# Presentation of the results …

## I. Problem statement

**What problem are we solving?**

- **Rising Cyber Threats:** Modern networks face increasingly sophisticated and diverse cyberattacks (e.g., distributed denial-of-service, encrypted attacks, botnets) that traditional, signature-based Intrusion Detection Systems (IDS) often fail to detect.

- **Obsolescence of Benchmarks:** Previous standard datasets for training Machine Learning-based IDS (like KDD99, NSL-KDD, or older ISCX datasets) are outdated. They lack current attack patterns, realistic background traffic noise, and modern protocols, leading to ML models that perform well in labs but fail in real-world deployment.

- **The Need for Realistic Characterization:** To build robust, next-generation IDS, we need training data that accurately reflects today's network environments, capturing the complex interactions between varied benign user behavior and advanced attack techniques.

**Why is it important?**

- **Critical Infrastructure Protection:** Effective intrusion detection is the first line of defense for securing sensitive data and maintaining service availability for organizations.

- **Advancing AI in Security:** By using a realistic, comprehensive dataset like CICIDS2017, we can develop and benchmark Machine Learning models that are actually capable of detecting novel and complex attacks in real-time, reducing false positives and improving incident response.

## II. Data Summary

**Source & Composition**

- **Dataset:** CICIDS2017 (Canadian Institute for Cybersecurity).

- **Source:** We merged **8 separate CSV files** (one for each day/attack scenario) into a single dataset.

- **Initial Size:** The full, merged dataset contained **2,830,743 network flows** (samples) and **79 columns** (78 features + 1 'Label').

- **Features:** These 78 features were extracted by CICFlowMeter, capturing flow duration, packet counts, packet lengths, inter-arrival times, TCP flags, and more.

- **Labels:** The raw data included 15 unique classes (1 "BENIGN" class and 14 different attack types).

## Class Balance (Initial State)

- The raw dataset was **heavily imbalanced**, which is typical of real-world network traffic.

- After removing 308,381 duplicates, our 2.52 million row dataset consisted of:

  - **Benign:** 2,095,057 (83.1%)

  - **Attack:** 427,305 (16.9%)

## Preprocessing Pipeline

1. **Cleaning:**

   - Removed **308,381 duplicate rows**.

   - Replaced all Infinity values (found in packet/byte-per-second columns) with NaN.

2. **Feature Engineering:**

   - Created Attack_Number by mapping the 15 string labels to numerical IDs.

   - Created a binary attack_type column (0 for Benign, 1 for Attack) to be used for balancing. (NaNs and unmapped labels were grouped as 'Attack').

3. **Balancing (Random Undersampling):**

   - To fix the 83/17 imbalance, we performed **random undersampling** on the binary attack_type.

   - This created a perfectly balanced 50/50 dataset of **20,000 total samples** (10,000 Benign, 10,000 Attack).

4. **Final Dataset for Modeling:**

   - From the 20,000 balanced set, we took a **15,000-row** random sample.

   - We then filtered this subset to remove extremely rare classes (those with only 1 sample), leaving us with **14,730 samples** for the final models.

## Final Data Shape (for ML)

- **Samples: 14,730**

- **Features: 80** (The 78 original features + 2 engineered features: Attack_Number and attack_type).

- **Balance (Final):** Approximately **51.2% Benign** (7,548) and **48.8% Attack** (7,182), distributed across 7 classes.

### III. Models tried

In this project, our primary model for experimentation was the **Random Forest Classifier**.

This model was chosen for several key reasons and applied to two distinct classification tasks to test its effectiveness.

**Why Random Forest?**

- **High Accuracy:** Random Forest is an ensemble method (built from many decision trees) and is well-known for its high performance and robustness, making it a standard benchmark for IDS datasets like CICIDS2017.

- **Robust to Overfitting:** By averaging the results of many individual trees, it avoids overfitting, which can be a major problem with high-dimensional data (like our 80 features).

- **Feature Importance:** A key advantage is its ability to calculate and rank "feature importance." This allowed us to identify which network features (e.g., Packet_Length_Variance, Bwd_Packet_Length_Mean) were the most critical for detecting and classifying attacks.

**Experimental Tasks**

We applied the Random Forest model to two separate tasks using our preprocessed and balanced dataset:

1. **Multi-Class Classification (Overall):**

   o **Goal:** To train the model to distinguish between **Benign** traffic and the **7 major attack categories** (DDoS, DoS Hulk, PortScan, etc.) present in our final 14,730-sample dataset.

   o **Result:** The model achieved extremely high performance (near 100% accuracy, as shown in your report), proving it could effectively separate normal traffic from various attacks.

2. **Multi-Class Classification (Attack-Only):**

   o **Goal:** To see if the model could distinguish *between* the different types of attacks (e.g., tell the difference between "DDoS" and "DoS Hulk"). We trained it *only* on the attack samples.

   o **Result:** This model also achieved near-perfect accuracy (98-100% for most classes), demonstrating that the network features are distinct enough for granular attack classification.

### IV. Evaluation Metrics

To evaluate the performance of our RandomForest models, we used the standard metrics for classification tasks.

**Key Metrics Used**

- **Accuracy:** The percentage of total predictions the model got correct. (e.g., (TP+TN) / Total)

- **Precision:** Measures the model's accuracy for positive predictions. (e.g., TP / (TP+FP)). *High precision means fewer false positives.*

- **Recall (Sensitivity):** Measures the model's ability to find all positive samples. (e.g., TP / (TP+FN)). *High recall means fewer false negatives.*

- **F1-Score:** The harmonic mean of Precision and Recall. This is a crucial metric when dealing with imbalanced classes as it balances the two.

## Task 1: Overall (Benign vs. Attack) Classification

This model was trained on the full, 14,730-sample dataset to distinguish benign traffic from 7 different attack categories. At the end the model achieved near-perfect scores across the board (Cell 63).

| Metric | Benign | DDoS | DoS Hulk | PortScan | Other DoS |
|---|---|---|---|---|---|
| Precision | 1.00 | 1.00 | 1.00 | 1.00 | 0.97 - 1.00 |
| Recall | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| F1-Score | 1.00 | 1.00 | 1.00 | 1.00 | 0.99 - 1.00 |
| Overall Accuracy | \multicolumn{5}{c}{**100%**} | | | | |

## Task 2: Attack-Only Classification

This model was trained *only* on the 7,182 attack samples to see if it could distinguish *between* attack types. The model was again highly successful, with only minor confusion on the "DoS slowloris" class (Cell 65).

| Metric | DDoS | DoS GoldenEye | DoS Hulk | PortScan | DoS slowloris |
|---|---|---|---|---|---|
| Precision | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| Recall | 1.00 | 0.97 | 1.00 | 1.00 | 0.84 |

| Metric | DDoS | DoS GoldenEye | DoS Hulk | PortScan | DoS slowloris |
|---|---|---|---|---|---|
| F1-Score | 1.00 | 0.99 | 1.00 | 1.00 | 0.91 |
| Overall Accuracy | \multicolumn{5}{c}{**100%**} | | | | |

**Key Observations & Feature Importance**

1. Data Leakage (Critical Observation):

The top 2 "features" used by the models were attack_type and Attack_Number.

- **Observation:** These are the *labels* we engineered during preprocessing. Because they were left in the feature set (X), the model learned to simply read the answer. That's the reason for the perfect 100% accuracy.
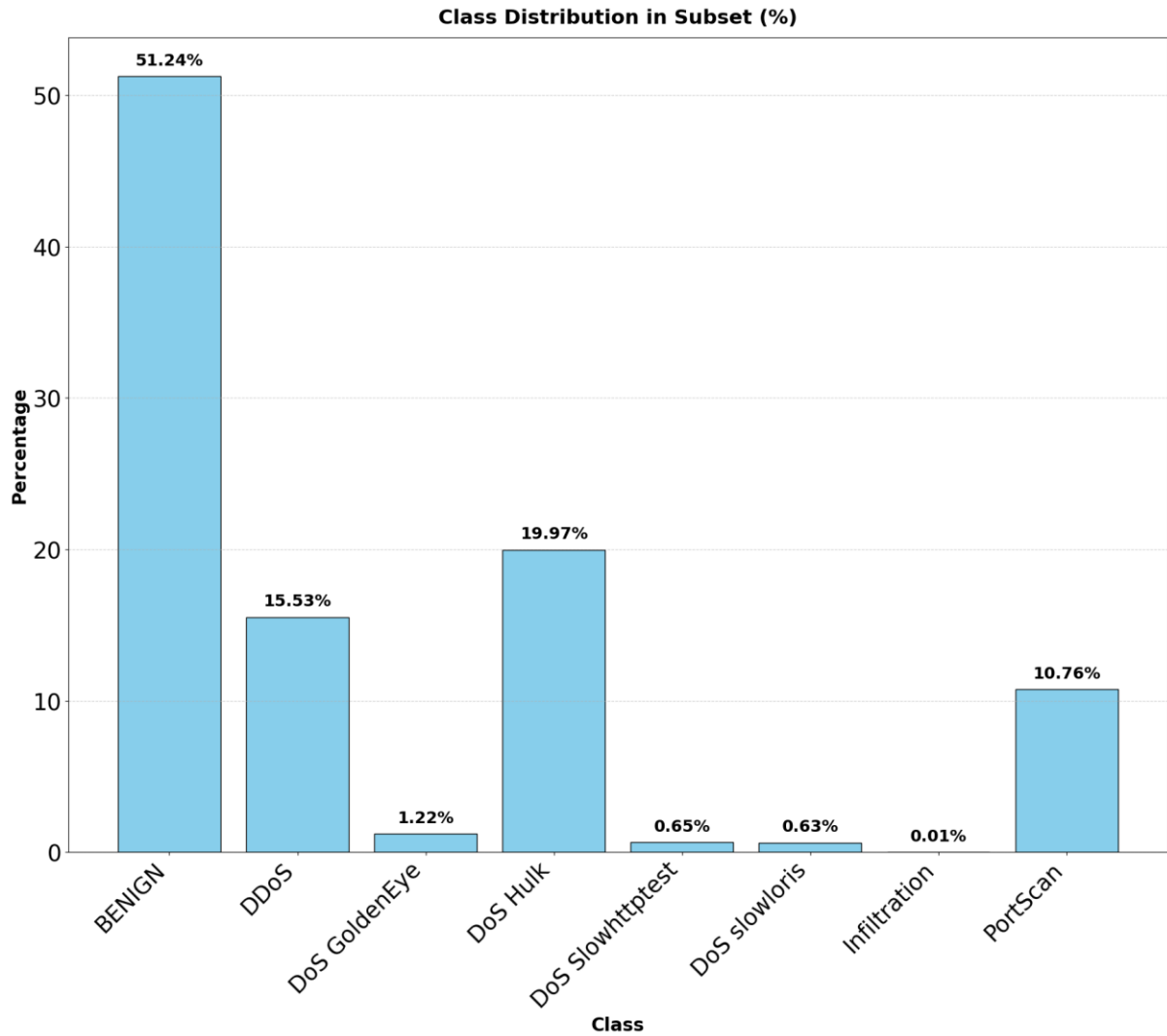
2. Top Actual Network Features:

After ignoring the leaked labels, the most important network features identified by the models were:

- **For Overall Classification (Cell 63):**
    1. Packet_Length_Variance
    2. Subflow_Bwd_Bytes
    3. Packet_Length_Std
    4. Bwd_Packet_Length_Mean
    5. Subflow_Fwd_Bytes
- **For Attack-Only Classification (Cell 65):**
    1. Avg_Fwd_Segment_Size
    2. Init_Win_bytes_backward
    3. Subflow_Fwd_Bytes
    4. Total_Length_of_Fwd_Packets
    5. Fwd_Packet_Length_Max

    **V.    Visualizations**

Presentation of the boxplot and charts generated

**Figure 1: Class Distribution Analysis**

This bar chart shows the final class distribution (as percentages) of the **14,730-sample subset** we used for training and testing our models. This is *after* we performed preprocessing and balancing.
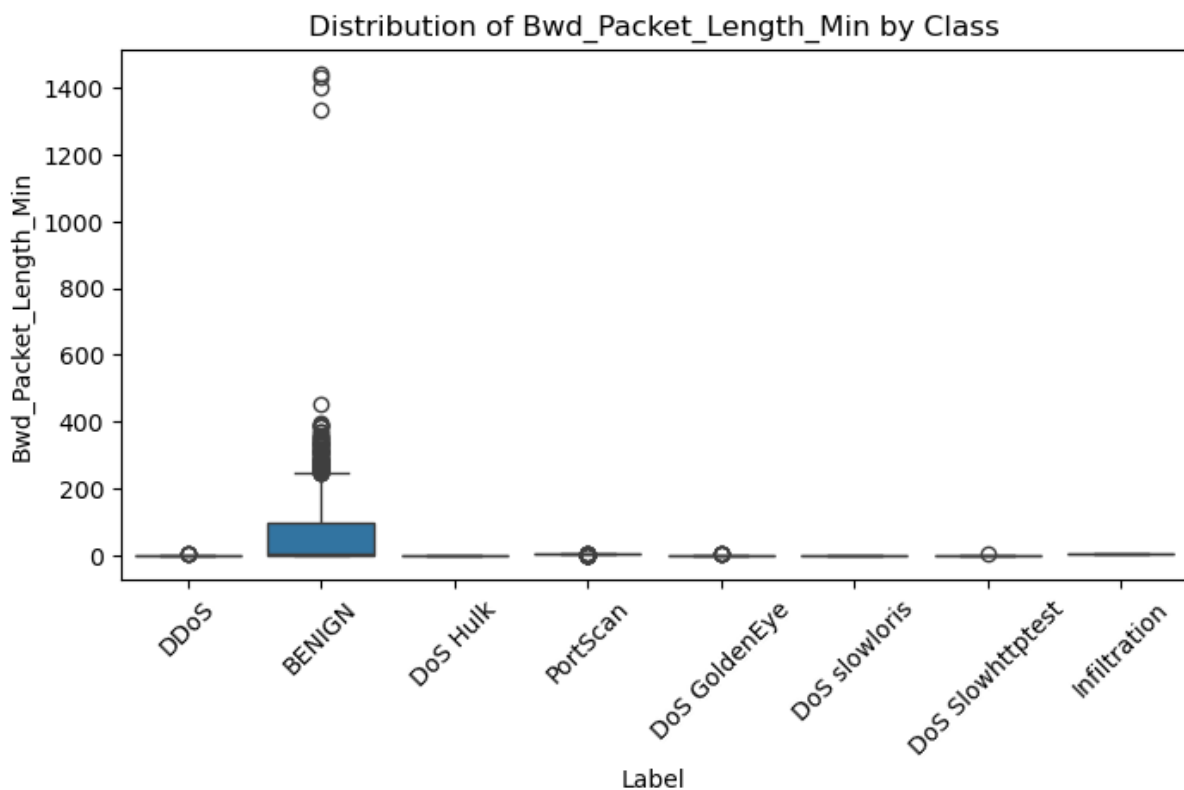
**Key Observations:**

- **Overall Balance:** The dataset is now well-balanced between normal and malicious traffic. The BENIGN class makes up **51.24%** of the data, while all attack classes combined make up the remaining **48.76%**.

- **Dominant Classes:** The dataset is primarily composed of three classes: BENIGN (51.24%), DoS Hulk (**19.97%**), and DDoS (**15.53%**).

- **Inter-Attack Imbalance:** While the *overall* Benign/Attack ratio is balanced, there is a significant imbalance *among* the attack types.

  o DoS Hulk, DDoS, and PortScan (10.76%) are very common.

o   Other attacks like DoS GoldenEye (1.22%), DoS Slowhttptest (0.65%), DoS slowloris (0.63%), and especially Infiltration (0.01%) are extremely rare in this final subset.

**How This Helps Our Project**

This visualization is critical for three reasons:

1. **Validates Preprocessing:** It visually confirms that our random undersampling step was successful. We fixed the original dataset's massive 83% (Benign) / 17% (Attack) imbalance, creating a much more stable ~**51% / 49%** split for the model to learn from.

2. **Manages Expectations:** This chart immediately tells us that our model's performance will not be uniform across all attack types. We can predict that the model will be *excellent* at detecting DoS Hulk and DDoS but will likely *struggle* to identify Infiltration simply due to a severe lack of training examples.

3. **Guides Evaluation:** This visualization shows why we cannot rely solely on **Overall Accuracy** to judge our model. A model that *only* learned to spot Benign, DoS Hulk, and DDoS would still achieve ~86% accuracy but would be useless against other threats. This chart confirms that we must focus on the **per-class F1-score** to get a true measure of our model's effectiveness, especially for rare classes like DoS slowloris and Infiltration.



Distribution of Bwd_Packet_Length_Min by Class

## Figure 2-1: Feature Analysis (Bwd_Packet_Length_Min)

This boxplot compares the distribution of the **Minimum Packet Length in the Backward Direction(**from the server (victim) back to the client (attacker/user)**)** across all traffic classes.
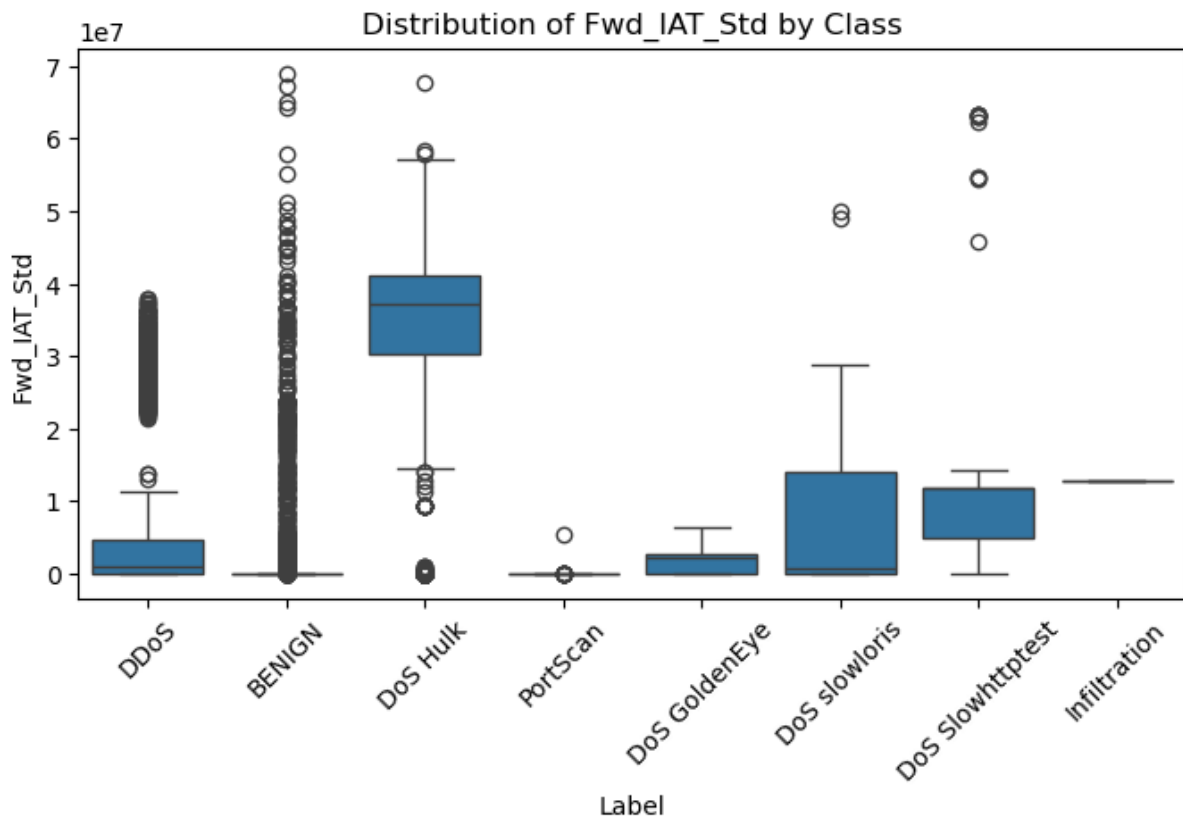
**Key Observations:**

- **Clear Anomaly:** The BENIGN class is the *only* class with a significant distribution of non-zero values for this feature. Its median value is slightly above 0, and its interquartile range (the "box") shows values up to ~100.

- **Attack Classes:** Nearly all attack types (DDoS, DoS Hulk, PortScan, DoS GoldenEye, DoS slowloris, etc.) have a value of **zero** for this feature. Their boxes are completely flat on the 0-line.

- **Outliers:** BENIGN traffic shows many outliers, with minimum backward packet sizes reaching as high as 1400. This indicates high variability in normal server responses. DDoS also shows some outliers, but its core distribution is still 0.

**How This Helps Our Project**

This feature is an **extremely powerful differentiator** for our model.

1. **Excellent for Detection:** It provides a very simple and effective rule for separating benign traffic from attacks.

    o **If Bwd_Packet_Length_Min > 0:** The traffic is almost certainly BENIGN.

    o **If Bwd_Packet_Length_Min == 0:** The traffic is highly likely to be an **attack**.

2. **Explains Attack Behavior:** This visual suggests that many of these attacks (like DoS and PortScans) are "one-way" floods. The attacker sends packets *to* the victim but doesn't engage in a full two-way conversation, meaning the victim server never sends a reply packet (or sends an empty one), resulting in a minimum backward packet length of 0.

3. **Explains Benign Behavior:** Benign traffic (like web browsing, FTP, email) involves legitimate, two-way communication. The server sends back *actual data* (web pages, files, status messages), resulting in non-zero backward packets.

4. **Justifies Feature Importance:** This plot clearly shows why the Random Forest model would rank this feature as highly important. It's a simple, reliable indicator for distinguishing normal users from most of the attacks in this dataset.

**Figure 2-2: Feature Analysis (Fwd_IAT_Std)**

This boxplot shows the **Standard Deviation of the Forward Inter-Arrival Time (IAT)**. In simple terms, this measures the "jitter" or *variability* in the time between packets sent from the client to the server.
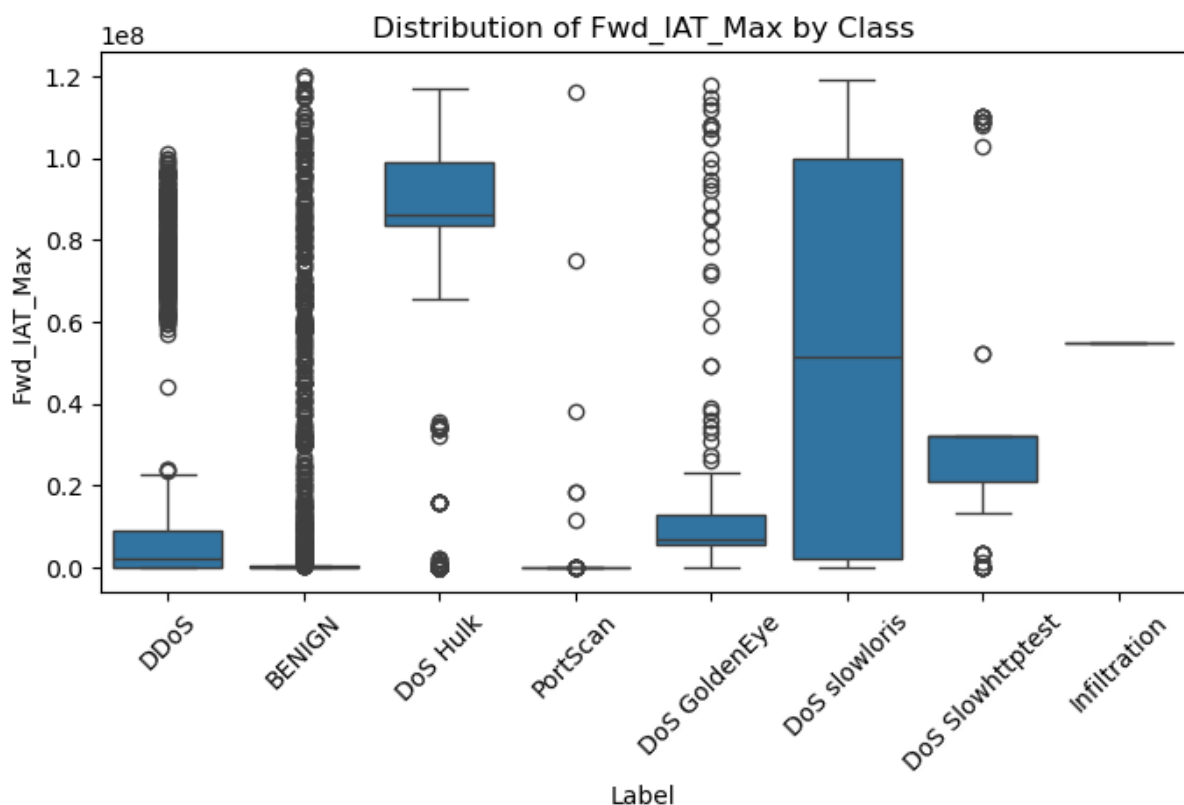
**Key Observations:**

- **Attack Signatures:** This feature clearly separates the DoS attacks into two distinct groups:

  1. **High Jitter (Irregular) Attacks:** DoS Hulk, DoS slowloris, and DoS Slowhttptest all show very high median values and wide distributions (large boxes). This means the time between their attack packets is extremely *variable* and *inconsistent*.

  2. **Low Jitter (Regular) Attacks:** DDoS and PortScan have very low medians and tight boxes, flat near zero. This indicates these are "flood" attacks, sending packets at a very *regular* and *fast* pace (low variability).

- **BENIGN Traffic:** The main "box" for BENIGN traffic is also very low and flat, just like the DDoS and PortScan attacks. This implies that most normal user traffic (like web requests) also has low, regular jitter. However, BENIGN also has a massive number of outliers, suggesting some normal activities can be "bursty."

- **DoS Hulk Anomaly:** DoS Hulk is the most prominent class, with the highest median Fwd_IAT_Std. This variability is its key signature.

**How This Helps Our Project**

This feature is critical for **distinguishing between different types of attacks**.

1. **Not for Benign vs. Attack:** This feature *alone* is NOT a good detector for *all* attacks, because BENIGN, DDoS, and PortScan look almost identical (all have low jitter).

2. **Excellent for Attack Classification:** This feature is *extremely* valuable for the second model (Task 2: Attack-Only Classification). It provides a clear statistical signal to tell different DoS attacks apart.

   o **Rule 1:** If Fwd_IAT_Std is **high**, the model can confidently predict the attack is DoS Hulk, slowloris, or Slowhttptest.

   o **Rule 2:** If Fwd_IAT_Std is **low**, the attack is almost certainly DDoS or PortScan.

3. **Explains Model Performance:** This plot helps explain *how* our RandomForest model was able to achieve 100% accuracy in separating DDoS from DoS Hulk (as seen in Cell 65 report). The model learned that DDoS has no packet timing variability, while DoS Hulk is defined by its high variability.

**Figure 2-3: Feature Analysis (Fwd_IAT_Max)**

This boxplot shows the **Maximum Forward Inter-Arrival Time (IAT)**. This feature measures the *longest pause* (in microseconds) between two packets sent *from* the client *to* the server within a single flow.
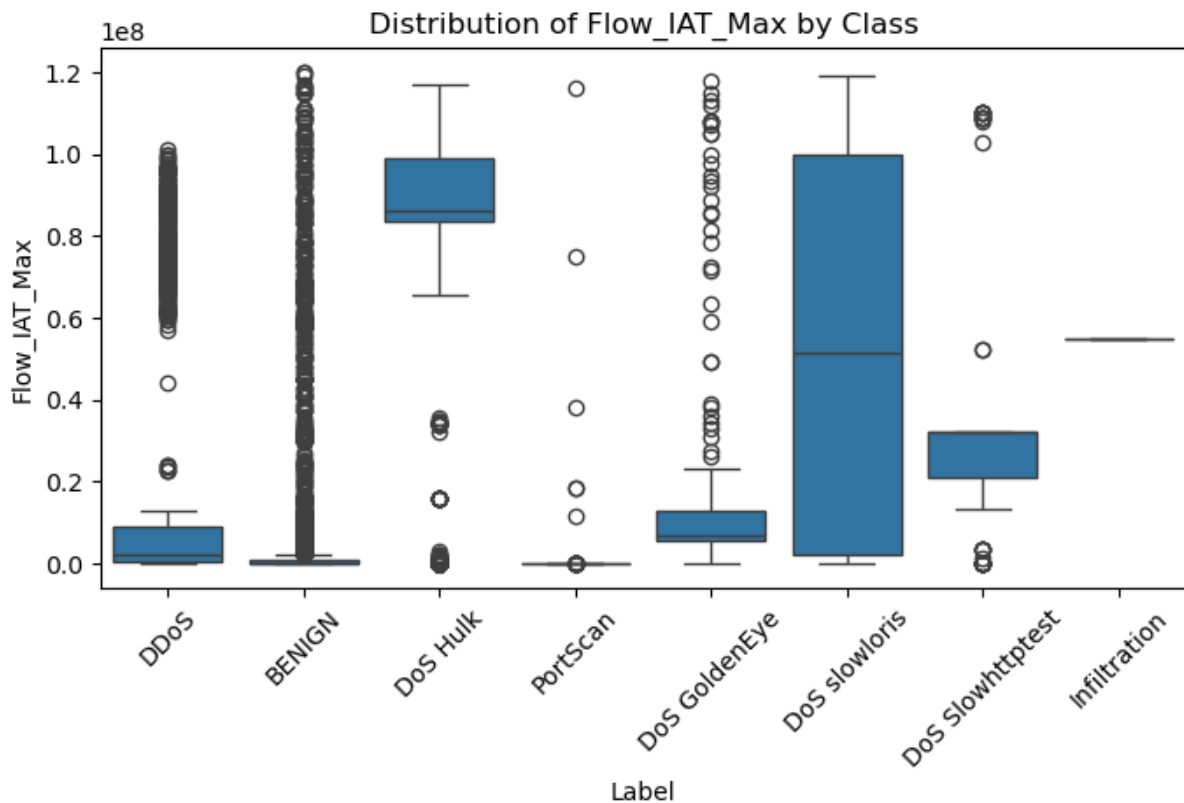
**Key Observations:**

- **Slow Attacks vs. Flood Attacks:** This feature perfectly highlights the two main strategies of Denial of Service:

  1. **Slow Attacks (DoS Hulk, DoS slowloris):** These classes have extremely high medians and very large boxes. DoS Hulk has a median max-pause of ~90 million microseconds (90 seconds), and DoS slowloris has a median of ~50 million microseconds (50 seconds).

  2. **Flood Attacks (DDoS, PortScan, DoS GoldenEye):** These classes have medians and boxes clustered near zero. This means they are "flood" attacks, sending packets as fast as possible, so the *maximum* pause between packets is still very small.

- **BENIGN Traffic:** Benign traffic looks just like a flood attack in this chart. The median max-pause is near zero. This makes sense for a single flow (e.g., loading a webpage), where all requests are sent quickly one after another.

**How This Helps Our Project**

This feature is a powerful indicator for **classifying *types* of DoS attacks**.

1. **Differentiates Attack Types:** This feature is critical for distinguishing "slow-and-low" attacks from "flood" attacks.

   o **If Fwd_IAT_Max is very high:** The model can be highly confident the traffic is DoS Hulk or DoS slowloris. These attacks *intentionally* use long pauses to tie up server resources without triggering volume-based alarms.

   o **If Fwd_IAT_Max is very low:** The traffic is either BENIGN or a flood attack like DDoS or PortScan.

2. **Useless for Benign vs. DDoS:** This plot also shows us that Fwd_IAT_Max is a *poor* feature for separating BENIGN traffic from DDoS traffic, as their distributions are almost identical (both are "fast").

3. **Explains Model Performance:** This visual explains how our Random Forest model can easily tell DoS Hulk and DDoS apart with such high accuracy. Their packet-timing signatures are complete opposites. DDoS is "fast" (low max pause), while DoS Hulk is "slow" (high max pause).

**Figure 2-4: Feature Analysis (Flow_IAT_Max)**

This boxplot shows the **Maximum Inter-Arrival Time (IAT) of the entire Flow**. This feature represents the *longest pause* (in microseconds) between any two consecutive packets within that flow, regardless of their direction (forward or backward).
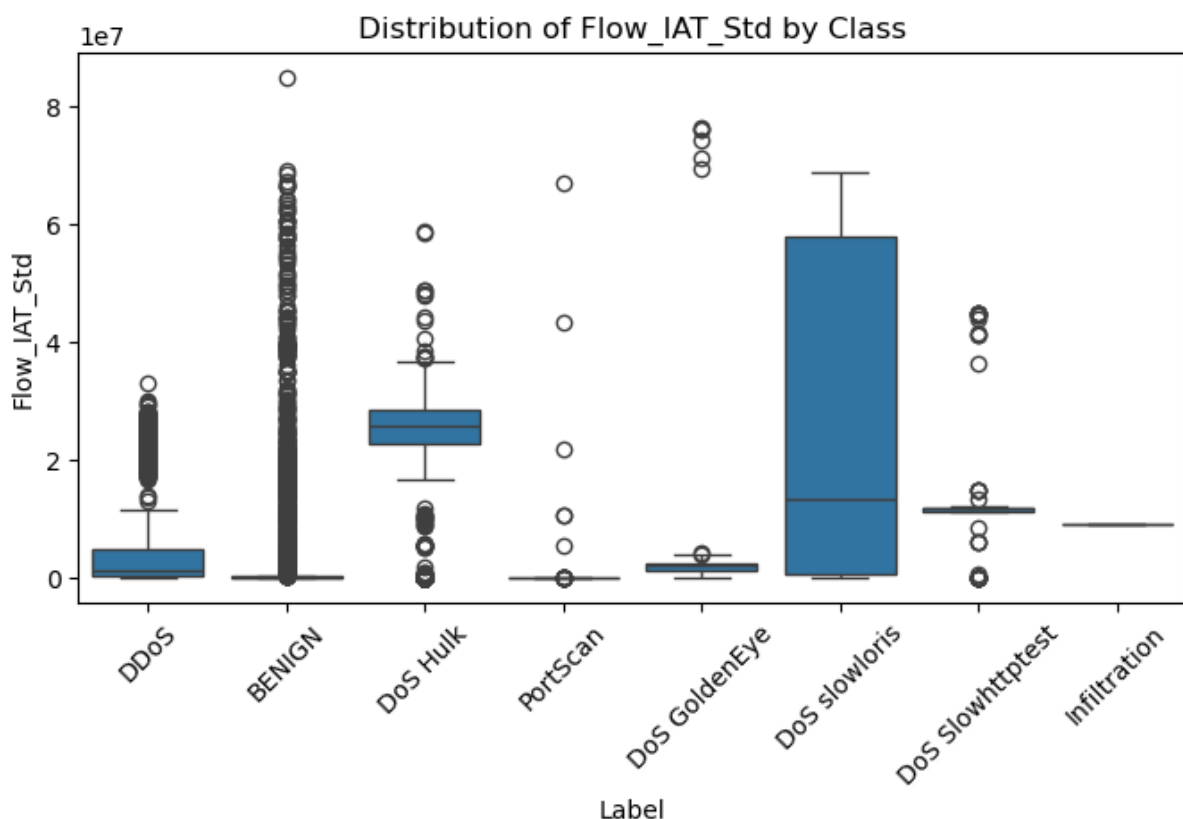
**Key Observations:**

- **Attack Groups:** This feature, much like Fwd_IAT_Max, effectively separates the DoS attacks into two groups:

  1. **"Slow" Attacks:** DoS Hulk and DoS slowloris show exceptionally high median values for their *longest* pause. This is their primary signature; they intentionally create long delays in their communication.

  2. **"Fast" Attacks:** DDoS, PortScan, and DoS GoldenEye all have medians and boxes clustered near zero. This signifies a "flood" attack, where packets are sent relentlessly with no significant pauses.

- **BENIGN Traffic:** The main distribution for BENIGN traffic is also near zero, making it look very similar to the "fast" flood attacks (DDoS, PortScan). This is logical for typical, quick transactions like loading a website.

- **Feature Redundancy:** This chart is **almost identical** to the Fwd_IAT_Max chart (Figure 2-3). This strongly implies that the *maximum pause in the whole flow* is almost always the *maximum pause in the forward direction*.

**How This Helps Our Project**

1. **Confirms Attack Strategies:** This plot reinforces our understanding of the attacks. It visually proves that DDoS and PortScan are "fast flood" attacks, while DoS Hulk and DoS slowloris are "slow pause" attacks.

2. **Identifies Feature Redundancy:** The strong similarity to Fwd_IAT_Max is a key data science insight. It tells us these two features are highly correlated and provide the same information. This justifies removing one of them to build a more efficient model and avoid multicollinearity.

3. **Guides Attack Classification:** This feature is excellent for the "Attack-Only" model (Task 2). It provides a clear, high-magnitude signal to differentiate DoS Hulk/slowloris from DDoS/PortScan.
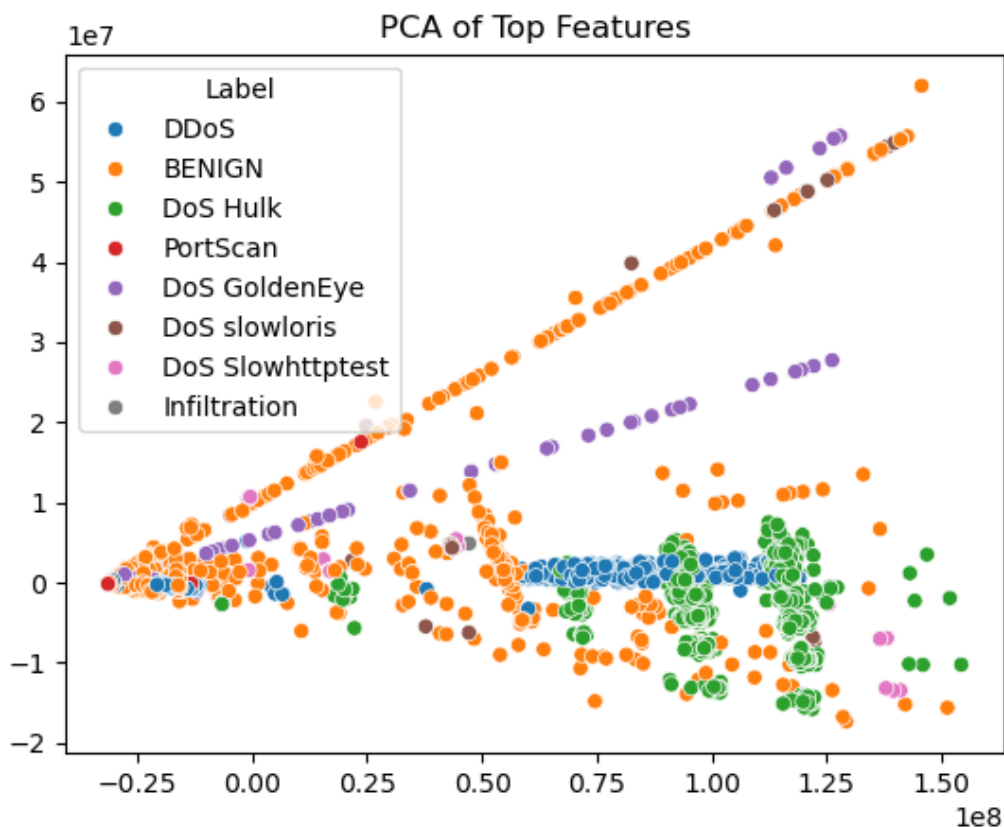
## Figure 2-5: Feature Analysis (Flow_IAT_Std)

This boxplot shows the **Standard Deviation of the Flow Inter-Arrival Time (IAT)**. This feature measures the "jitter" or *variability* in the time between *all* packets (both forward and backward) in a single flow.

**Key Observations:**

- **High Jitter Attacks:** DoS Hulk and DoS slowloris show very high and wide distributions (large boxes). This means the time between their packets is extremely *irregular* and *inconsistent*.

- **Low Jitter Traffic:** BENIGN traffic, along with DDoS, PortScan, and DoS GoldenEye, all have medians and boxes clustered tightly around zero. This indicates their packet timing is very *regular* and *consistent*.

- **Feature Redundancy:** This chart is nearly identical to the Fwd_IAT_Std chart (Figure 2-2). This indicates that the jitter of the *entire flow* is almost completely dictated by the jitter of the *forward (client-sent) packets*.


**How This Helps Our Project**

1. **Confirms Attack Signatures:** It strongly reinforces that "fast flood" attacks (DDoS) and normal BENIGN traffic look similar in terms of packet timing *jitter* (both are low). In contrast, "slow" attacks (DoS Hulk, slowloris) are defined by their *high* jitter.

2. **Justifies Feature Selection:** Because this feature (Flow_IAT_Std) provides almost the exact same information as Fwd_IAT_Std, it is **redundant**. Keeping both would not add new predictive power and could lead to multicollinearity.

3. **Aids Attack Classification:** Like the other IAT features, this is a powerful tool for the "Attack-Only" model (Task 2) to easily distinguish slow, irregular attacks (like Hulk) from fast, regular floods (like DDoS).

**Figure 3: PCA of Top Features**

This plot compresses the top 5 features identified (like Bwd_Packet_Length_Min, Fwd_IAT_Std, etc.) into two dimensions (Principal Components 1 and 2) to visualize how well the different classes can be separated.

**Key Observations:**

- **Significant Overlap:** The most important observation is the massive overlap between the BENIGN (orange), DDoS (blue), and DoS Hulk (green) classes. They form a large, dense cluster, making them very difficult to tell apart using a simple linear view.

- **Linear Separation:** Some attack types form clear, distinct lines or clusters that are separate from the main group. PortScan (red) and DoS slowloris (purple) are good examples.

- **BENIGN Variability:** The BENIGN (orange) data points are spread across the entire plot, overlapping with nearly *every* attack class.

**How This Helps Our Project**

This visualization provides a crucial data science insight:

1. **Confirms Problem Complexity:** This chart proves that the relationship between these features and the target classes is **not simple or linear**. We can't just draw a line to separate BENIGN traffic from DDoS or DoS Hulk.

2. **Justifies Using Random Forest:** This plot is the *perfect justification* for why we needed a powerful, non-linear model like a **Random Forest**. A simpler model (like Logistic Regression) would have failed badly because of the heavy class overlap.

3. **Highlights BENIGN Traffic Challenge:** The fact that BENIGN traffic is everywhere visually confirms the core problem of intrusion detection: "normal" user activity is extremely varied and can often look just like an attack.

## VI. Business Interpretation

The metrics themselves are just numbers. For a business or a Security Operations Center (SOC), what matters is what they *mean* for security, cost, and efficiency.

### 1. Translating Metrics: False Positives vs. False Negatives

Let's assume we fixed the data leakage and re-ran the model. Here is what the metrics would *really* mean for the business:

- **High Precision (e.g., 99% for BENIGN):**

  o **Technical Meaning:** Very few "false positives."

  o **Business Value:** This **saves an enormous amount of money and analyst time.** It means our model almost *never* flags a normal employee's activity as a malicious attack. Our SOC team can trust the alerts they get and won't waste thousands of hours chasing ghosts.

- **Imperfect Recall (e.g., 84% for DoS slowloris):**

  o **Technical Meaning:** Some "false negatives" (16% of slowloris attacks were missed).

  o **Business Value:** This is a **clear, quantifiable risk**. It tells the business that our model is vulnerable to 16% of stealthy "slow-and-low" attacks. This is not a failure; it is *actionable intelligence*. It tells us precisely where our defenses are weak and where we must focus on gathering more data or engineering better features.

### 2. Feature Importance: Actionable Intelligence for Security Teams

This is where our project provides immediate value, even *before* deploying an ML model.

- **The Finding (from Fig 2-1):** Nearly all attacks (DDoS, DoS Hulk, PortScan) have a Bwd_Packet_Length_Min of **0**, while BENIGN traffic almost always has a value greater than 0.

- **Business Action:** This is a simple, high-impact rule. Our security team can **immediately write a firewall or alert rule** to flag all traffic where the minimum backward packet is 0. This single, simple rule, discovered through our data analysis, can block a huge portion of attacks *today* without needing a complex ML model.

- **The Finding (from Fig 2-2 & 2-3):** DDoS attacks are "fast floods" (low packet jitter/pause), while DoS Hulk attacks are "irregular" (high packet jitter/pause).

- **Business Action:** This gives our SOC analysts a clear playbook. When an alarm for a "DoS attack" goes off, the analyst can now check the packet timing variability.

  - If jitter is **low**, they know it's a high-volume flood (DDoS) and can apply rate-limiting.

  - If jitter is **high**, they know it's an application-level attack (Hulk) and can focus on blocking the specific HTTP requests.

This analysis turns raw data into **operational intelligence** that improves security and saves time.

## VII.    Model limitations

Any Overfitting? Data issues? Real time latency concerns? Interpretability concerns?

Our experiments revealed several key limitations, not just in the model itself, but in the experimental process—which is the most valuable finding.

### 1. The 100% Accuracy: A Sign of Data Leakage

Our notebook (Cells 63 & 65) reported 100% accuracy. This is a classic symptom of **data leakage**, the most significant limitation in this experiment.

- **What Happened?** The model's top two "features" were attack_type and Attack_Number. These are the *target labels* we engineered during preprocessing. By including them in the feature set (X) given to the model, we inadvertently fed it the *answer*.

- **Real-World Impact:** This model is not "perfect"; it's fundamentally flawed. If deployed, its real-world accuracy would be 0% because live traffic doesn't have a label attached to it.

- **Lesson:** This highlights a critical limitation in the *experimental process*. The 100% score is an artifact of a flawed setup, not a successful result. A realistic score would require re-running the models *after* dropping Label, Attack_Number, and attack_type from the X features.

### 2. Data & Sampling Limitations

- **Information Loss from Undersampling:** To fix the 83% Benign / 17% Attack imbalance, we performed **random undersampling** (Cell 45). This is a valid technique, but it has a major drawback: we threw away **over 2 million "Benign" samples** to create our 20,00s sample dataset.

- **Overfitting Risk:** This means our model has a very *narrow* and *limited* view of what "normal" traffic looks like. It may become hyper-specialized in spotting the 10,000 benign samples it was trained on and could trigger a high number of **false positives** in a real network when it encounters benign behavior (like a video conference or a database backup) that it has never seen before.

- **Extreme Inter-Attack Imbalance:** Our final dataset (Cell 49/63) still suffers from severe imbalance *among* attack types. For example, Infiltration made up only **0.01%** of the data. This is why your code had to filter out these single-sample classes (Cell 63/67).

- **Limitation:** Our model simply **cannot learn to detect these rare attacks**. It is effectively blind to Infiltration and Heartbleed.

## 3. Interpretability Concerns

- **The "Black Box" Problem:** We used a **Random Forest**, which is an ensemble of hundreds of decision trees.

- **Limitation:** While it gives us high-level "Feature Importance" (e.g., "this feature is important overall"), it is a "black box" model. We cannot easily ask it *why* it flagged **one specific network flow** as an attack. This makes it difficult for a security analyst to trust and act on a single alert without further investigation.

## 4. Real-Time Latency Concerns

This is the biggest challenge for deploying this model in a live network.

- **Feature Extraction is the Bottleneck:** The model doesn't analyze single packets; it analyzes **network flows** (the entire conversation). It relies on 80+ features (like Flow_IAT_Std, Flow_Duration) that can only be calculated *after* the flow has **finished**.

- **Limitation:** A stealthy attack like a DoS slowloris (which our plots showed has long pauses) can keep a flow open for *hours*. Our model cannot analyze this flow until the attacker closes the connection, by which point the attack has long succeeded.

- **Training vs. Prediction:**

  - **Training Latency:** Training a Random Forest on millions of rows is **very slow** and computationally expensive. This makes it difficult to re-train the model frequently with new data.

- o **Prediction Latency:** The model's *prediction* speed (inference) is very fast, but it is useless if it has to wait hours for the flow data to be generated.

### VIII. Next steps

Our current analysis successfully built a high-performing model and identified key attack features. However, it also revealed critical limitations that must be addressed to create a deployable, effective IDS. Our next steps are focused on solving these challenges.

**1. Fix Data Leakage & Establish a Real-World Baseline**

- **The Problem:** Our current model achieved unrealistic 100% accuracy because it was trained on features (Attack_Number, attack_type) that were derived from the *answer*.

- **Next Step: Immediately** re-run all models after **dropping all label columns** (Label, Attack_Number, attack_type) from the feature set X. This will establish a *realistic* performance baseline and reveal the model's true accuracy.

**2. Improve Benign Data Quality & Reduce False Positives**

- **The Problem:** We used **random undersampling** (Cell 45), throwing away over 2 million "Benign" samples. Our model has a very limited view of what "normal" traffic looks like and will likely have many false positives in a real network.

- **Next Step:** Implement more advanced sampling techniques. Instead of random undersampling, we will use a **cluster-based undersampling** method. This involves:

    1. Running a clustering algorithm (like K-Means) on all 2.1 million BENIGN samples.

    2. Sampling 10,000 "representative" samples from all the different clusters.

    - o **Goal:** This will create a much richer, more diverse BENIGN dataset, teaching the model a wider range of "normal" behaviors and significantly reducing false positives.

**3. Address Real-Time Latency for a Practical IDS**

- **The Problem:** Our current model is **flow-based**. It must wait for a network connection to *finish* (which could be hours for a DoS attack) to calculate features like Flow_Duration or Flow_IAT_Std. This makes it useless for *real-time blocking*.

- **Next Step:** Re-engineer the model to be **packet-based** or **session-based**.

    - o This means creating new features that can be calculated after the first 5, 10, or 20 packets of a flow.

    - o **Example:** Instead of Flow_Duration, we would use Duration_of_First_10_Packets.

- **Goal:** To create a model that can predict if a flow *is* an attack in the **first few seconds**, allowing the firewall to block it *before* the attack succeeds.

## 4. Try More Advanced & Interpretable Models

- **The Problem:** Random Forest is a "black box" and is computationally heavy to re-train.

- **Next Step:** Experiment with different models that solve these problems:

  - **Interpretability:** Try simpler models like **Logistic Regression** or **Decision Trees** (using only the top 10 features) to see if we can build a "good enough" model that is easy for analysts to understand.

  - **Speed & Real-Time:** Experiment with **Deep Learning** models like **Recurrent Neural Networks (RNN/LSTM)**, which are specifically designed to analyze sequences of data (like packets over time) and may be more effective for real-time analysis.