# Intrusion Detection Systems (IDS) with IoT Systems

## 1   Description

In today's data-driven world, the ability to leverage big data for insights and security decision-making is a critical skill, particularly in the context of IoT systems. This project offers an exciting opportunity for students to dive into the world of cybersecurity by working with real-world datasets designed to simulate Intrusion Detection Systems (IDS) within IoT environments. Students will explore the multifaceted world of data science and security by working with distinct types of datasets focusing on network traffic and intrusion events.

## 2   Project Goals

This project is designed to provide a holistic understanding of data science, machine learning, and their applicability to cybersecurity in IoT systems. More details:

- **Data Exploration**: You will learn how to navigate and understand complex IoT network traffic datasets, gaining insights into data structures, network behavior patterns, and potential intrusion or anomaly markers.

- **Data Pre-processing**: You will discover techniques to clean, preprocess, and transform raw network traffic data into a usable format, ensuring data quality and removing noise for effective analysis.

- **Visualization**: You will create informative visualizations to effectively present your findings, including network traffic flows, intrusion events, and anomaly patterns.

- **Interpretation**: You will develop the skills to interpret the results of your intrusion detection models, drawing meaningful insights about the effectiveness of the system in detecting and responding to security threats in IoT environments.

## 3   Datasets

The provided datasets have been collected to simulate and analyze Intrusion Detection Systems (IDS) within IoT environments. These datasets are designed to help you develop solutions for detecting and mitigating potential intrusions or anomalies in IoT networks.

The dataset has been divided into several parts, with each part representing a unique section of the overall dataset. These sections focus on specific intrusion scenarios, types of IoT devices, or network events. This structure ensures that each group or individual works on a manageable dataset part, allowing for in-depth exploration and analysis of a particular aspect of IoT security.

Based on your assigned dataset part, you will develop an IDS model using machine learning techniques to detect intrusions. By the end of the project, you will be required to test and validate your model's performance using your specific dataset section. The goal is to evaluate your IDS model's ability to accurately identify potential threats within an IoT environment.

# 4  Group or Individual Work

You have the flexibility to either form groups or work individually. Dataset Distribution. Every group (or individual student) will be assigned a distinct portion of the dataset to work with. The goal of partitioning the dataset is to ensure that each group focuses on a manageable subset of data, allowing for in-depth analysis and experimentation.

# 5  Required Tasks

The following presents a set of questions related to various aspects of data pre-processing, machine learning, and post-processing for both datasets.

## 5.1  Data Loading

1. Load the selected dataset (e.g. into a Pandas DataFrame using appropriate functions like "***pd.read_csv()***" in the case of tabular/time series datasets).

2. Combine data from multiple sources if possible to enrich the dataset.

3. Perform basic data preprocessing tasks such as handling missing values, data type conversions, and data cleaning.

   - Check for missing values using "*df.isna()*" and handle them by either imputing or removing them. Make sure to explain the reasons behind your decision to use imputation or removal techniques to address missing data. If you choose imputation, be sure to clarify which specific method you find most appropriate for your dataset.
   - Convert data types as needed using "*df.astype()*".
   - Clean the missing data by removing duplicates using "*df.drop_duplicates()*" and correcting inconsistent values. You have to indicate which technique was applied and why.

## 5.2  Exploratory Data Analysis (EDA)

Profile the data to gain insights into its distribution, relation, summary statistics, and potential data quality issues (outlier data).

   Therefore, in this part, you will use Matplotlib Seaborn, or Plotly to create a variety of plots that can be :

1. Line charts to visualize trends over time.

2. Scatter plots to identify relationships between numerical variables.

3. Bar plots to compare categorical data.

4. Heatmaps to show correlations between feature/variables.

   Then, you have to :

1. Calculate summary statistics (using functions like *df.describe()*) to understand the dataset's central tendencies and distributions by using histograms, KDE plots, and other visualizations to visualize the data distribution. In addition, you should identify categorical variables and create bar plots to understand their distributions.

2. Calculate the correlation matrices and plot them to identify the relationships between numerical variables. Another plot is required between the variables and the target one (classification case).

## 5.3 Data Manipulation

1. Apply grouping and aggregation to calculate summary statistics for specific categories. You can apply *groupby()* and aggregation functions like *sum()*, *mean()*, and *count()* to create summary tables.

2. *Feature engineering*: Experiment with creating new features or transforming existing ones to enhance class separation. Investigate whether feature engineering results in significant improvements.

## 5.4 Documentation and Presentation

- *Comprehensive Jupyter Notebook*: Create a detailed and well-documented Jupyter Notebook that captures all the steps taken throughout the project. This should include:

  - Data preprocessing steps.
  - Performance metrics and comparisons of models.
  - Key insights and conclusions drawn from the analysis.

  Ensure that the notebook is easy to follow, with explanations and comments to guide readers through the entire process.

- **Clear and concise presentation**: Prepare a professional presentation that summarizes the key findings, insights, and improvements achieved throughout the project. The presentation should focus on:

  - Key steps taken in the data processing and model building.
  - Major results from model performance and comparisons.
  - Highlights of any significant insights or improvements.
  - Visualization of results and performance metrics.
  - Lessons learned and potential future improvements.

  Keep the presentation concise, clear, and visually engaging, ensuring that the core achievements are effectively communicated.

# 6 Learning Outcomes:

- Proficiency in Python for data analysis, machine learning, and cybersecurity tasks.

- Practical experience working with real-world IoT network traffic datasets for intrusion detection.

- Data visualization skills for effectively communicating network behaviors and intrusion detection results.

- Ability to apply scalable techniques for analyzing large IoT datasets and detecting anomalies.

- Enhanced problem-solving and critical thinking skills in the context of cybersecurity and IoT systems.

# 7 conclusion

By the end of this project, you will have the skills and confidence to tackle real-world big data and machine learning challenges in cybersecurity, particularly within the context of Intrusion Detection Systems for IoT environments.

# Good luck