# Final Project Part 3: Final Report Submission

**Team 3**
Panshul Saraswat
Nikhil Nagendra
Raj Mehta
Vyshnavi Gokavarapu
Rutuja Lohakare

**Problem Statement:**
In professional basketball, particularly within the NBA, real-time data analytics are pivotal for enhancing game analysis, optimizing player performance, and boosting fan engagement. Despite the potential of live data, its effective utilization is often hampered by the complexity of processing and visualizing this data in real-time. Our project accesses live NBA data using the nba_api, which provides a rich set of game statistics through endpoints like https://nba-prod-us-east-1-mediaops-stats.s3.amazonaws.com/NBA/liveData/scoreboard/todaysScoreboard_00.json. The API outputs data in JSON format, encapsulating game details such as scores, player stats, and team performance metrics across various game periods.

However, challenges persist in swiftly ingesting, processing, and displaying these data points. The structure of the data includes nested elements under keys like games, homeTeam, and awayTeam, detailing real-time updates like game scores, individual player performances, and game status. The goal of our proposed data pipeline—integrating Apache Kafka, Apache Spark, InfluxDB, and Grafana—is to streamline these processes. By setting up Kafka for efficient data ingestion, using Spark for real-time data analytics, storing processed data in InfluxDB, and visualizing insights through Grafana, we aim to facilitate quick, informed decision-making for coaches and teams, and enhance the interactive experience for fans. This pipeline promises to revolutionize how stakeholders interact with live game data, turning raw stats into actionable insights instantaneously.

## Task 1: Data Ingestion

Apache Kafka is the platform used for handling real-time data feeds. It is a distributed event streaming platform capable of handling trillions of events a day.

Apache Kafka Requires Java to run, as Kafka itself is a Java-based system.

### Install Java:
When setting up a data pipeline involving Apache Kafka, Java is a necessary dependency. Go to Oracle's official Java download page. Choose the version of Java Development Kit (JDK) you need. For most users working with Kafka and Spark, JDK 22 or the latest version will be suitable. Run the installer and check if the java has been installed or not from the terminal.

echo $JAVA_HOME
java -version



```
(base) nik@Nikhils-MacBook-Air-3 ~ % echo $JAVA_HOME
[java -version
/Library/Java/JavaVirtualMachines/jdk-22.jdk/Contents/Home
java version "22.0.1" 2024-04-16
Java(TM) SE Runtime Environment (build 22.0.1+8-16)
Java HotSpot(TM) 64-Bit Server VM (build 22.0.1+8-16, mixed mode, sharing)
(base) nik@Nikhils-MacBook-Air-3 ~ %
```
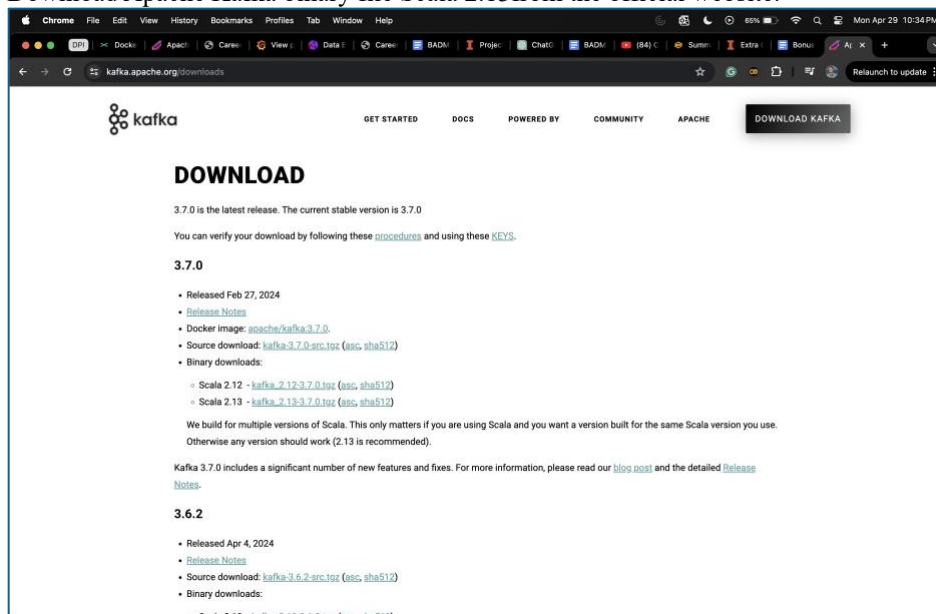
**Install Zookeeper:**

Go to the Apache Zookeeper releases page to download the latest stable release of Zookeeper. Look for the binary archive, which typically has an extension .tar.gz. Extract the file and have it ready.



**Install Kafka:**

Download Apache Kafka binary file Scala 2.13from the official website.



Once the file is downloaded place the file in a folder. Extract the file in the same folder by giving the following command on the terminal.

*tar -xzvf kafka_2.13-3.7.0.tar.gz*

**Kafka Setup:**

Next objective is to create a virtual environment (venv) so that the Kafka server can be run in it.

After the successful creation of the environment venv navigate to the extracted Kafka folder in the venv.



Now we run the below script to launch the zookeeper which is required to run Kafka.

*bin/zookeeper-server-start.sh config/zookeeper.properties*

After the above command is run in the terminal successfully, below results will be displayed indicating the zookeeper service is running



**Note: Do not terminate the zookeeper terminal as we need to keep this service running**

Next step is run a kafka server, open another terminal and navigate to extracted Kafka directory and give the below command to run a script which will launch the Kafka server.

*bin/kafka-server-start.sh config/server.properties*

Once the command is run successfully, below results will be shown indicating Kafka server is now setup and is now ready for data to be sent to it.

In an another terminal deactivate virtual environment venv, as this terminal will be later used to run a python file.

Create a python script in which you provide the data source and the Kafka topic for it to be written to.

The kafka-python library, which is a Kafka client for Python, allows for easy interaction with the Kafka system. This library can be installed using below command.

*pip install kafka-python*

Below is the script which used to send NBA data from NBA api. This code also contains a scheduler which ensures that the data is updated every 60 seconds in the Kafka topic/server. The scheduler can be removed as well depending on the use case.

**Code:**

```python
from kafka import KafkaProducer
import json
import time
from nba_api.live.nba.endpoints import scoreboard
from apscheduler.schedulers.blocking import BlockingScheduler
import logging

# Set up logging
logging.basicConfig(level=logging.INFO)
logger = logging.getLogger(__name__)

# Kafka Configuration
kafka_broker = 'localhost:9092'
kafka_topic = 'nba-scoreboard-1' <topic name

# Create a Kafka producer
producer = KafkaProducer(
bootstrap_servers=[kafka_broker],
value_serializer=lambda x: json.dumps(x).encode('utf-8')
)

def fetch_nba_scores():
try:
# Fetch today's NBA scoreboard
games = scoreboard.ScoreBoard()
data = games.get_dict()  # Fetch the games data as a dictionary
return data
except Exception as e:
logger.error(f"Failed to fetch NBA scores: {e}")
return None

def send_data_to_kafka(data):
try:
if data:
producer.send(kafka_topic, data)
producer.flush()
logger.info("Data sent to Kafka successfully.")
else:
logger.info("No data to send.")
except Exception as e:
logger.error(f"Failed to send data to Kafka: {e}")

def fetch_and_send():
data = fetch_nba_scores()  # Fetch data from NBA API
send_data_to_kafka(data)   # Send the data to Kafka

def main():
scheduler = BlockingScheduler()
scheduler.add_job(fetch_and_send, 'interval', minutes=1)  # Adjust the interval as needed
try:
scheduler.start()
except (KeyboardInterrupt, SystemExit):
pass
```

```
if __name__ == "__main__":
    main()
```

Save the python file and run in the terminal as shown below.

Python kaf.py



Once the data is sent successfully you can see a print statement in the terminal:

        INFO:__main__:Data sent to Kafka successfully.



**Viewing the data that is sent to Kafka:**
Open a new terminal and ensure that you are in the same venv and then navigate to the Kafka folder.
Use the below command to run a script that shows the data that is there in the Kafka topic from beginning.

*bin/kafka-console-consumer.sh --topic <your_topic_name> --from-beginning --bootstrap-server localhost:9092*

Below you can see a screenshot of the NBA data that was sent to Kafka.

0429/OKCNOP", "gameStatus": 3, "gameStatusText": "Final", "period": 4, "gameClock": "", "gameTimeUTC": "2024-04-30T00:30:00Z", "gameEt": "2024-04-29T20:3
0:00Z", "regulationPeriods": 4, "ifNecessary": false, "seriesGameNumber": "Game 4", "gameLabel": "West - First Round", "gameSubLabel": "Game 4", "seriesT
ext": "OKC wins 4-0", "seriesConference": "West", "poRoundDesc": "First Round", "gameSubtype": "", "homeTeam": {"teamId": 1610612740, "teamName": "Pelica
ns", "teamCity": "New Orleans", "teamTricode": "NOP", "wins": 0, "losses": 4, "score": 89, "seed": 8, "inBonus": null, "timeoutsRemaining": 0, "periods":
[{"period": 1, "periodType": "REGULAR", "score": 21}, {"period": 2, "periodType": "REGULAR", "score": 22}, {"period": 3, "periodType": "REGULAR", "score
": 28}, {"period": 4, "periodType": "REGULAR", "score": 18}]}, "awayTeam": {"teamId": 1610612760, "teamName": "Thunder", "teamCity": "Oklahoma City", "te
amTricode": "OKC", "wins": 4, "losses": 0, "score": 97, "seed": 1, "inBonus": null, "timeoutsRemaining": 1, "periods": [{"period": 1, "periodType": "REGU
LAR", "score": 21}, {"period": 2, "periodType": "REGULAR", "score": 23}, {"period": 3, "periodType": "REGULAR", "score": 26}, {"period": 4, "periodType":
"REGULAR", "score": 27}]}, "gameLeaders": {"homeLeaders": {"personId": 202685, "name": "Jonas Valanciunas", "jerseyNum": "17", "position": "C", "teamTri
code": "NOP", "playerSlug": null, "points": 19, "rebounds": 13, "assists": 1}, "awayLeaders": {"personId": 1628983, "name": "Shai Gilgeous-Alexander", "j
erseyNum": "2", "position": "G", "teamTricode": "OKC", "playerSlug": null, "points": 24, "rebounds": 10, "assists": 3}}, "pbOdds": {"team": null, "odds":
0.0, "suspended": 0}}]}}
{"meta": {"version": 1, "request": "https://nba-prod-us-east-1-mediaops-stats.s3.amazonaws.com/NBA/liveData/scoreboard/todaysScoreboard_00.json", "time":
"2024-04-30 12:12:13.1213", "code": 200}, "scoreboard": {"gameDate": "2024-04-29", "leagueId": "00", "leagueName": "National Basketball Association", "g
ames": [{"gameId": "0042300155", "gameCode": "20240429/LALDEN", "gameStatus": 2, "gameStatusText": "Q4 4:07", "period": 4, "gameClock": "PT04M07.00S", "g
ameTimeUTC": "2024-04-30T02:00:00Z", "gameEt": "2024-04-29T22:00:00-04:00", "regulationPeriods": 4, "ifNecessary": false, "seriesGameNumber": "Game 5", "
gameLabel": null, "gameSubLabel": null, "seriesText": "DEN leads 3-1", "seriesConference": "West", "poRoundDesc": "First Round", "gameSubtype": "", "home
Team": {"teamId": 1610612743, "teamName": "Nuggets", "teamCity": "Denver", "teamTricode": "DEN", "wins": 3, "losses": 1, "score": 97, "seed": 2, "inBonus
": "0", "timeoutsRemaining": 2, "periods": [{"period": 1, "periodType": "REGULAR", "score": 28}, {"period": 2, "periodType": "REGULAR", "score": 22}, {"p
eriod": 3, "periodType": "REGULAR", "score": 31}, {"period": 4, "periodType": "REGULAR", "score": 16}]}, "awayTeam": {"teamId": 1610612747, "teamName": "
Lakers", "teamCity": "Los Angeles", "teamTricode": "LAL", "wins": 1, "losses": 3, "score": 95, "seed": 7, "inBonus": "0", "timeoutsRemaining": 2, "period
s": [{"period": 1, "periodType": "REGULAR", "score": 24}, {"period": 2, "periodType": "REGULAR", "score": 29}, {"period": 3, "periodType": "REGULAR", "sc
ore": 26}, {"period": 4, "periodType": "REGULAR", "score": 16}]}, "gameLeaders": {"homeLeaders": {"personId": 1629008, "name": "Michael Porter Jr.", "jer
seyNum": "1", "position": "F", "teamTricode": "DEN", "playerSlug": "michael-porter-jr", "points": 26, "rebounds": 4, "assists": 1}, "awayLeaders": {"pers
onId": 2544, "name": "LeBron James", "jerseyNum": "23", "position": "F", "teamTricode": "LAL", "playerSlug": "lebron-james", "points": 26, "rebounds": 7,
"assists": 10}}, "pbOdds": {"team": null, "odds": 0.0, "suspended": 0}}, {"gameId": "0042300104", "gameCode": "20240429/BOSMIA", "gameStatus": 3, "gameS
tatusText": "Final", "period": 4, "gameClock": "", "gameTimeUTC": "2024-04-29T23:30:00Z", "gameEt": "2024-04-29T19:30:00Z", "regulationPeriods": 4, "ifNe
cessary": false, "seriesGameNumber": "Game 4", "gameLabel": "East - First Round", "gameSubLabel": "Game 4", "seriesText": "BOS leads 3-1", "seriesConfere
nce": "East", "poRoundDesc": "First Round", "gameSubtype": "", "homeTeam": {"teamId": 1610612748, "teamName": "Heat", "teamCity": "Miami", "teamTricode":
"MIA", "wins": 1, "losses": 3, "score": 88, "seed": 8, "inBonus": null, "timeoutsRemaining": 0, "periods": [{"period": 1, "periodType": "REGULAR", "scor
e": 24}, {"period": 2, "periodType": "REGULAR", "score": 12}, {"period": 3, "periodType": "REGULAR", "score": 23}, {"period": 4, "periodType": "REGULAR",
"score": 29}]}, "awayTeam": {"teamId": 1610612738, "teamName": "Celtics", "teamCity": "Boston", "teamTricode": "BOS", "wins": 3, "losses": 1, "score": 1
02, "seed": 1, "inBonus": null, "timeoutsRemaining": 1, "periods": [{"period": 1, "periodType": "REGULAR", "score": 34}, {"period": 2, "periodType": "REG
ULAR", "score": 19}, {"period": 3, "periodType": "REGULAR", "score": 28}, {"period": 4, "periodType": "REGULAR", "score": 21}]}, "gameLeaders": {"homeLea
ders": {"personId": 1628389, "name": "Bam Adebayo", "jerseyNum": "13", "position": "C-F", "teamTricode": "MIA", "playerSlug": null, "points": 25, "reboun
ds": 17, "assists": 5}, "awayLeaders": {"personId": 1628401, "name": "Derrick White", "jerseyNum": "9", "position": "G", "teamTricode": "BOS", "playerSlu
g": null, "points": 38, "rebounds": 4, "assists": 3}}, "pbOdds": {"team": null, "odds": 0.0, "suspended": 0}}, {"gameId": "0042300144", "gameCode": "2024
0429/OKCNOP", "gameStatus": 3, "gameStatusText": "Final", "period": 4, "gameClock": "", "gameTimeUTC": "2024-04-30T00:30:00Z", "gameEt": "2024-04-29T20:3
0:00Z", "regulationPeriods": 4, "ifNecessary": false, "seriesGameNumber": "Game 4", "gameLabel": "West - First Round", "gameSubLabel": "Game 4", "seriesT
ext": "OKC wins 4-0", "seriesConference": "West", "poRoundDesc": "First Round", "gameSubtype": "", "homeTeam": {"teamId": 1610612740, "teamName": "Pelica
ns", "teamCity": "New Orleans", "teamTricode": "NOP", "wins": 0, "losses": 4, "score": 89, "seed": 8, "inBonus": null, "timeoutsRemaining": 0, "periods":
[{"period": 1, "periodType": "REGULAR", "score": 21}, {"period": 2, "periodType": "REGULAR", "score": 22}, {"period": 3, "periodType": "REGULAR", "score
": 28}, {"period": 4, "periodType": "REGULAR", "score": 18}]}, "awayTeam": {"teamId": 1610612760, "teamName": "Thunder", "teamCity": "Oklahoma City", "te
amTricode": "OKC", "wins": 4, "losses": 0, "score": 97, "seed": 1, "inBonus": null, "timeoutsRemaining": 1, "periods": [{"period": 1, "periodType": "REGU
LAR", "score": 21}, {"period": 2, "periodType": "REGULAR", "score": 23}, {"period": 3, "periodType": "REGULAR", "score": 26}, {"period": 4, "periodType":
"REGULAR", "score": 27}]}, "gameLeaders": {"homeLeaders": {"personId": 202685, "name": "Jonas Valanciunas", "jerseyNum": "17", "position": "C", "teamTri
code": "NOP", "playerSlug": null, "points": 19, "rebounds": 13, "assists": 1}, "awayLeaders": {"personId": 1628983, "name": "Shai Gilgeous-Alexander", "j
erseyNum": "2", "position": "G", "teamTricode": "OKC", "playerSlug": null, "points": 24, "rebounds": 10, "assists": 3}}, "pbOdds": {"team": null, "odds":
0.0, "suspended": 0}}]}}

This Completes Apache Kafka setup.

**Task 2: Data Processing, and Storing**

**Setup Apache Spark:**

PySpark is used to process the data streaming from Kafka. It handles reading, transforming, and extracting useful information from the Kafka topics. The PySpark library in Python, installed from the below command.

*pip install pyspark*

Spark-Kafka connector allows PySpark to read from and write to Kafka topics directly. This is crucial for integrating Kafka streams with Spark processing.
**Spark-SQL-Kafka-0-10 Package** must be specified in your Spark session to connect Spark with Kafka. It's included with Spark but needs to be explicitly enabled with configuration options when initializing the Spark session:

.config("spark.jars.packages", "org.apache.spark:spark-sql-kafka-0-10_2.12:<spark_version>")

Replace <spark_version> with your actual Spark version.

**Setup Influx DB:**

InfluxDB is a time-series database designed to handle high write and query loads. It is particularly useful in scenarios involving real-time analytics.
InfluxDB Client for Python: influxdb_client library, allows you to connect, write, and query data from an InfluxDB instance. This can be installed using following command.

*pip install influxdb-client*

Creating a bucket in InfluxDB can be done through the InfluxDB UI (User Interface). Open your web browser and navigate to the InfluxDB instance. The default URL is usually *http://localhost:8086*

Once Logged in we need to create and setup bucket in the influx db where the data is stored. As explained above install the dependency if not installed earlier.



Create and save the token, later to be used to find the bucket

Now we create organization and set up server url as they are needed to have an initial connection, here 'UIUC' is the name given to organization.



Click on create bucket, below are the 2 buckets created in this process, once done test data can be sent to check, and below images shows the same setup

## Screenshot 1

**Setting Up** Python

⏱ 5 minutes

- ✓ Overview
- ✓ Install Dependencies
- ✓ Get Token
- ✓ Initialize Client
- ✓ Write Data
- ▶ Execute a Simple Query
- ▶ Execute an Aggregate Query
- ☆ Finish

```
|> range(start: -10m)
```

In this query, we are looking for data points within the last 10 minutes with a measurement of "measurement1".

Let's use that Flux query in our Python code!

Run the following:

```python
query_api = client.query_api()

query = """from(bucket: "<BUCKET>")
 |> range(start: -10m)
 |> filter(fn: (r) => r._measurement == "measurement1")"""
tables = query_api.query(query, org="UIUC")

for table in tables:
  for record in table.records:
    print(record)
```

COPY TO CLIPBOARD

PREVIOUS    NEXT

---

## Screenshot 2

**Setting Up** Python

⏱ 5 minutes

- ✓ Overview
- ✓ Install Dependencies
- ✓ Get Token
- ✓ Initialize Client
- ✓ Write Data
- ✓ Execute a Simple Query
- ▶ Execute an Aggregate Query
- ☆ Finish

In this example, we use the `mean()` function to calculate the average value of data points in the last 10 minutes.

Run the following:

```python
query_api = client.query_api()

query = """from(bucket: "<BUCKET>")
   |> range(start: -10m)
   |> filter(fn: (r) => r._measurement == "measurement1")
   |> mean()"""
tables = query_api.query(query, org="UIUC")

for table in tables:
    for record in table.records:
        print(record)
```

COPY TO CLIPBOARD

This will return the mean of the five values. ( (0+1+2+3+4) / 5 = 2 )

PREVIOUS    NEXT

In the data explorer, flux query can be run to retrieve the data.



From this set up below are the things to be noted:
influxdb_token = "CGJpFEeqwj1XiZe0HXFyT38_UOCuxrX-qEhIx7RyA4KY8L5HlCN-pvvdghxStvXuoN56RmrIfW0N8JGyU2B1MA=="

influxdb_org = "UIUC"

influxdb_bucket = "n_score"

influxdb_url = http://localhost:8086

**Python Shell Code:**
Once all the above set up is done, create a new python script which consumes the data from our kafka topic consumer and defines schema for our data to convert from json to tabular format and send it to our influx db at every minute interval, below is the code:

[Yesterday 9:23 PM] Saraswat, Panshul

Spark code :

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, from_json, explode
from pyspark.sql.types import *
from influxdb_client import InfluxDBClient, Point, WriteOptions

# InfluxDB Client Configuration
url = "http://localhost:8086"
token = "CGJpFEeqwj1XiZe0HXFyT38_UOCuxrX-qEhIx7RyA4KY8L5HlCN-
pvvdghxStvXuoN56RmrIfW0N8JGyU2B1MA=="
org = "UIUC"
bucket = "n_score"
client = InfluxDBClient(url=url, token=token, org=org)
write_api = client.write_api(write_options=WriteOptions(batch_size=1000, flush_interval=10_000))

# Define the schema based on the nested JSON structure
schema = StructType([
    StructField("meta", StructType([
        StructField("version", IntegerType()),
        StructField("request", StringType()),
        StructField("time", StringType()),
        StructField("code", IntegerType())
    ])),
    StructField("scoreboard", StructType([
        StructField("gameDate", StringType()),
        StructField("leagueId", StringType()),
        StructField("leagueName", StringType()),
        StructField("games", ArrayType(StructType([
            StructField("gameId", StringType()),
            StructField("gameCode", StringType()),
            StructField("gameStatus", IntegerType()),
            StructField("gameStatusText", StringType()),
            StructField("period", IntegerType()),
            StructField("gameClock", StringType()),
            StructField("gameTimeUTC", StringType()),
            StructField("gameEt", StringType()),
            StructField("regulationPeriods", IntegerType()),
            StructField("ifNecessary", BooleanType()),
            StructField("seriesGameNumber", StringType()),
            StructField("gameLabel", StringType()),
            StructField("gameSubLabel", StringType()),
            StructField("seriesText", StringType()),
            StructField("seriesConference", StringType()),
```

```
StructField("poRoundDesc", StringType()),
StructField("gameSubtype", StringType()),
StructField("homeTeam", StructType([
    StructField("teamId", IntegerType()),
    StructField("teamName", StringType()),
    StructField("teamCity", StringType()),
    StructField("teamTricode", StringType()),
    StructField("wins", IntegerType()),
    StructField("losses", IntegerType()),
    StructField("score", IntegerType()),
    StructField("seed", IntegerType()),
    StructField("inBonus", StringType()),
    StructField("timeoutsRemaining", IntegerType()),
    StructField("periods", ArrayType(StructType([
        StructField("period", IntegerType()),
        StructField("periodType", StringType()),
        StructField("score", IntegerType())
    ])))
])),
StructField("awayTeam", StructType([
    StructField("teamId", IntegerType()),
    StructField("teamName", StringType()),
    StructField("teamCity", StringType()),
    StructField("teamTricode", StringType()),
    StructField("wins", IntegerType()),
    StructField("losses", IntegerType()),
    StructField("score", IntegerType()),
    StructField("seed", IntegerType()),
    StructField("inBonus", StringType()),
    StructField("timeoutsRemaining", IntegerType()),
    StructField("periods", ArrayType(StructType([
        StructField("period", IntegerType()),
        StructField("periodType", StringType()),
        StructField("score", IntegerType())
    ])))
])),
StructField("gameLeaders", StructType([
    StructField("homeLeaders", StructType([
        StructField("personId", IntegerType()),
        StructField("name", StringType()),
        StructField("jerseyNum", StringType()),
        StructField("position", StringType()),
        StructField("teamTricode", StringType()),
        StructField("playerSlug", StringType()),
        StructField("points", IntegerType()),
        StructField("rebounds", IntegerType()),
        StructField("assists", IntegerType())
    ])),
    StructField("awayLeaders", StructType([
        StructField("personId", IntegerType()),
        StructField("name", StringType()),
        StructField("jerseyNum", StringType()),
        StructField("position", StringType()),
        StructField("teamTricode", StringType()),
        StructField("playerSlug", StringType()),
        StructField("points", IntegerType()),
```

```python
                    StructField("rebounds", IntegerType()),
                    StructField("assists", IntegerType())
                ]))
            ])),
            StructField("pbOdds", StructType([
                StructField("team", StringType()),
                StructField("odds", DoubleType()),
                StructField("suspended", IntegerType())
            ]))
        ])))
    ]))
])


# Initialize Spark Session
spark = SparkSession.builder.appName("NBA Scores Kafka Consumer").getOrCreate()


# Read from Kafka
df = spark.readStream.format("kafka") \
    .option("kafka.bootstrap.servers", "localhost:9092") \
    .option("subscribe", "nba-scoreboard-1") \
    .load()


# Select message value and convert from bytes to string
df = df.selectExpr("CAST(value AS STRING)")


# Parse the JSON string using the defined schema
parsed_df = df.select(from_json(col("value"), schema).alias("data"))


# Explode the games array into individual game records
games_df = parsed_df.select(explode(col("data.scoreboard.games")).alias("game"))


# Select relevant fields for processing
flattened_df = games_df.select(
    col("game.gameId").alias("gameId"),
    col("game.gameCode").alias("gameCode"),
    col("game.gameStatusText").alias("gameStatusText"),
    col("game.homeTeam.teamName").alias("homeTeamName"),
    col("game.homeTeam.score").alias("homeScore"),
    col("game.awayTeam.teamName").alias("awayTeamName"),
    col("game.awayTeam.score").alias("awayScore"),
    col("game.period").alias("period"),
```

```python
        col("game.gameTimeUTC").alias("gameTimeUTC")
)




# Function to write each batch of DataFrame to InfluxDB
def write_to_influx(batch_df, epoch_id):
    points = []
    for row in batch_df.collect():
        point = Point("nba_games") \
            .tag("gameId", row.gameId) \
            .tag("gameCode", row.gameCode) \
            .field("gameStatusText", row.gameStatusText) \
            .field("homeTeamName", row.homeTeamName) \
            .field("homeScore", int(row.homeScore)) \
            .field("awayTeamName", row.awayTeamName) \
            .field("awayScore", int(row.awayScore)) \
            .field("period", int(row.period)) \
            .time(row.gameTimeUTC)
        points.append(point)
    write_api.write(bucket=bucket, record=points)




# Write the stream to InfluxDB
query = flattened_df.writeStream.outputMode("append") \
    .foreachBatch(write_to_influx) \
    .start()

query.awaitTermination()
```

Save the python file, and in new terminal run the spark job from the below command.

***spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.5.1 spark.py***

Here while sending the data, we look for numoutputRows= -1 which means that the data was successfully sent to InfluxDB

**Task 3: Visualize with Grafana**

**Setup Grafana:**
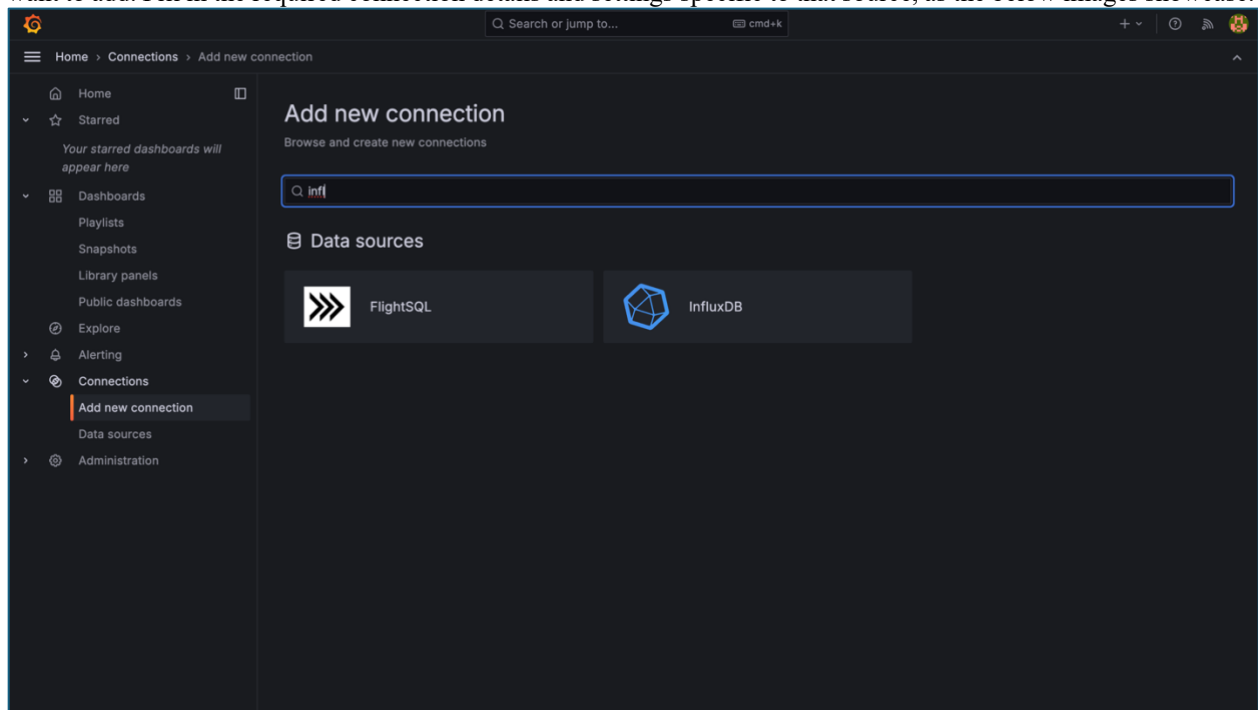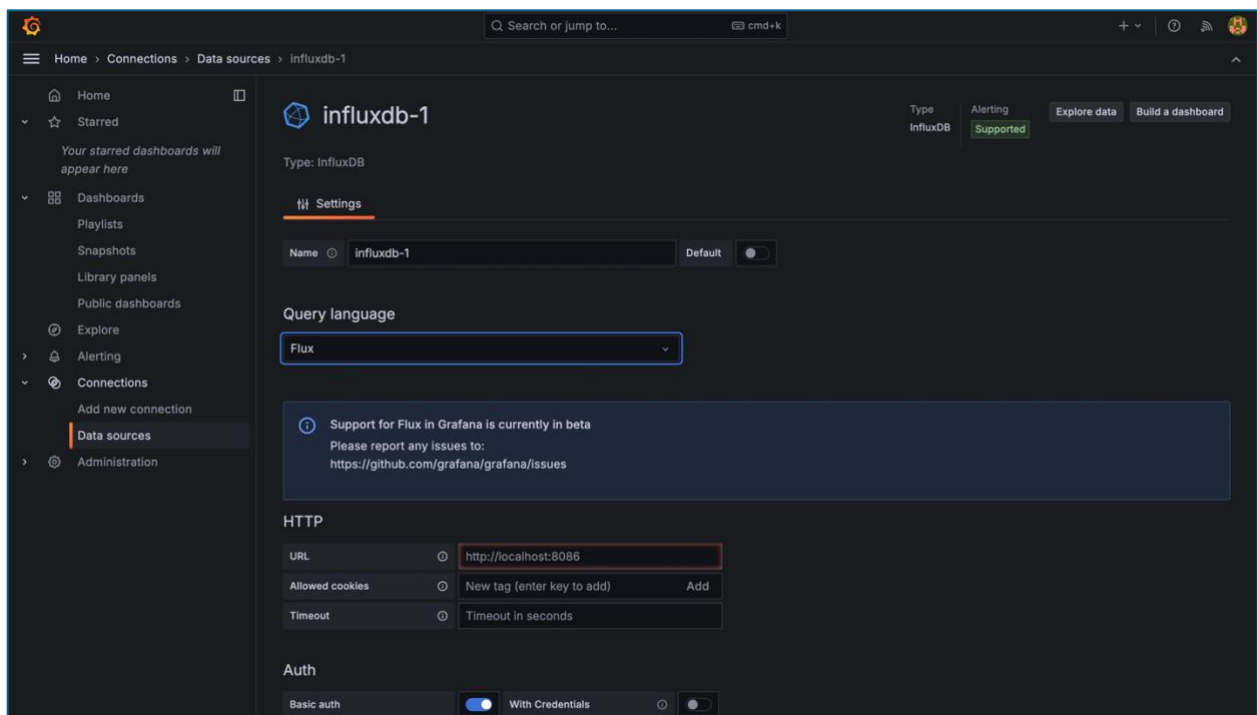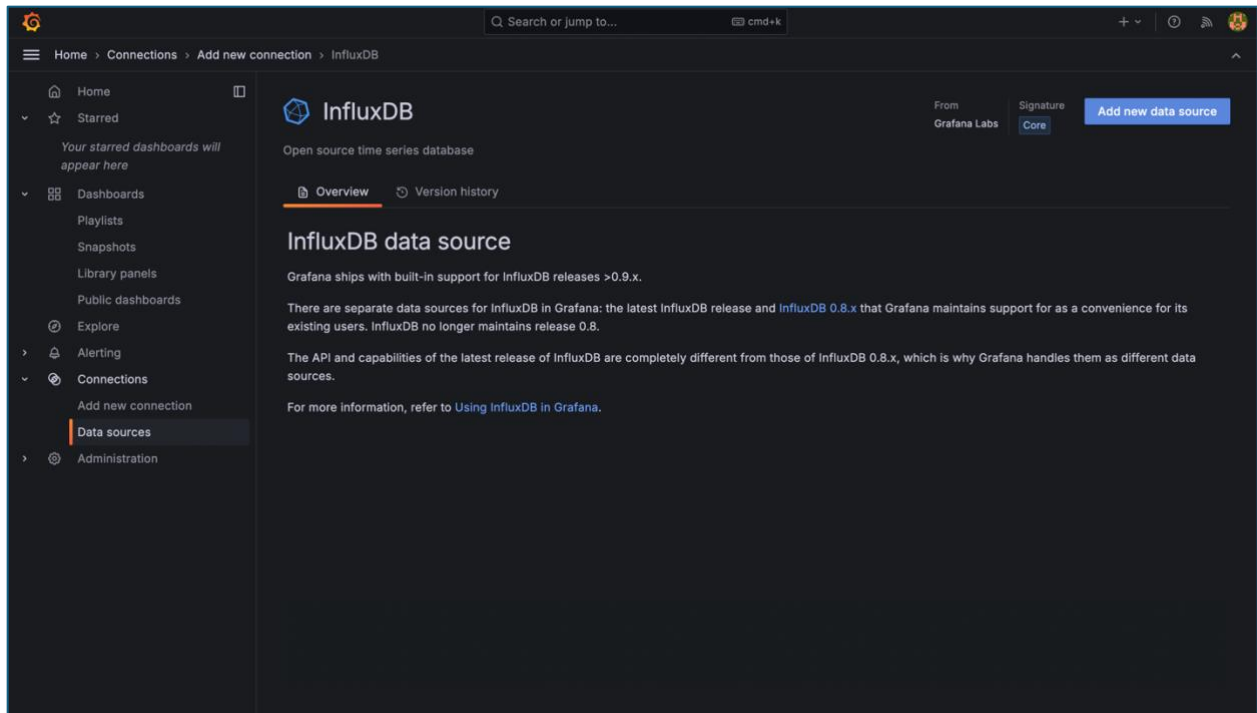Install and run Grafana by running below 2 commands in terminal

*brew install grafana*
*brew services start grafana*

Once Grafana is installed and running, open your web browser and navigate to http://localhost:3000/. The default port for Grafana is 3000.

The default login credentials are Username: admin and Password: admin. Upon first login, you will be prompted to change the password.

After logging in, you can configure data sources (like InfluxDB, MySQL, PostgreSQL, etc.) that Grafana will use to pull data for visualization. Go to Configuration > Data Sources > Add data source and select the type of source you want to add. Fill in the required connection details and settings specific to that source, as the below images showcase.

**Visulization:**
Once your data sources are configured, you can start creating dashboards to visualize your data.
Go to + Create > Dashboard and begin adding panels.

We can utilize flux query to fetch the data we need in our desired format to create and add visuals
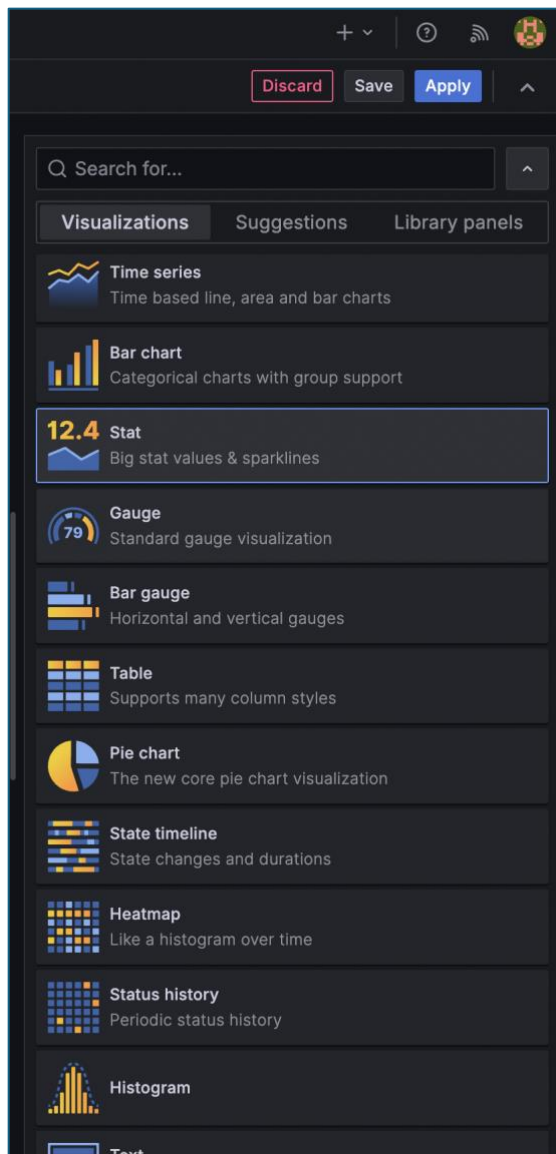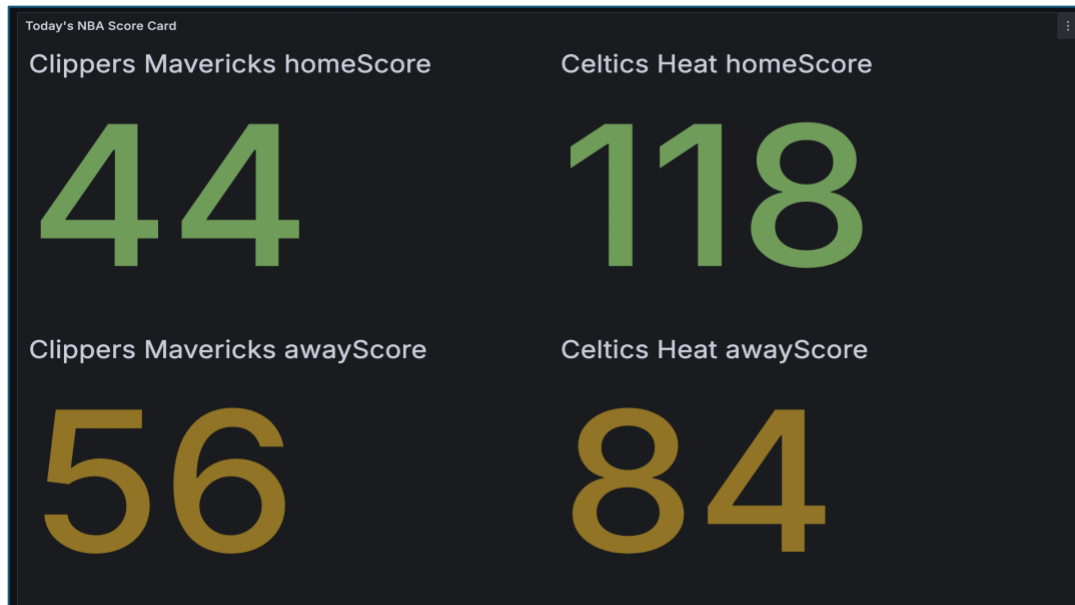
Below is the query used:

```
from(bucket: "n_score")
  |> range(start: -1d)
  |> filter(fn: (r) => r._field == "awayScore" or r._field == "homeScore" or r._field == "awayT
  |> pivot(rowKey:["_time"], columnKey: ["_field"], valueColumn: "_value")
  |> map(fn: (r) => ({
      _time: r._time,
      homeTeam: r.homeTeamName,
      homeScore: r.homeScore,
      awayTeam: r.awayTeamName,
      awayScore: r.awayScore
    })
  )
  |> sort(columns: ["_time"], desc: true)
  |> limit(n: 10)
```

After the query is run, select the desired visualization from the right hand side and configure it to your use case.

Once done Dashboard is ready, this dashboard updates the live score of the NBA game at every one-minute interval

**Today's NBA Score Card**

Clippers Mavericks homeScore

**44**

Celtics Heat homeScore

**118**

Clippers Mavericks awayScore

**56**

Celtics Heat awayScore

**84**

**Challenges:**

- o    New Technology Learning Curve
- o    Dependency Setup Challenges
- o    Complexities with GTFS Data
- o    Influx DB Data Storage and Querying
- o    Real-Time NBA Data Limitations