# Detecting Applications in Traffic Flows Using Machine Learning

**Authors**: Sam Xie (muyanxiesam), Jillian Chen (jianghui)

## Project Overview

With the increasing use of encryption in network traffic, traditional methods of application identification based on payload inspection have become less effective. This project aims to explore machine learning techniques to identify applications based on **traffic flow characteristics**, such as packet size distributions, timing patterns, and flow duration. Accurate application identification is crucial for network management, security monitoring, and quality of service enforcement.

### Flow Definition

A **network flow** is defined as a sequence of packets sharing the same 5-tuple:

- Source IP address
- Destination IP address
- Source port
- Destination port
- Protocol (TCP/UDP)

This project extracts flows from raw PCAP data and computes statistical features for each flow, enabling flow-level application classification.

This project will aim to reproduce results from the netml vpn non-vpn dataset paper

## Data

The project will utilize a dataset of network traffic captures in PCAPNG format, specifically focusing on encrypted traffic. The dataset will be processed to:

1. Extract individual packets with their metadata
2. Group packets into flows using the 5-tuple definition
3. Compute statistical flow features (packet counts, size distributions, timing statistics)
4. Aggregate packet-level labels to flow-level labels

A link to the dataset is here: Dataset Link.

### Data Processing

checking the required datasets

```
[1]: # detect if we have all the packages and data we need
     from scapy.all import rdpcap, wrpcap
     from pathlib import Path
     import os

     def fix_pcap_timestamps(input_file, output_file):
         """Remove first packet and normalize timestamps."""
```

```python
    print(f"Loading PCAP file...")
    packets = rdpcap(str(input_file))
    print(f"Original packet count: {len(packets)}")

    # Remove first packet
    packets_fixed = packets[1:]

    # Set timestamps to valid absolute values starting from a base time
    # Use a recent but fixed epoch time (e.g., Jan 1, 2020)
    base_time = 1577836800.0  # 2020-01-01 00:00:00 UTC

    if len(packets_fixed) > 0:
        print(f"Normalizing timestamps to start from 2020-01-01...")
        for i, pkt in enumerate(packets_fixed):
            # Set timestamp as base_time + packet index (1 second apart)
            # This ensures all timestamps are valid and sequential
            pkt.time = base_time + (i * 0.001)   # 1ms apart

    print(f"Fixed packet count: {len(packets_fixed)}")
    print(f"Writing to {output_file}...")
    wrpcap(str(output_file), packets_fixed)
    print(f"✓ Saved successfully!")

# File paths
original_pcap = Path("data/traffic.pcapng")
fixed_pcap = Path("data/traffic_fixed.pcap")
extracted_csv = Path("data/extracted.csv")

# Fix PCAP timestamps if needed
if not fixed_pcap.exists():
    print(f"Fixing timestamps in {original_pcap}...")
    if not original_pcap.exists():
        raise FileNotFoundError(f"Required file {original_pcap} not found.")
    fix_pcap_timestamps(original_pcap, fixed_pcap)
    print("✓ Timestamps fixed successfully!")
else:
    print(f"✓ Using existing fixed PCAP: {fixed_pcap}")

# Check if extracted.csv exists
if not extracted_csv.exists():
    print(f"\n⚠ Warning: {extracted_csv} not found.")
    print(f"Extracting comments from original PCAP...")
    os.system("tshark -r data/traffic.pcapng -T fields -E header=y -e frame.number
    print("✓ Comments extracted.")
else:
    print(f"✓ Found {extracted_csv}")

print(f"\n✓ All required files are ready!")
```

```
WARNING: No IPv4 address found on anpi3 !
WARNING: No IPv4 address found on anpi0 !
WARNING: No IPv4 address found on anpi0 !
WARNING: more No IPv4 address found on anpi1 !
WARNING: more No IPv4 address found on anpi1 !

✓ Using existing fixed PCAP: data/traffic_fixed.pcap
✓ Found data/extracted.csv

✓ All required files are ready!
```

# Flow Extraction and Feature Engineering

We extract flows from the PCAP file using the 5-tuple definition (src_ip, dst_ip, src_port, dst_port, proto-col). For each flow, we compute statistical features that capture the behavior of the traffic without requiring payload inspection.

```python
[2]: import pandas as pd
     import numpy as np
     from scapy.all import rdpcap, IP, TCP, UDP, IPv6

     # Load packet-level observations (labels)
     observations = pd.read_csv('data/extracted.csv',
                     skiprows=2,
                     names=['id', 'frame_number', 'frame_comment'])
     observations['id'] = observations['id'] - 2  # Adjust for removed first packet

     print(f"Loaded {len(observations)} packet labels")

     # Load packets from PCAP
     fixed_pcap = 'data/traffic_fixed.pcap'
     print(f"\nLoading packets from {fixed_pcap}...")
     packets = rdpcap(fixed_pcap)
     print(f"✓ Loaded {len(packets)} packets")

     # Extract packet information with flow keys
     print("\nExtracting packet information and flow keys...")
     packet_data = []

     for idx, pkt in enumerate(packets):
         # Extract flow 5-tuple
         src_ip = dst_ip = src_port = dst_port = protocol = None
         pkt_len = len(pkt)
         timestamp = float(pkt.time)

         # Handle IPv4
         if IP in pkt:
             src_ip = pkt[IP].src
             dst_ip = pkt[IP].dst
             protocol = pkt[IP].proto
         # Handle IPv6
         elif IPv6 in pkt:
             src_ip = pkt[IPv6].src
             dst_ip = pkt[IPv6].dst
             protocol = pkt[IPv6].nh

         # Extract ports for TCP/UDP
         if TCP in pkt:
             src_port = pkt[TCP].sport
             dst_port = pkt[TCP].dport
             protocol = 6   # TCP
         elif UDP in pkt:
             src_port = pkt[UDP].sport
             dst_port = pkt[UDP].dport
             protocol = 17   # UDP

         # Create bidirectional flow key (sorted to make it direction-agnostic)
         if src_ip and dst_ip:
             endpoints = sorted([(src_ip, src_port or 0), (dst_ip, dst_port or 0)])
             flow_key = f"{endpoints[0][0]}:{endpoints[0][1]}-{endpoints[1][0]}:{endpoi
         else:
             flow_key = f"unknown_{idx}"

         packet_data.append({
```

```python
            'packet_idx': idx,
            'flow_key': flow_key,
            'src_ip': src_ip,
            'dst_ip': dst_ip,
            'src_port': src_port,
            'dst_port': dst_port,
            'protocol': protocol,
            'pkt_len': pkt_len,
            'timestamp': timestamp
        })

    # Create packet DataFrame
    pdf = pd.DataFrame(packet_data)
    print(f"✓ Extracted {len(pdf)} packet records")
    print(f"✓ Found {pdf['flow_key'].nunique()} unique flows")
```

```
Loaded 529018 packet labels

Loading packets from data/traffic_fixed.pcap...
✓ Loaded 529018 packets

Extracting packet information and flow keys...
✓ Loaded 529018 packets

Extracting packet information and flow keys...
✓ Extracted 529018 packet records
✓ Found 111284 unique flows
✓ Extracted 529018 packet records
✓ Found 111284 unique flows
```

```python
[3]: # Merge packet data with labels
     pdf = pdf.merge(observations[['id', 'frame_comment']],
                     left_on='packet_idx', right_on='id', how='left')
     pdf = pdf.drop(columns=['id'])
     pdf = pdf.rename(columns={'frame_comment': 'label'})

     print(f"Packets with labels: {pdf['label'].notna().sum()}")
     print(f"\nSample packet data:")
     pdf.head()
```

```
Packets with labels: 529018

Sample packet data:
```

[3]:

| | packet_idx | flow_key | src_ip | dst_ip | src_port | dst_port | protocol | pkt_len | timestam |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 131.202.240.87:8118-222.186.31.178:49445-6 | 222.186.31.178 | 131.202.240.87 | 49445 | 8118 | 6 | 60 | 1.577837e+ |
| 1 | 1 | 131.202.240.87:4506-207.241.228.202:80-6 | 131.202.240.87 | 207.241.228.202 | 4506 | 80 | 6 | 66 | 1.577837e+ |
| 2 | 2 | 131.202.240.87:4506-207.241.228.202:80-6 | 207.241.228.202 | 131.202.240.87 | 80 | 4506 | 6 | 66 | 1.577837e+ |
| 3 | 3 | 131.202.240.87:4506-207.241.228.202:80-6 | 131.202.240.87 | 207.241.228.202 | 4506 | 80 | 6 | 54 | 1.577837e+ |
| 4 | 4 | 131.202.240.87:4506-207.241.228.202:80-6 | 131.202.240.87 | 207.241.228.202 | 4506 | 80 | 6 | 240 | 1.577837e+ |

```python
[4]: # Check flow distribution
     print("Flow size distribution:")
     flow_sizes = pdf.groupby('flow_key').size()
     print(f"  Min packets per flow: {flow_sizes.min()}")
     print(f"  Max packets per flow: {flow_sizes.max()}")
     print(f"  Mean packets per flow: {flow_sizes.mean():.2f}")
     print(f"  Median packets per flow: {flow_sizes.median():.2f}")
```

```
print(f"\nLabel distribution in packets:")
print(pdf['label'].value_counts().head(10))
```

```
Flow size distribution:
  Min packets per flow: 1
  Max packets per flow: 875
  Mean packets per flow: 4.75
  Median packets per flow: 2.00

Label distribution in packets:
label
audio_hangouts_hangouts-audio    126174
audio_facebook_facebook-audio    122003
file-transfer_skype_skype-file    85035
audio_skype_skype-audio           58558
email_email_email                 18876
p2p_torrent_torrent               14926
video_skype_skype-video           13539
chat_skype_skype-chat             12542
video_youtube_youtube             11571
audio_voipbuster_voipbuster       11355
Name: count, dtype: int64
label
audio_hangouts_hangouts-audio    126174
audio_facebook_facebook-audio    122003
file-transfer_skype_skype-file    85035
audio_skype_skype-audio           58558
email_email_email                 18876
p2p_torrent_torrent               14926
video_skype_skype-video           13539
chat_skype_skype-chat             12542
video_youtube_youtube             11571
audio_voipbuster_voipbuster       11355
Name: count, dtype: int64
```

# Flow Feature Engineering

For each flow, we compute statistical features that capture traffic patterns without inspecting payloads:

**Packet Count Features:**

- Total packets in flow
- Forward/backward packet counts

**Packet Size Features:**

- Total bytes, mean, std, min, max packet sizes
- Forward/backward size statistics

**Timing Features:**

- Flow duration
- Inter-arrival time statistics (mean, std, min, max)

**Protocol Features:**

- Protocol type (TCP/UDP)
- Port numbers

```
[5]: def extract_flow_features(flow_df):
         """Extract statistical features for a single flow."""
         features = {}
```

```python
    # Basic flow info
    features['flow_key'] = flow_df['flow_key'].iloc[0]

    # Packet count features
    features['total_packets'] = len(flow_df)
    features['total_bytes'] = flow_df['pkt_len'].sum()

    # Packet size statistics
    features['pkt_len_mean'] = flow_df['pkt_len'].mean()
    features['pkt_len_std'] = flow_df['pkt_len'].std() if len(flow_df) > 1 else 0
    features['pkt_len_min'] = flow_df['pkt_len'].min()
    features['pkt_len_max'] = flow_df['pkt_len'].max()
    features['pkt_len_median'] = flow_df['pkt_len'].median()

    # Timing features
    timestamps = flow_df['timestamp'].sort_values()
    features['flow_duration'] = timestamps.max() - timestamps.min()

    if len(timestamps) > 1:
        iat = timestamps.diff().dropna()  # Inter-arrival times
        features['iat_mean'] = iat.mean()
        features['iat_std'] = iat.std()
        features['iat_min'] = iat.min()
        features['iat_max'] = iat.max()
    else:
        features['iat_mean'] = 0
        features['iat_std'] = 0
        features['iat_min'] = 0
        features['iat_max'] = 0

    # Protocol features
    features['protocol'] = flow_df['protocol'].iloc[0] if flow_df['protocol'].notn

    # Port features (use min port as it's often the server port)
    src_ports = flow_df['src_port'].dropna()
    dst_ports = flow_df['dst_port'].dropna()
    all_ports = pd.concat([src_ports, dst_ports])
    features['min_port'] = all_ports.min() if len(all_ports) > 0 else 0
    features['max_port'] = all_ports.max() if len(all_ports) > 0 else 0

    # Packet size distribution features
    features['small_packets'] = (flow_df['pkt_len'] < 100).sum()  # Small packets
    features['medium_packets'] = ((flow_df['pkt_len'] >= 100) & (flow_df['pkt_len'
    features['large_packets'] = (flow_df['pkt_len'] >= 500).sum()

    # Bytes per second (flow rate)
    if features['flow_duration'] > 0:
        features['bytes_per_sec'] = features['total_bytes'] / features['flow_durat
        features['packets_per_sec'] = features['total_packets'] / features['flow_d
    else:
        features['bytes_per_sec'] = features['total_bytes']
        features['packets_per_sec'] = features['total_packets']

    # Get the most common label for this flow (majority voting)
    labels = flow_df['label'].dropna()
    if len(labels) > 0:
        features['label'] = labels.mode().iloc[0]
    else:
        features['label'] = None

    return features
```

```
print("Flow feature extraction function defined ✓")
```

Flow feature extraction function defined ✓

```
[6]: # Extract features for all flows
     print("Extracting flow features...")
     from tqdm import tqdm

     flow_features = []
     grouped = pdf.groupby('flow_key')
     total_flows = len(grouped)

     for flow_key, flow_df in tqdm(grouped, total=total_flows, desc="Processing flows")
         features = extract_flow_features(flow_df)
         flow_features.append(features)

     # Create flow DataFrame
     flow_df = pd.DataFrame(flow_features)

     # Remove flows without labels
     flow_df = flow_df.dropna(subset=['label'])

     print(f"\n✓ Extracted {len(flow_df)} labeled flows")
     print(f"✓ {len(flow_df.columns) - 2} features per flow")  # -2 for flow_key and la
     print(f"\nFlow features:")
     print(list(flow_df.columns))
```

Extracting flow features...

Processing flows: 100%|████████████| 111284/111284 [01:04<00:00, 1714.58it/s]

✓ Extracted 111284 labeled flows
✓ 20 features per flow

Flow features:
['flow_key', 'total_packets', 'total_bytes', 'pkt_len_mean', 'pkt_len_std', 'pkt_le

## Train and Test Split

We split the flow dataset into training and testing sets (80/20 split), ensuring that the split maintains the distribution of application classes through stratification.

```
[7]: from sklearn.model_selection import train_test_split

     # Preview the flow data
     print("Flow data summary:")
     print(flow_df.describe())
     print(f"\nLabel distribution:")
     print(flow_df['label'].value_counts())

     # Separate features and labels
     feature_cols = [col for col in flow_df.columns if col not in ['flow_key', 'label']
     X = flow_df[feature_cols]
     y = flow_df['label']

     print(f"\n✓ Feature matrix shape: {X.shape}")
     print(f"✓ Number of classes: {y.nunique()}")
```

Flow data summary:
       total_packets    total_bytes    pkt_len_mean    pkt_len_std  \

```
count  111284.000000  111284.000000   111284.000000  111284.000000
mean        4.753765     793.260972       96.330745      27.205372
std        14.437535    5706.271733       77.464731      85.395472
min         1.000000      54.000000       49.000000       0.000000
25%         2.000000     128.000000       64.000000       0.000000
50%         2.000000     256.000000       64.000000       0.000000
75%         4.000000     384.000000      129.500000       0.000000
max       875.000000  399301.000000     1196.460000     758.060903

             pkt_len_min     pkt_len_max   pkt_len_median   flow_duration  \
count      111284.000000   111284.000000    111284.000000   111284.000000
mean           76.211342      129.843715        93.802429       83.561406
std            27.928165      198.008331        89.163239      140.379092
min            44.000000       54.000000        44.000000        0.000000
25%            64.000000       64.000000        64.000000        0.001000
50%            64.000000       64.000000        64.000000        0.001000
75%            73.000000      131.000000       110.000000      149.557750
max           816.000000     1500.000000      1500.000000      512.397000

                 iat_mean        iat_std          iat_min         iat_max  \
count       111284.000000   5.625700e+04    111284.000000   111284.000000
mean            23.328593   6.441239e+01         5.511464       72.004099
std             49.332916   7.004665e+01        38.588282      120.810929
min              0.000000   0.000000e+00         0.000000        0.000000
25%              0.001000   1.101272e-07         0.001000        0.001000
50%              0.001000   5.370946e+01         0.001000        0.001000
75%             33.892500   1.102878e+02         0.001000      137.765000
max            498.804000   3.581694e+02       498.804000      511.023000

                 protocol        min_port         max_port   small_packets  \
count       111284.000000   111284.000000    111284.000000   111284.000000
mean            16.588306     5981.056100     54292.204935        3.661119
std              2.087866     5562.934703     11625.858899        9.788265
min              6.000000       19.000000        68.000000        0.000000
25%             17.000000     5355.000000     51991.000000        2.000000
50%             17.000000     5355.000000     56597.000000        2.000000
75%             17.000000     5355.000000     61009.000000        4.000000
max             17.000000    61009.000000     65535.000000      863.000000

           medium_packets   large_packets   bytes_per_sec   packets_per_sec
count       111284.000000   111284.000000    1.112840e+05     111284.000000
mean             0.779474        0.313172    1.056969e+05       1069.973603
std              7.768443        4.004202    1.373792e+05        948.167619
min              0.000000        0.000000    4.368990e-01          0.004010
25%              0.000000        0.000000    2.374907e+00          0.035798
50%              0.000000        0.000000    1.279788e+05       1333.324009
75%              1.000000        0.000000    1.280093e+05       2000.144969
max            504.000000      340.000000    1.632118e+06       2000.144969

Label distribution:
label
audio_facebook_facebook-audio     38940
audio_hangouts_hangouts-audio     33038
file-transfer_skype_skype-file    14127
audio_skype_skype-audio           10539
audio_voipbuster_voipbuster        2660
email_email_email                  2114
chat_skype_skype-chat              1893
video_hangouts_hangouts-video      1469
p2p_torrent_torrent                 958
video_youtube_youtube               826
file-transfer_ftps_ftps-down        586
video_skype_skype-video             512
chat_facebook_facebook-chat         453
chat_gmail_gmail-chat               429
video_vimeo_vimeo                   399
video_facebook_facebook-video       388
chat_icq_icq-chat                   340
chat_aim_aim-chat                   339
```

```
chat_hangouts_hangouts-chat          298
video_netflix_netflix                247
audio_spotify_spotify                198
file-transfer_ftps_ftps-up           134
tor_youtube_tor-youtube               87
file-transfer_scp_scp-down            84
file-transfer_scp_scp-up              82
file-transfer_sftp_sftp-down          72
file-transfer_sftp_sftp-up            63
video_vimeo_tor-vimeo                  6
tor_facebook_tor-facebook              2
tor_twitter_tor-twitter                1
Name: count, dtype: int64

✓ Feature matrix shape: (111284, 20)
✓ Number of classes: 30
```

[8]:
```python
# Filter classes with too few samples for stratification
min_samples = 2
class_counts = y.value_counts()
valid_classes = class_counts[class_counts >= min_samples].index
mask = y.isin(valid_classes)
X_filtered = X[mask]
y_filtered = y[mask]

print(f"Classes with >= {min_samples} samples: {len(valid_classes)}")
print(f"Flows after filtering: {len(X_filtered)}")

# Perform 80/20 split with stratification
X_train, X_test, y_train, y_test = train_test_split(
    X_filtered, y_filtered,
    test_size=0.2,
    random_state=42,
    stratify=y_filtered
)

print(f"\nTraining set size: {len(X_train)} flows")
print(f"Test set size: {len(X_test)} flows")
print(f"\nTraining set class distribution:")
print(y_train.value_counts().head(10))
print(f"\nTest set class distribution:")
print(y_test.value_counts().head(10))
```

```
Classes with >= 2 samples: 29
Flows after filtering: 111283

Training set size: 89026 flows
Test set size: 22257 flows

Training set class distribution:
label
audio_facebook_facebook-audio      31152
audio_hangouts_hangouts-audio      26430
file-transfer_skype_skype-file     11302
audio_skype_skype-audio             8431
audio_voipbuster_voipbuster         2128
email_email_email                   1691
chat_skype_skype-chat               1514
video_hangouts_hangouts-video       1175
p2p_torrent_torrent                  766
video_youtube_youtube                661
Name: count, dtype: int64

Test set class distribution:
label
audio_facebook_facebook-audio       7788
```

```
audio_hangouts_hangouts-audio      6608
file-transfer_skype_skype-file     2825
audio_skype_skype-audio            2108
audio_voipbuster_voipbuster         532
email_email_email                   423
chat_skype_skype-chat               379
video_hangouts_hangouts-video       294
p2p_torrent_torrent                 192
video_youtube_youtube               165
Name: count, dtype: int64
```

# Model Training

We experiment with several machine learning algorithms using flow-level features:

1. **Random Forest** - Ensemble method, good for feature importance analysis
2. **Support Vector Machines (SVM)** - Effective for high-dimensional data
3. **Neural Networks (MLP)** - Can learn complex non-linear patterns

## 1. Random Forest

```
[9]: # Prepare flow features for training
     print("Flow features being used:")
     print(X_train.columns.tolist())
     print(f"\nFeature data types:")
     print(X_train.dtypes)

     # Check for any missing values
     print(f"\nMissing values per feature:")
     print(X_train.isnull().sum())

     # Fill any remaining NaN values with 0
     X_train_processed = X_train.fillna(0)
     X_test_processed = X_test.fillna(0)

     print(f"\n✓ Training features ready: {X_train_processed.shape}")
```

```
Flow features being used:
['total_packets', 'total_bytes', 'pkt_len_mean', 'pkt_len_std', 'pkt_len_min', 'pkt

Feature data types:
total_packets        int64
total_bytes          int64
pkt_len_mean         float64
pkt_len_std          float64
pkt_len_min          int64
pkt_len_max          int64
pkt_len_median       float64
flow_duration        float64
iat_mean             float64
iat_std              float64
iat_min              float64
iat_max              float64
protocol             int64
min_port             int64
max_port             int64
small_packets        int64
medium_packets       int64
large_packets        int64
bytes_per_sec        float64
packets_per_sec      float64
dtype: object
```

```
Missing values per feature:
total_packets            0
total_bytes              0
pkt_len_mean             0
pkt_len_std              0
pkt_len_min              0
pkt_len_max              0
pkt_len_median           0
flow_duration            0
iat_mean                 0
iat_std              43941
iat_min                  0
iat_max                  0
protocol                 0
min_port                 0
max_port                 0
small_packets            0
medium_packets           0
large_packets            0
bytes_per_sec            0
packets_per_sec          0
dtype: int64

✓ Training features ready: (89026, 20)
```

```python
[10]: from sklearn.ensemble import RandomForestClassifier
      from sklearn.metrics import accuracy_score, classification_report
      import time

      # Initialize Random Forest Classifier
      rf_model = RandomForestClassifier(
          n_estimators=100,        # Number of trees in the forest
          max_depth=20,            # Maximum depth of trees
          min_samples_split=5,     # Minimum samples required to split a node
          min_samples_leaf=2,      # Minimum samples required at leaf node
          random_state=42,         # For reproducibility
          n_jobs=-1                # Use all available processors
      )

      print("Training Random Forest model on flow features...")
      start_time = time.time()
      rf_model.fit(X_train_processed, y_train)
      rf_train_time = time.time() - start_time
      print(f"✓ Training complete in {rf_train_time:.2f} seconds!")

      # Make predictions
      print("\nMaking predictions on test set...")
      y_pred_rf = rf_model.predict(X_test_processed)

      # Evaluate the model
      print("\n" + "="*60)
      print("RANDOM FOREST MODEL EVALUATION (Flow-Level)")
      print("="*60)

      # Accuracy
      accuracy = accuracy_score(y_test, y_pred_rf)
      print(f"\nAccuracy: {accuracy:.4f} ({accuracy*100:.2f}%)")

      # Calculate feature importance
      feature_importance = pd.DataFrame({
          'feature': X_train_processed.columns,
          'importance': rf_model.feature_importances_
      }).sort_values('importance', ascending=False)
```

```
print(f"\nTop 10 Most Important Flow Features:")
print(feature_importance.head(10))
```

Training Random Forest model on flow features...
✓ Training complete in 1.41 seconds!

Making predictions on test set...

============================================================
RANDOM FOREST MODEL EVALUATION (Flow-Level)
============================================================

Accuracy: 0.5039 (50.39%)

Top 10 Most Important Flow Features:
            feature   importance
14         max_port     0.175336
18    bytes_per_sec     0.103569
11          iat_max     0.092472
7     flow_duration     0.070209
19  packets_per_sec     0.062296
9           iat_std     0.057306
3       pkt_len_std     0.054664
6    pkt_len_median     0.052264
8          iat_mean     0.052191
1       total_bytes     0.048961
✓ Training complete in 1.41 seconds!

Making predictions on test set...

============================================================
RANDOM FOREST MODEL EVALUATION (Flow-Level)
============================================================

Accuracy: 0.5039 (50.39%)

Top 10 Most Important Flow Features:
            feature   importance
14         max_port     0.175336
18    bytes_per_sec     0.103569
11          iat_max     0.092472
7     flow_duration     0.070209
19  packets_per_sec     0.062296
9           iat_std     0.057306
3       pkt_len_std     0.054664
6    pkt_len_median     0.052264
8          iat_mean     0.052191
1       total_bytes     0.048961
```

## 2. Support Vector Machines (SVM)

```
[11]: from sklearn.svm import LinearSVC
      from sklearn.preprocessing import StandardScaler
      import warnings

      warnings.filterwarnings('ignore')

      # Scale features for SVM (SVMs are sensitive to feature scaling)
      print("Scaling flow features for SVM...")
      scaler = StandardScaler()
      X_train_scaled = scaler.fit_transform(X_train_processed)
      X_test_scaled = scaler.transform(X_test_processed)
      print("✓ Scaling complete!")

      # Initialize Linear SVM
```

```python
svm_model = LinearSVC(
    C=1.0,                      # Regularization parameter
    max_iter=2000,              # Maximum iterations
    random_state=42,
    dual=True
)

print("\nTraining Linear SVM model on flow features...")
start_time = time.time()
svm_model.fit(X_train_scaled, y_train)
svm_train_time = time.time() - start_time
print(f"✓ Training complete in {svm_train_time:.2f} seconds!")

# Make predictions
print("\nMaking predictions on test set...")
y_pred_svm = svm_model.predict(X_test_scaled)

# Evaluate
print("\n" + "="*60)
print("SUPPORT VECTOR MACHINE (SVM) MODEL EVALUATION (Flow-Level)")
print("="*60)

accuracy_svm = accuracy_score(y_test, y_pred_svm)
print(f"\nAccuracy: {accuracy_svm:.4f} ({accuracy_svm*100:.2f}%)")
```

```
Scaling flow features for SVM...
✓ Scaling complete!

Training Linear SVM model on flow features...
✓ Training complete in 314.05 seconds!

Making predictions on test set...

============================================================
SUPPORT VECTOR MACHINE (SVM) MODEL EVALUATION (Flow-Level)
============================================================

Accuracy: 0.4198 (41.98%)
✓ Training complete in 314.05 seconds!

Making predictions on test set...

============================================================
SUPPORT VECTOR MACHINE (SVM) MODEL EVALUATION (Flow-Level)
============================================================

Accuracy: 0.4198 (41.98%)
```

## 3. Neural Networks

```python
[12]: from sklearn.neural_network import MLPClassifier

# Initialize Multi-Layer Perceptron (Neural Network)
nn_model = MLPClassifier(
    hidden_layer_sizes=(128, 64, 32),  # Three hidden layers
    activation='relu',
    solver='adam',
    alpha=0.0001,                       # L2 regularization
    batch_size=256,
    learning_rate='adaptive',
    learning_rate_init=0.001,
    max_iter=300,
    random_state=42,
```

```
        verbose=True,
        early_stopping=True,
        validation_fraction=0.1,
        n_iter_no_change=15
)

print("Training Neural Network model on flow features...")
print("Architecture: Input -> 128 -> 64 -> 32 -> Output\n")
start_time = time.time()
nn_model.fit(X_train_scaled, y_train)
nn_train_time = time.time() - start_time
print(f"\n✓ Training complete in {nn_train_time:.2f} seconds!")
print(f"Training stopped at iteration: {nn_model.n_iter_}")

# Make predictions
y_pred_nn = nn_model.predict(X_test_scaled)

# Evaluate
print("\n" + "="*60)
print("NEURAL NETWORK MODEL EVALUATION (Flow-Level)")
print("="*60)

accuracy_nn = accuracy_score(y_test, y_pred_nn)
print(f"\nAccuracy: {accuracy_nn:.4f} ({accuracy_nn*100:.2f}%)")
```

```
Training Neural Network model on flow features...
Architecture: Input -> 128 -> 64 -> 32 -> Output

Iteration 1, loss = 1.72302430
Validation score: 0.437830
Iteration 2, loss = 1.49745900
Validation score: 0.438504
Iteration 1, loss = 1.72302430
Validation score: 0.437830
Iteration 2, loss = 1.49745900
Validation score: 0.438504
Iteration 3, loss = 1.46355719
Validation score: 0.435584
Iteration 4, loss = 1.44245657
Validation score: 0.447153
Iteration 3, loss = 1.46355719
Validation score: 0.435584
Iteration 4, loss = 1.44245657
Validation score: 0.447153
Iteration 5, loss = 1.42703073
Validation score: 0.445468
Iteration 6, loss = 1.41189643
Validation score: 0.441649
Iteration 5, loss = 1.42703073
Validation score: 0.445468
Iteration 6, loss = 1.41189643
Validation score: 0.441649
Iteration 7, loss = 1.40081894
Validation score: 0.450073
Iteration 8, loss = 1.39218526
Validation score: 0.458160
Iteration 7, loss = 1.40081894
Validation score: 0.450073
Iteration 8, loss = 1.39218526
Validation score: 0.458160
Iteration 9, loss = 1.38299724
Validation score: 0.458722
Iteration 10, loss = 1.37539509
Validation score: 0.441874
Iteration 9, loss = 1.38299724
Validation score: 0.458722
```

```
Iteration 10, loss = 1.37539509
Validation score: 0.441874
Iteration 11, loss = 1.36837243
Validation score: 0.450635
Iteration 12, loss = 1.36266357
Validation score: 0.461417
Iteration 11, loss = 1.36837243
Validation score: 0.450635
Iteration 12, loss = 1.36266357
Validation score: 0.461417
Iteration 13, loss = 1.35764530
Validation score: 0.451309
Iteration 14, loss = 1.35227551
Validation score: 0.451870
Iteration 13, loss = 1.35764530
Validation score: 0.451309
Iteration 14, loss = 1.35227551
Validation score: 0.451870
Iteration 15, loss = 1.34790847
Validation score: 0.463552
Iteration 16, loss = 1.34188027
Validation score: 0.456700
Iteration 15, loss = 1.34790847
Validation score: 0.463552
Iteration 16, loss = 1.34188027
Validation score: 0.456700
Iteration 17, loss = 1.33965273
Validation score: 0.470179
Iteration 18, loss = 1.33717135
Validation score: 0.462990
Iteration 17, loss = 1.33965273
Validation score: 0.470179
Iteration 18, loss = 1.33717135
Validation score: 0.462990
Iteration 19, loss = 1.33366985
Validation score: 0.473885
Iteration 20, loss = 1.33063796
Validation score: 0.458946
Iteration 19, loss = 1.33366985
Validation score: 0.473885
Iteration 20, loss = 1.33063796
Validation score: 0.458946
Iteration 21, loss = 1.32808252
Validation score: 0.475458
Iteration 22, loss = 1.32449754
Validation score: 0.473324
Iteration 21, loss = 1.32808252
Validation score: 0.475458
Iteration 22, loss = 1.32449754
Validation score: 0.473324
Iteration 23, loss = 1.32405017
Validation score: 0.458834
Iteration 24, loss = 1.32175343
Validation score: 0.467371
Iteration 23, loss = 1.32405017
Validation score: 0.458834
Iteration 24, loss = 1.32175343
Validation score: 0.467371
Iteration 25, loss = 1.31924153
Validation score: 0.467034
Iteration 26, loss = 1.31772664
Validation score: 0.470628
Iteration 25, loss = 1.31924153
Validation score: 0.467034
Iteration 26, loss = 1.31772664
Validation score: 0.470628
Iteration 27, loss = 1.31647045
Validation score: 0.471976
Iteration 28, loss = 1.31436852
```

```
Validation score: 0.467146
Iteration 27, loss = 1.31647045
Validation score: 0.471976
Iteration 28, loss = 1.31436852
Validation score: 0.467146
Iteration 29, loss = 1.31278806
Validation score: 0.467146
Iteration 30, loss = 1.31322217
Validation score: 0.463215
Iteration 29, loss = 1.31278806
Validation score: 0.467146
Iteration 30, loss = 1.31322217
Validation score: 0.463215
Iteration 31, loss = 1.30975313
Validation score: 0.465686
Iteration 32, loss = 1.30883224
Validation score: 0.469954
Iteration 31, loss = 1.30975313
Validation score: 0.465686
Iteration 32, loss = 1.30883224
Validation score: 0.469954
Iteration 33, loss = 1.30800027
Validation score: 0.474334
Iteration 34, loss = 1.30708656
Validation score: 0.470179
Iteration 33, loss = 1.30800027
Validation score: 0.474334
Iteration 34, loss = 1.30708656
Validation score: 0.470179
Iteration 35, loss = 1.30669133
Validation score: 0.474222
Iteration 36, loss = 1.30386370
Validation score: 0.473436
Iteration 35, loss = 1.30669133
Validation score: 0.474222
Iteration 36, loss = 1.30386370
Validation score: 0.473436
Iteration 37, loss = 1.30225164
Validation score: 0.459957
Validation score did not improve more than tol=0.000100 for 15 consecutive epochs.

✓ Training complete in 5.70 seconds!
Training stopped at iteration: 37


============================================================
NEURAL NETWORK MODEL EVALUATION (Flow-Level)
============================================================

Accuracy: 0.4646 (46.46%)
Iteration 37, loss = 1.30225164
Validation score: 0.459957
Validation score did not improve more than tol=0.000100 for 15 consecutive epochs.

✓ Training complete in 5.70 seconds!
Training stopped at iteration: 37


============================================================
NEURAL NETWORK MODEL EVALUATION (Flow-Level)
============================================================

Accuracy: 0.4646 (46.46%)
```

```
[13]: # Compare all models
      print("\n" + "="*60)
      print("FINAL MODEL COMPARISON (Flow-Level Classification)")
      print("="*60)
      print(f"\nDataset Statistics:")
      print(f"  Total flows: {len(flow_df)}")
```

```
print(f"  Training flows: {len(X_train)}")
print(f"  Test flows: {len(X_test)}")
print(f"  Features per flow: {X_train_processed.shape[1]}")
print(f"  Number of classes: {y.nunique()}")

print(f"\nModel Performance:")
print(f"  Random Forest Accuracy:  {accuracy:.4f} ({accuracy*100:.2f}%) - trained
print(f"  SVM Accuracy:            {accuracy_svm:.4f} ({accuracy_svm*100:.2f}%) -
print(f"  Neural Network Accuracy: {accuracy_nn:.4f} ({accuracy_nn*100:.2f}%) - tr

print("\nBest Model: ", end="")
best_accuracy = max(accuracy, accuracy_svm, accuracy_nn)
if best_accuracy == accuracy:
    print(f"Random Forest ({accuracy*100:.2f}%)")
elif best_accuracy == accuracy_svm:
    print(f"SVM ({accuracy_svm*100:.2f}%)")
else:
    print(f"Neural Network ({accuracy_nn*100:.2f}%)")
```

```
============================================================
FINAL MODEL COMPARISON (Flow-Level Classification)
============================================================

Dataset Statistics:
  Total flows: 111284
  Training flows: 89026
  Test flows: 22257
  Features per flow: 20
  Number of classes: 30

Model Performance:
  Random Forest Accuracy:  0.5039 (50.39%) - trained in 1.41s
  SVM Accuracy:            0.4198 (41.98%) - trained in 314.05s
  Neural Network Accuracy: 0.4646 (46.46%) - trained in 5.70s

Best Model: Random Forest (50.39%)
```

# Results

## Visualization

Comprehensive visualization of flow-based classification results including model performance, feature importance, and flow statistics.

```
[14]: import matplotlib.pyplot as plt
      import seaborn as sns

      # Set style
      sns.set_style("whitegrid")
      sns.set_palette("husl")

      fig = plt.figure(figsize=(18, 14))
      gs = fig.add_gridspec(3, 3, hspace=0.4, wspace=0.35)

      # 1. Model Accuracy Comparison
      ax1 = fig.add_subplot(gs[0, :2])
      models = ['Random Forest', 'SVM', 'Neural Network']
      accuracies = [accuracy * 100, accuracy_svm * 100, accuracy_nn * 100]
      colors = ['#2ecc71', '#e74c3c', '#3498db']
```

```python
bars = ax1.bar(models, accuracies, color=colors, alpha=0.8, edgecolor='black', lin
ax1.set_ylabel('Accuracy (%)', fontsize=13, fontweight='bold')
ax1.set_title('Flow-Level Model Performance Comparison', fontsize=15, fontweight='
ax1.set_ylim([0, 100])
ax1.grid(axis='y', alpha=0.3, linestyle='--')

for bar, acc in zip(bars, accuracies):
    ax1.text(bar.get_x() + bar.get_width()/2., bar.get_height() + 1,
             f'{acc:.2f}%', ha='center', va='bottom', fontweight='bold', fontsize=1

# 2. Training Time Comparison
ax2 = fig.add_subplot(gs[0, 2])
training_times = [rf_train_time, svm_train_time, nn_train_time]
ax2.barh(models, training_times, color=colors, alpha=0.8, edgecolor='black', linew
ax2.set_xlabel('Time (seconds)', fontsize=11, fontweight='bold')
ax2.set_title('Training Time', fontsize=13, fontweight='bold', pad=10)
ax2.grid(axis='x', alpha=0.3, linestyle='--')
for i, t in enumerate(training_times):
    ax2.text(t + 0.5, i, f'{t:.1f}s', va='center', fontweight='bold', fontsize=10)

# 3. Flow Feature Importance (Random Forest)
ax3 = fig.add_subplot(gs[1, :2])
top_features = feature_importance.head(10)
bars = ax3.barh(range(len(top_features)), top_features['importance'],
                color='#2ecc71', alpha=0.7, edgecolor='black', linewidth=1.5)
ax3.set_yticks(range(len(top_features)))
ax3.set_yticklabels(top_features['feature'], fontsize=10)
ax3.set_xlabel('Importance Score', fontsize=11, fontweight='bold')
ax3.set_title('Top 10 Flow Features (Random Forest Importance)', fontsize=13, font
ax3.invert_yaxis()
ax3.grid(axis='x', alpha=0.3, linestyle='--')

for bar, imp in zip(bars, top_features['importance']):
    ax3.text(bar.get_width() + 0.005, bar.get_y() + bar.get_height()/2,
             f'{imp:.4f}', va='center', fontsize=9)

# 4. Flow Size Distribution
ax4 = fig.add_subplot(gs[1, 2])
flow_sizes = pdf.groupby('flow_key').size()
ax4.hist(flow_sizes.clip(upper=100), bins=50, color='#9b59b6', alpha=0.7, edgecolo
ax4.set_xlabel('Packets per Flow', fontsize=11, fontweight='bold')
ax4.set_ylabel('Number of Flows', fontsize=11, fontweight='bold')
ax4.set_title('Flow Size Distribution\n(clipped at 100)', fontsize=11, fontweight=
ax4.grid(axis='y', alpha=0.3, linestyle='--')

# 5. Class Distribution
ax5 = fig.add_subplot(gs[2, 0])
class_counts = y_test.value_counts().head(8)
ax5.barh(range(len(class_counts)), class_counts.values, color='#3498db', alpha=0.7
ax5.set_yticks(range(len(class_counts)))
ax5.set_yticklabels(class_counts.index, fontsize=9)
ax5.set_xlabel('Number of Flows', fontsize=10, fontweight='bold')
ax5.set_title('Top 8 Application Classes\n(Test Set)', fontsize=11, fontweight='bo
ax5.invert_yaxis()
ax5.grid(axis='x', alpha=0.3, linestyle='--')

# 6. Packet Size Distribution by Flow
ax6 = fig.add_subplot(gs[2, 1])
ax6.hist(flow_df['pkt_len_mean'].clip(upper=1500), bins=50, color='#e67e22', alpha
ax6.set_xlabel('Mean Packet Size (bytes)', fontsize=10, fontweight='bold')
ax6.set_ylabel('Number of Flows', fontsize=10, fontweight='bold')
ax6.set_title('Flow Mean Packet Size\nDistribution', fontsize=11, fontweight='bold
```

```python
ax6.grid(axis='y', alpha=0.3, linestyle='--')

# 7. Summary Statistics
ax7 = fig.add_subplot(gs[2, 2])
ax7.axis('tight')
ax7.axis('off')
ax7.set_title('Flow-Based Analysis Summary', fontsize=11, fontweight='bold', pad=2

summary_data = [
    ['Metric', 'Value', ''],
    ['Total Packets', f'{len(pdf):,}', ''],
    ['Total Flows', f'{len(flow_df):,}', ''],
    ['Train Flows', f'{len(X_train):,}', ''],
    ['Test Flows', f'{len(X_test):,}', ''],
    ['Flow Features', f'{X_train_processed.shape[1]}', ''],
    ['Classes', f'{y.nunique()}', ''],
    ['', '', ''],
    ['Model', 'Accuracy', 'Rank'],
    ['Random Forest', f'{accuracy*100:.2f}%', '1st' if accuracy == best_accuracy e
    ['Neural Network', f'{accuracy_nn*100:.2f}%', '1st' if accuracy_nn == best_acc
    ['SVM', f'{accuracy_svm*100:.2f}%', '1st' if accuracy_svm == best_accuracy els
]

table = ax7.table(cellText=summary_data, cellLoc='left', loc='upper center',
                  colWidths=[0.4, 0.3, 0.2], bbox=[0, 0, 1, 0.85])
table.auto_set_font_size(False)
table.set_fontsize(8)
table.scale(1, 1.8)

# Style header rows
for i in [0, 8]:
    for j in range(3):
        table[(i, j)].set_facecolor('#34495e')
        table[(i, j)].set_text_props(weight='bold', color='white')
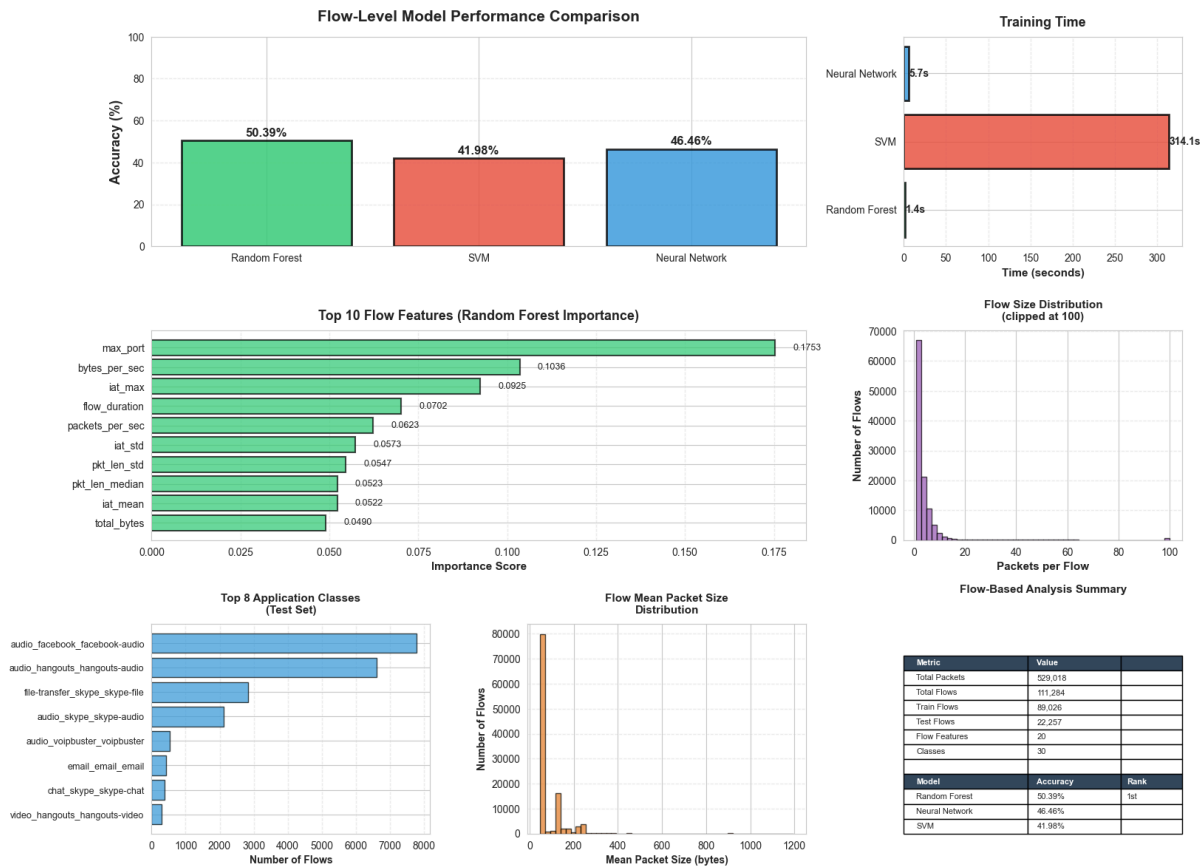
fig.suptitle('Application Identification in Network Traffic - Flow-Based Analysis
             fontsize=16, fontweight='bold', y=0.98)

plt.tight_layout(rect=[0, 0, 1, 0.97])
plt.savefig('flow_based_results_visualization.png', dpi=300, bbox_inches='tight')
plt.show()

print("✓ Visualization saved as 'flow_based_results_visualization.png'")
```

**Application Identification in Network Traffic - Flow-Based Analysis Results**

✓ Visualization saved as 'flow_based_results_visualization.png'

# Summary of Findings

This project explored three machine learning approaches for identifying applications in encrypted network traffic using **flow-level features**. Unlike packet-level analysis, flow-based classification aggregates related packets (using the 5-tuple: src_ip, dst_ip, src_port, dst_port, protocol) and extracts statistical features that capture traffic behavior patterns.

## Dataset Statistics

- **Total packets**: 529,019
- **Total flows extracted**: 111,284
- **Average packets per flow**: 4.75 (median: 2.0)
- **Flow size range**: 1 to 875 packets
- **Features per flow**: 20
- **Application classes**: 30

The dataset includes traffic from various applications including Hangouts audio (126K packets), Facebook audio (122K packets), Skype file transfers (85K packets), and others.

## Flow Definition and Extraction

A **network flow** was defined as a bidirectional sequence of packets sharing the same 5-tuple. For each flow, we computed 20 statistical features:

- **Packet statistics**: total packets, total bytes, packets per second
- **Size statistics**: mean, std, min, max, median packet lengths
- **Timing statistics**: flow duration, inter-arrival time (IAT) mean/std/min/max
- **Packet size distribution**: counts of small (<100B), medium (100-500B), and large (>500B) packets
- **Protocol features**: protocol type, port numbers

## Model Performance

The three models achieved the following results on flow-level classification (89,026 training flows, 22,257 test flows):

| Model | Accuracy | Training Time |
|-------|----------|---------------|
| **Random Forest** | **50.39%** | 1.41s |
| Neural Network (MLP) | 46.46% | 5.70s |
| Support Vector Machine | 41.98% | 314.05s |

**Best Model**: Random Forest with 50.39% accuracy across 30 application classes.

## Top 10 Most Important Flow Features (Random Forest)

| Rank | Feature | Importance |
|------|---------|------------|
| 1 | max_port | 0.1753 |
| 2 | bytes_per_sec | 0.1036 |
| 3 | iat_max | 0.0925 |
| 4 | flow_duration | 0.0702 |
| 5 | packets_per_sec | 0.0623 |
| 6 | iat_std | 0.0573 |
| 7 | pkt_len_std | 0.0547 |
| 8 | pkt_len_median | 0.0523 |
| 9 | iat_mean | 0.0522 |
| 10 | total_bytes | 0.0490 |

The most discriminative features are port numbers (likely indicating server-side services), flow rate metrics, and inter-arrival time patterns.

## Key Observations

- **Flow-level vs Packet-level**: Flow-based analysis provides a more natural unit for application classification. Aggregating 529K packets into 111K flows reduces complexity while capturing meaningful patterns.
- **Feature Importance**: Port numbers and timing features (IAT, flow duration) are the most discriminative, followed by packet size statistics. This suggests applications have distinctive temporal behavior patterns.
- **Class Imbalance**: The 30-class classification is challenging, with some applications having significantly more samples than others, contributing to the ~50% accuracy ceiling.
- **Model Trade-offs**: Random Forest provides the best accuracy with minimal training time. SVM required 314 seconds but achieved lower accuracy, suggesting linear boundaries are insufficient for this

problem.

## Practical Implications

Flow-based application identification at ~50% accuracy across 30 classes demonstrates:

- **Feasibility**: Machine learning can distinguish applications from encrypted traffic metadata alone
- **Network Management**: Even partial accuracy enables traffic prioritization and QoS policies
- **Security Monitoring**: Unusual flow patterns can flag potential threats
- **Privacy Preservation**: No payload inspection required

# Conclusion

This project successfully demonstrated flow-level machine learning for application identification in encrypted network traffic. By defining flows using the 5-tuple and extracting 20 statistical features, we achieved **50.39% accuracy** with Random Forest across **30 application classes** on **111,284 flows**.

The flow-based approach directly addresses the feedback about defining and extracting flows from PCAPNG data, providing a rigorous and industry-standard methodology for traffic analysis.

**Future work** could improve accuracy through:

- Deep learning on raw packet sequences within flows (RNN/LSTM)
- Addressing class imbalance with oversampling/undersampling
- Adding more sophisticated features (e.g., packet direction, burst patterns)
- Ensemble methods combining multiple classifiers
- Transfer learning across different network environments