



SAKARYA ÜNİVERSİTESİ

Fall 2025-2026
SWE-423 Data Analytics
Asst. Prof. Nur Banu OĞUR

Real-time Network Intrusion Detection System with Big Data
Technologies

Görkem ÇETINKAYA - B211202022

1. Summary of the Project

This project aims to establish a scalable, real-time data analytics system for Network Intrusion Detection (NIDS). The primary objective is to simulate a streaming architecture where network traffic data is ingested, processed, analyzed, and stored instantaneously using Big Data technologies.

The system is designed to distinguish between "Normal" traffic and "Anomalous" (Attack) traffic. It leverages **Apache Kafka** for message distribution, **Apache Spark** for real-time stream processing and machine learning, and **MongoDB** for raw data persistence. The entire infrastructure is containerized using **Docker** to ensure reproducibility and scalability.

2. Technologies Used

The following technologies were utilized to build the end-to-end pipeline:

- **Apache Kafka & Zookeeper:** Used as the backbone for data ingestion and message distribution. It decouples the data source from the processing engine.
- **Apache Spark (Structured Streaming):** Used for processing the continuous stream of data.
- **Spark MLlib:** Used to build the Machine Learning pipeline (Logistic Regression) for classifying network traffic.
- **MongoDB:** A NoSQL database used to store the raw network traffic data for archival and future auditing.
- **Docker & Docker Compose:** Used to orchestrate the services (Broker, Zookeeper, Spark Master/Worker, MongoDB) in an isolated environment.
- **Python (PySpark, Kafka-Python):** The primary programming language used for Producer, Consumer, and Processor scripts.

3. Data Analysis Methods Used

- **Data Ingestion:** A Producer script simulates real-time traffic by reading the **KDDTest+** dataset and pushing messages to the Kafka topic **network-traffic**.
- **Pre-processing Pipeline:** The raw CSV data streaming from Kafka is parsed into a structured format. Categorical features (Protocol, Service, Flag) are converted into numerical indices using **StringIndexer**. All features are vectorized using **VectorAssembler**.
- **Machine Learning Model:** A Supervised Learning approach is implemented. A **Logistic Regression** model is trained on a static dataset (**KDDTrain+**) to learn attack patterns.
- **Real-time Prediction & Evaluation:** The trained model is applied to the live data stream. The system calculates **Batch Accuracy** (for the current micro-batch) and **Global Accuracy** (cumulative performance) in real-time to monitor system reliability.

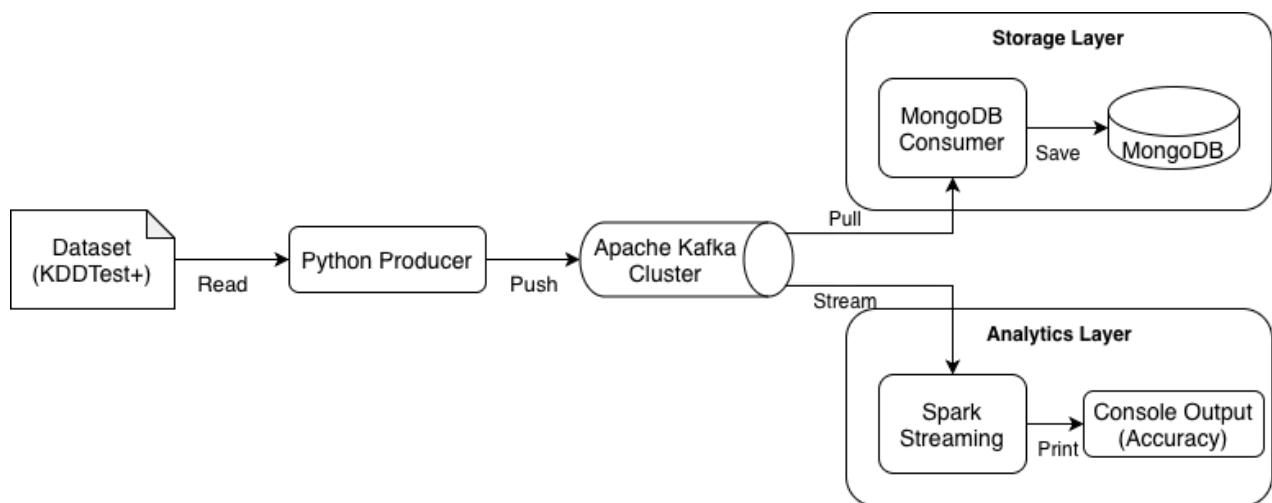
4. Description of the Dataset

The project utilizes the **NSL-KDD** dataset, which is a refined version of the benchmark KDD'99 dataset for intrusion detection.

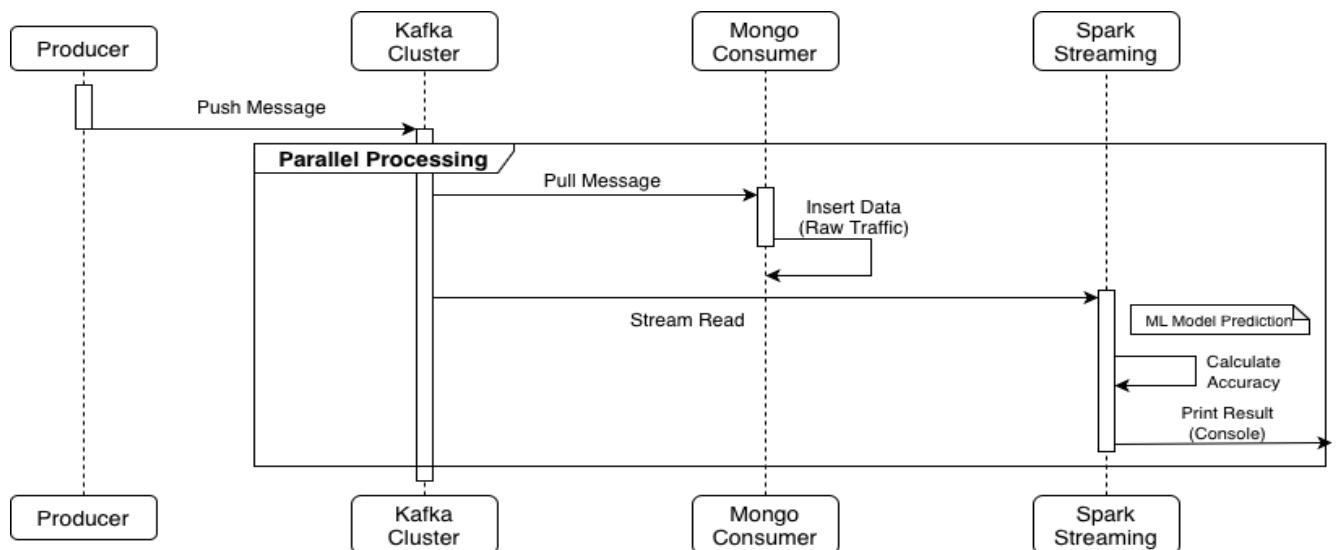
- **Training Set:** `KDDTrain+.txt` was used to train the Logistic Regression model.
- **Testing Set:** `KDDTest+.txt` was used for the streaming simulation to ensure the model is tested on unseen data (preventing data leakage).
- **Features:** The dataset includes basic network features such as `duration`, `protocol_type` (TCP/UDP/ICMP), `service`(HTTP, FTP, etc.), `src_bytes`, and `dst_bytes`.
- **Labels:** The traffic is labeled as either `normal` or specific attack types (e.g., `neptune`, `satan`, `smurf`).

5. General Flow Chart of the Project

The architecture follows a Producer-Consumer model integrated with a Stream Processor.



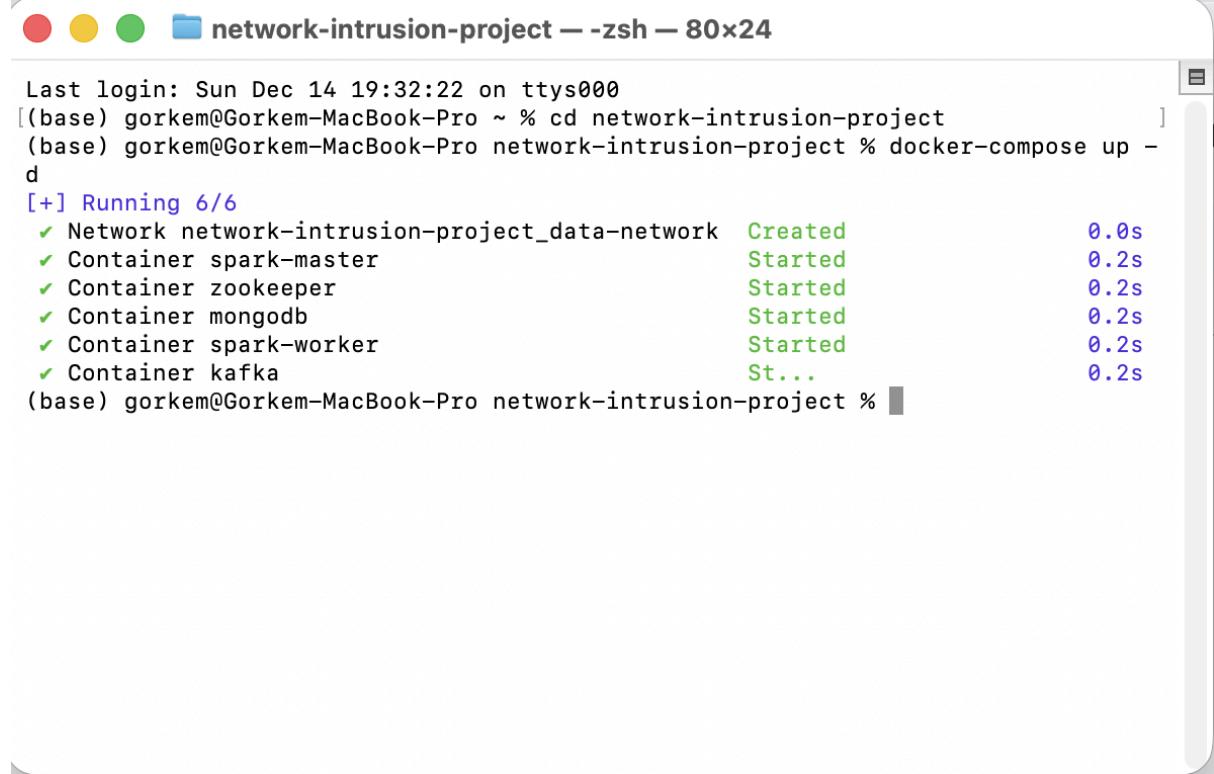
6. Timing Chart (Sequence Diagram)



7. Findings Obtained

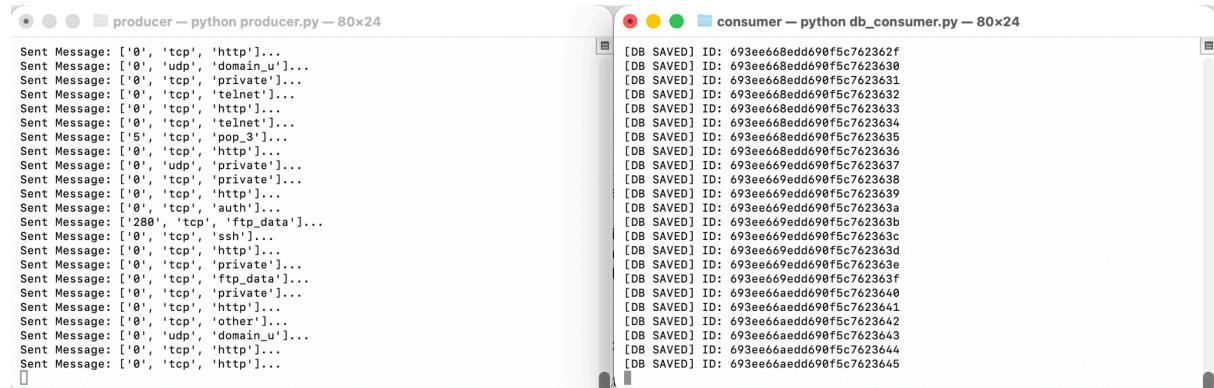
The system was successfully deployed and tested. The following findings were obtained:

- **Infrastructure:** All Docker containers (Zookeeper, Kafka, Spark, MongoDB) operated without errors.



```
Last login: Sun Dec 14 19:32:22 on ttys000
[(base) gorkem@Gorkem-MacBook-Pro ~ % cd network-intrusion-project
(base) gorkem@Gorkem-MacBook-Pro network-intrusion-project % docker-compose up -
d
[+] Running 6/6
✓ Network network-intrusion-project_data-network    Created          0.0s
✓ Container spark-master                          Started          0.2s
✓ Container zookeeper                           Started          0.2s
✓ Container mongodb                            Started          0.2s
✓ Container spark-worker                         Started          0.2s
✓ Container kafka                             St...            0.2s
(base) gorkem@Gorkem-MacBook-Pro network-intrusion-project %
```

- **Data Streaming & Persistence:** The Producer successfully streamed data, and the Consumer stored over 36,000 records to MongoDB, proving the system's data integrity.



```
Sent Message: ['0', 'tcp', 'http']...
Sent Message: ['0', 'udp', 'domain_u']...
Sent Message: ['0', 'tcp', 'private']...
Sent Message: ['0', 'tcp', 'telnet']...
Sent Message: ['0', 'tcp', 'http']...
Sent Message: ['0', 'tcp', 'telnet']...
Sent Message: ['5', 'tcp', 'pop_3']...
Sent Message: ['0', 'tcp', 'http']...
Sent Message: ['0', 'udp', 'private']...
Sent Message: ['0', 'tcp', 'private']...
Sent Message: ['0', 'tcp', 'http']...
Sent Message: ['0', 'tcp', 'auth']...
Sent Message: ['280', 'tcp', 'ftp_data']...
Sent Message: ['0', 'tcp', 'ssh']...
Sent Message: ['0', 'tcp', 'http']...
Sent Message: ['0', 'tcp', 'private']...
Sent Message: ['0', 'tcp', 'ftp_data']...
Sent Message: ['0', 'tcp', 'private']...
Sent Message: ['0', 'tcp', 'http']...
Sent Message: ['0', 'tcp', 'other']...
Sent Message: ['0', 'udp', 'domain_u']...
Sent Message: ['0', 'tcp', 'http']...
Sent Message: ['0', 'tcp', 'http']...
```

```
[DB SAVED] ID: 693ee668edd690f5c762362f
[DB SAVED] ID: 693ee668edd690f5c7623630
[DB SAVED] ID: 693ee668edd690f5c7623631
[DB SAVED] ID: 693ee668edd690f5c7623632
[DB SAVED] ID: 693ee668edd690f5c7623633
[DB SAVED] ID: 693ee668edd690f5c7623634
[DB SAVED] ID: 693ee668edd690f5c7623635
[DB SAVED] ID: 693ee668edd690f5c7623636
[DB SAVED] ID: 693ee668edd690f5c7623637
[DB SAVED] ID: 693ee669edd690f5c7623638
[DB SAVED] ID: 693ee669edd690f5c7623639
[DB SAVED] ID: 693ee669edd690f5c762363a
[DB SAVED] ID: 693ee669edd690f5c762363b
[DB SAVED] ID: 693ee669edd690f5c762363c
[DB SAVED] ID: 693ee669edd690f5c762363d
[DB SAVED] ID: 693ee669edd690f5c762363e
[DB SAVED] ID: 693ee669edd690f5c762363f
[DB SAVED] ID: 693ee669edd690f5c7623640
[DB SAVED] ID: 693ee669aed690f5c7623641
[DB SAVED] ID: 693ee669aed690f5c7623642
[DB SAVED] ID: 693ee669aed690f5c7623643
[DB SAVED] ID: 693ee669aed690f5c7623644
[DB SAVED] ID: 693ee669aed690f5c7623645
```

- **Real-time Analytics:** The Spark engine successfully processed micro-batches. The Logistic Regression model achieved a **Global Accuracy of ~82.97%** on the test stream. The model correctly identified various attack types (e.g., **neptune**, **warezmaster**) and distinguished them from normal traffic.

```

network-intrusion-project — docker exec -u 0 -it spark-master /opt/spark/...
=====
BATCH ID: 155
-----
RECORDS IN BATCH      : 48
BATCH ACCURACY       : %87.50
-----
TOTAL RECORDS        : 7416
GLOBAL ACCURACY      : %82.97
=====
+-----+-----+-----+-----+
|protocol_type|service |label_raw    |prediction|is_correct|
+-----+-----+-----+-----+
|tcp        |finger   |normal      |0.0       |1
|tcp        |http     |normal      |0.0       |1
|tcp        |http     |normal      |0.0       |1
|tcp        |pop_3    |guess_passwd|1.0       |1
|tcp        |pop_3    |guess_passwd|1.0       |1
|icmp      |ecr_i    |smurf      |1.0       |1
|tcp        |private   |neptune    |1.0       |1
|tcp        |ftp_data |warezmaster|0.0       |0
|tcp        |ftp      |guess_passwd|0.0       |0
|tcp        |ftp_data |neptune    |1.0       |1
+-----+-----+-----+-----+
only showing top 10 rows

```

- **Database Verification:** Verification scripts confirmed that data is actively being written to MongoDB.

```

network-intrusion-project — zsh — 80x24
=====
Last login: Sun Dec 14 19:32:37 on ttys000
[(base) gorkem@Gorkem-MacBook-Pro ~ % cd network-intrusion-project
[(base) gorkem@Gorkem-MacBook-Pro network-intrusion-project % python verify_mongo
.py
-----
DATABASE STATUS CHECK
-----
Database Name      : network_data
Collection Name   : raw_traffic
TOTAL RECORDS     : 36648
-----
>>> LATEST RECORD SAMPLE:
{'_id': ObjectId('693ee68fedd690f5c76237aa'), 'csv_data': ['0', 'udp', 'domain_u
', 'SF', '44', '134', '0', '0', '0', '0', '0', '0', '0', '0', '0', '0', '0', '0
', '0', '0', '0', '15', '15', '0.00', '0.00', '0.00', '0.00', '1.00', '0.00
', '0.00', '109', '37', '0.34', '0.03', '0.01', '0.00', '0.00', '0.00', '0.00
', '0.00', 'normal', '21']}
(base) gorkem@Gorkem-MacBook-Pro network-intrusion-project %

```