# Automated Alert Classification Using Machine learning and NLP

## Problem Specification

Security Operations Centers (SOC) and Network Operation Center (NOC) are inundated with high volumes of alerts, which causes analysts to experience alert fatigue and struggle with timely incident response. The current reliance on manual triage processes amplifies false positive rates and risks critical security incidents being overlooked. Existing automated triage tools often fail to integrate unstructured alert descriptions with structured metadata, thereby missing opportunities to enhance classification accuracy through Natural Language Processing (NLP).As stated in Figure 1.1 Majority of NOC/SOC alerts are unable to handle and insufficient information also with Estimated false positive rate 40 % and false negative rate (failure to take action) 30 %
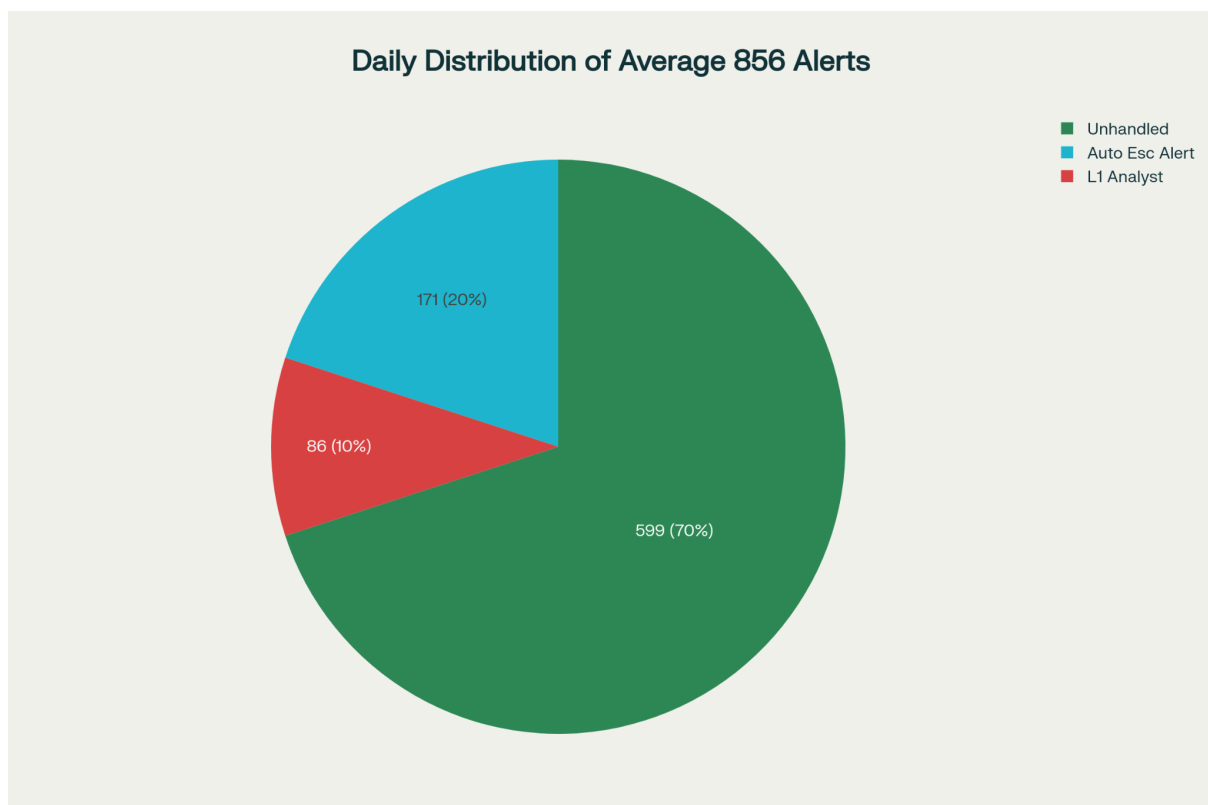


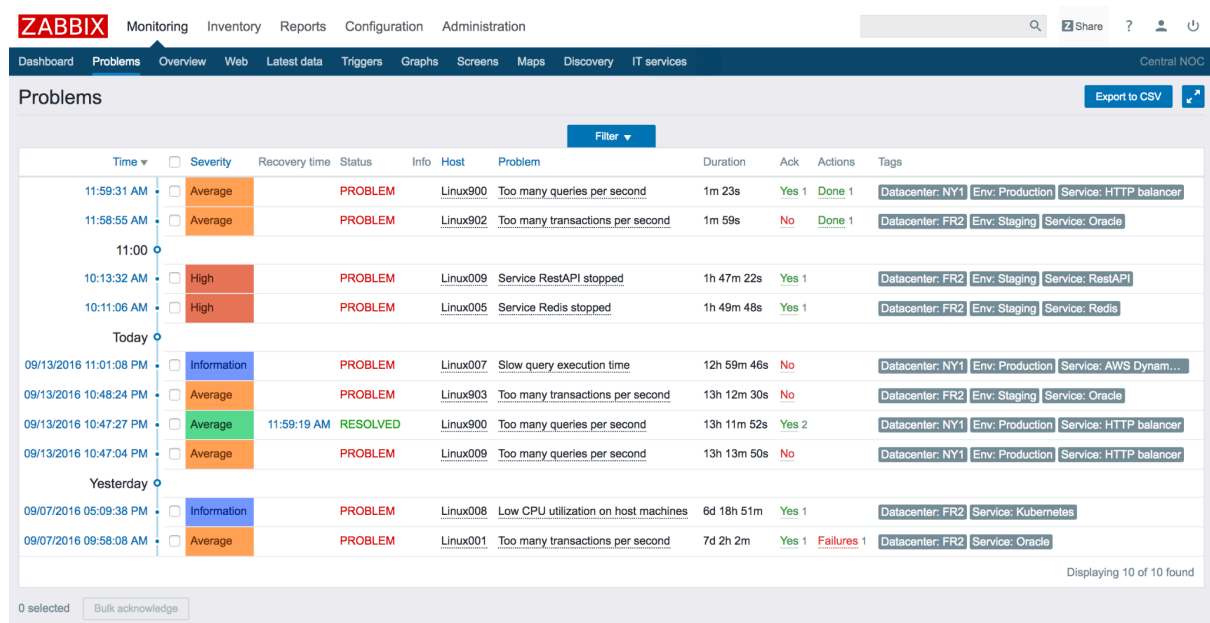Figure 1.1 Daily Alert Distribution of NOC/SOC

Figure 1.2 Typical events from monitoring system

Figure 1.2 explain traditional alerts are base on objects and trigger on certain threshold which work well on time series data . The **objective of this project is to develop a hybrid system that combines traditional machine learning with NLP to automatically classify cybersecurity and performance problem alerts**. By leveraging both structured features—such as severity, timestamp, source and destination IP addresses, user and system entities and NLP‑derived features—such as alert title and description embeddings, TF-IDF vectors, named entities, and sentiment or urgency scores—this solution aims to reduce analyst workload while preserving high accuracy and low false negative rates.

**Data Source**

Data will be sourced through 2 possible approach. Option 1 Public datasets such as the AIT alert collection, CICIDS network intrusion logs, the NapierOne threat indicator repository, will provide a foundational corpus. Option 2, partnerships with private NOC providers and managed SOC services will supply anonymized,

**Sample Data**

This system will process a stratified sample of approx 2,000 alerts drawn from enterprise SOC/NOC environments,. Stratification will ensure balanced representation across critical, high, medium, and low-severity categories. The model's output will consist of a severity prediction—categorized as Minor, Major, or Critical—accompanied by probability scores.

| Feature | Encoding Method |
|---------|-----------------|
| Time | Extract → Hour, DayOfWeek, Month; |
| Severity | Ordinal Encoding (0=Info→3=Average) |
| Status | One-Hot Encoding |
| Host | Label Encoding (Sensitive Data) |
| Problem | TF-IDF (max 100) with NER |

Figure 1.3 Alert Features and target encoding method

**Expectation**

The alert-classification model—built on static features and enhanced with basic NLP signals—can be expected. with precision of .70−.80, while a recall of .70 This operating point corresponds to an overall F1-score of approximately 0.65−0.75,