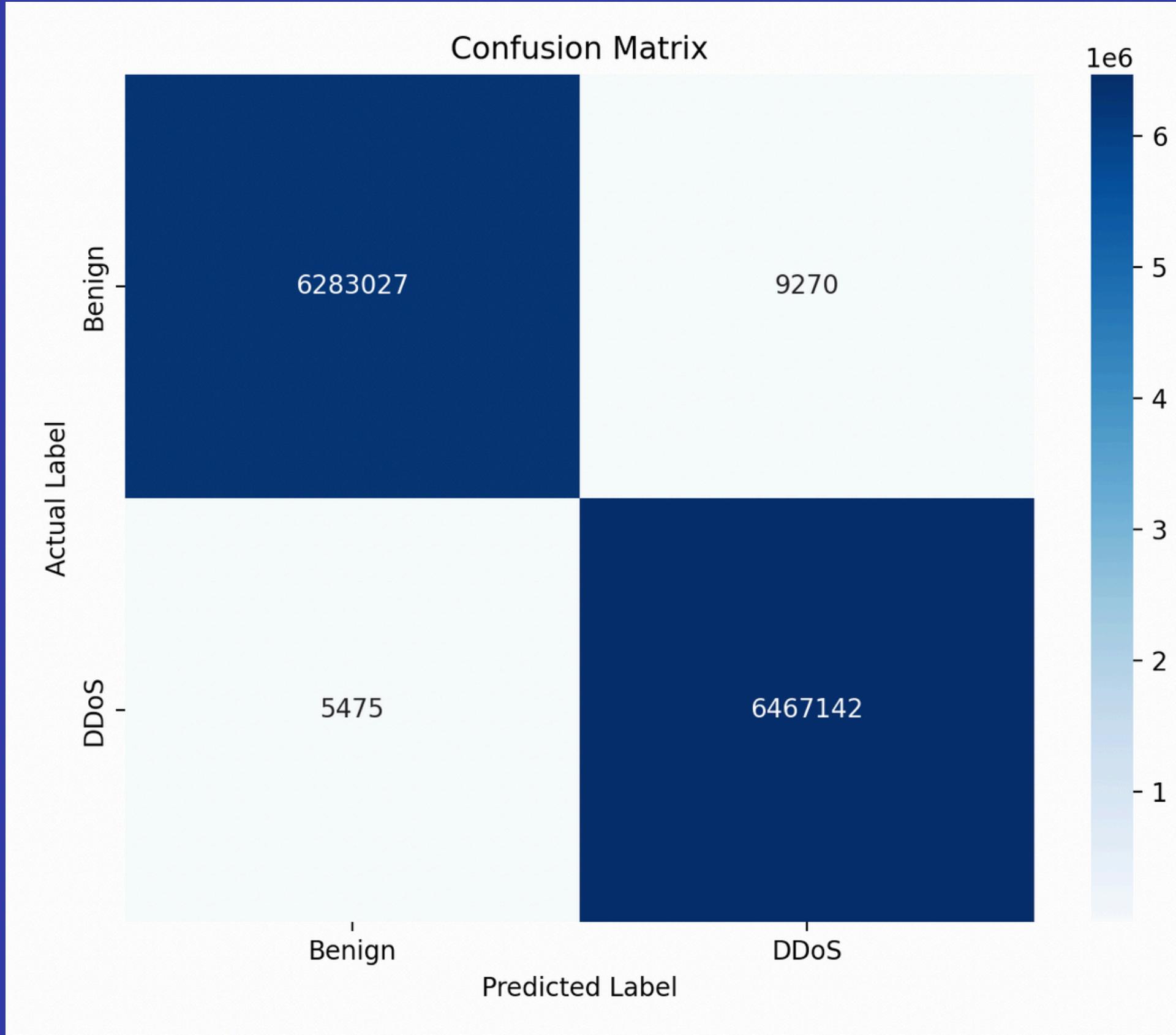
- Fwd Seg Size Min
- Initial fwd bytes
(the top two features)

SANITY CHECKS

Cross validation average: 99%

Accuracy after label shuffling: 50%

DECISION TREE



METRICS:

Accuracy: 99.88%

still better than all the other models!

Complete dataset
(13M rows)

CONCLUSION

Hard to work with DDoS attacks replicated in a controlled env
→ Not taking into account human factor

Size and **Info Distribution** are the most useful features for identifying a DDoS attack

DDoS and Benign flows have distinct characteristics, yet, it's difficult to stop in real time

THANK YOU VERY MUCH!

THOUGHTS, COMMENTS, CONCERNS?

