# *Abstract*

## *Anomaly Detection in Network Intrusion Data using Unsupervised Learning*

In the domain of cybersecurity, the timely detection of anomalous network activities is paramount for safeguarding digital assets. This project addresses the critical challenge of identifying unusual and potentially malicious patterns within network traffic by leveraging **unsupervised learning**, specifically the **Isolation Forest** algorithm. The **UNSW-NB15 dataset**, a comprehensive collection of real modern normal activities and synthetic attack behaviors, serves as the foundation for this study.

The methodology encompasses rigorous data preprocessing, including handling missing values, encoding categorical features, and scaling numerical attributes to prepare the data for modeling. The Isolation Forest model is then trained on this clean, unlabeled dataset to identify outliers based on their unique characteristics. The project involves predicting anomaly scores and binary anomaly labels, followed by an in-depth analysis of the detected anomalies. This includes examining the distribution of anomaly scores, presenting samples of flagged connections, and performing a crucial qualitative assessment against the dataset's original ground truth labels to understand which known attack categories, as well as 'normal' instances, are identified as anomalous. Visualizations, such as anomaly score histograms and feature-space scatter plots, further aid in interpreting the model's findings.

The outcome of this project demonstrates the Isolation Forest's efficacy in autonomously uncovering both known network intrusions and previously unseen, unusual behaviors that warrant further investigation. This capability is vital for enhancing proactive threat intelligence and potentially detecting zero-day attacks in real-world network security operations