# Stream Data Analytics over Data Lake Presentation

**Submitted By**

Rajat Dutt Sharma & Prasoon Dwivedi

**Project Guide**

Prof. R. Chandrashekhar

Mr. Vivek Yadav
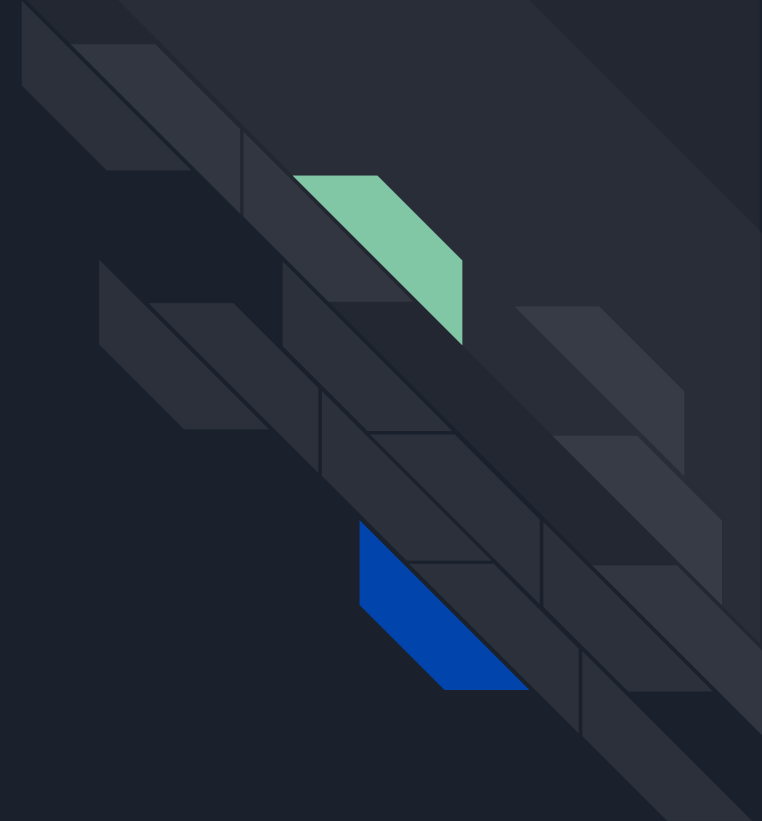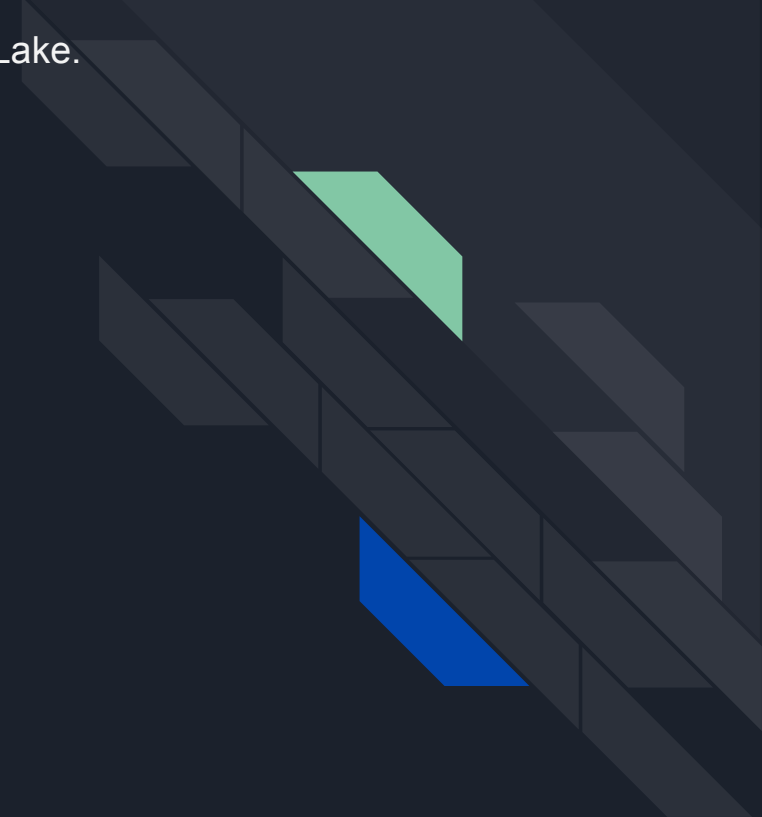
# Presentation Overview

# Overview

- Enable ingestion of Stream Data on the Data Lake.
- Perform Stream Data Analytics.

# Project Goals

- Perform Stream Data Analytics using :
  - Apache Flink
  - Apache Spark

  on the BMTC gprs trace data to calculate average speed of BMTC buses at any point of time.

# Challenges

- Develop data emitter to generate Stream Data from the static data.
- Select Tools for :
  - Stream Data Ingestion
  - Stream Data Analytics
- Develop Consumers for Stream Data :
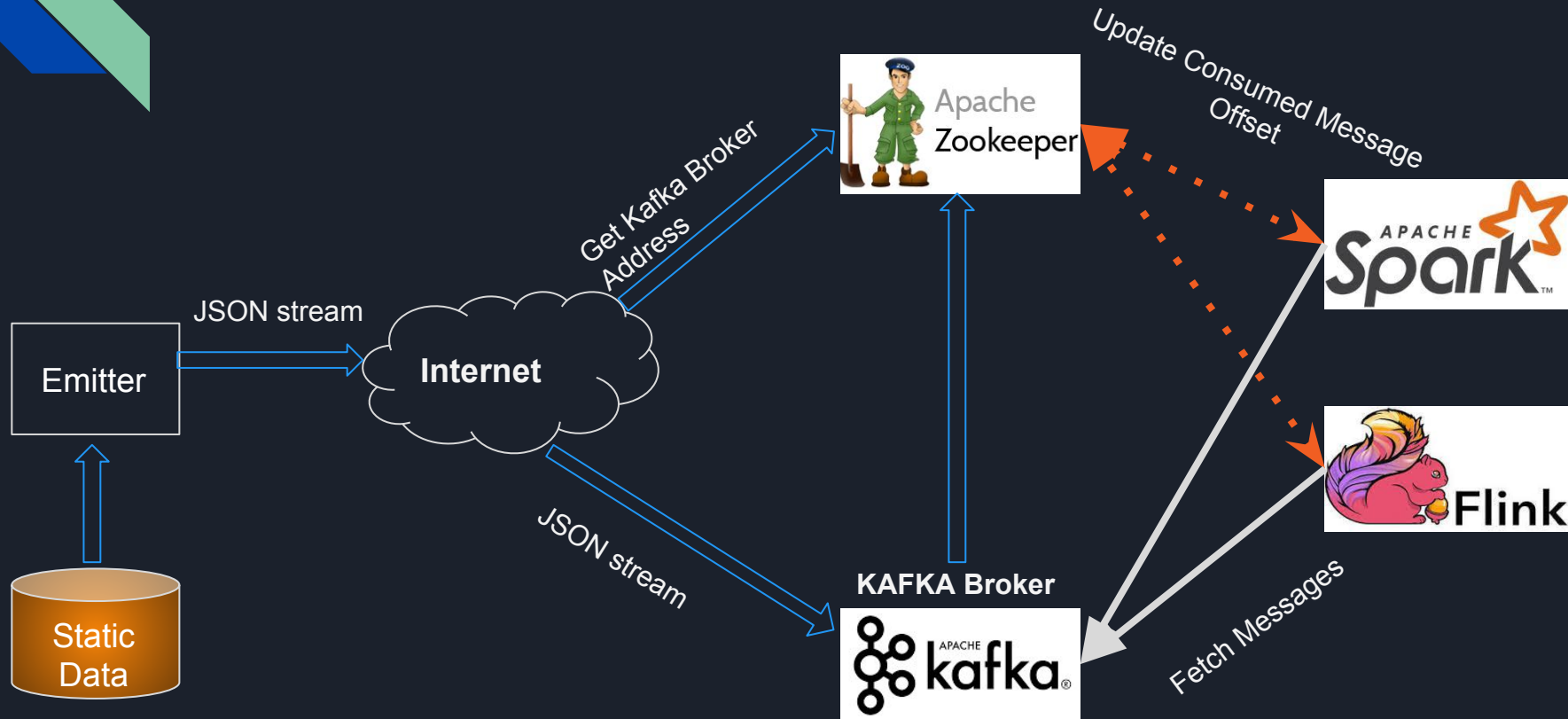  - Spark Consumer
  - Flink Consumer

# Work Overview

- Emitter :
    - Streams data at the same rate as the original data stream by mimicking the time difference between consecutive stream records from the timestamp of the records.
- The Data Analytics Flow in the consumers :
    - The stream data is first converted to object stream .
    - The object stream is then mapped to Sliding Keyed Window of size 10 for vehicle specific calculation of Average Speed.

# Architecture

# Tools

- Stream Data Ingestion
  - Apache Flume Vs Apache Kafka
    - Choice Made : Apache Kafka
      - API's to work with Spark And Flink.
      - Stream of records can be categorised into topics.
- Stream Data Analytics
  - Apache Spark Vs Apache Flink
  - Both have been tried

# UI
# (Apache Flink)

# UI
# (Apache Spark)



**Status:** SUCCEEDED
**Completed Stages:** 2
**Skipped Stages:** 1

▸ Event Timeline
▾ DAG Visualization

| Stage 16 (skipped) | Stage 17 | Stage 18 |

kafka direct stream [0]
@ 16:24:00

map @ 16:24:00

mapPartitions
@ 16:24:00

kafka direct stream [0]
@ 16:24:05

map @ 16:24:05

mapPartitions
@ 16:24:05

groupByKeyAndWindow
@ 16:24:00

groupByKeyAndWindow
@ 16:24:05

mapValues @ 16:24:05

flatMap @ 16:24:05

# UI
# (Apache Spark)

# Technologies Used

- Core Processing Modules : Java for Spark and Flink consumer programs.

- Tools Used :

    - Apache Kafka

    - Apache Flink

    - Apache Spark

    - Apache Zookeeper

- Kafka APIs for Spark and Flink

thank you!