# PES UNIVERSITY

100 feet Ring Road, BSK 3rd Stage
Bengaluru 560085

Department of Computer Science and Engineering
B. Tech. CSE - 6th Semester
Jan – May 2023

## UE20CS343
## DATABASE TECHNOLOGIES (DBT)
# Project Report

# Performing Stream Processing and Batch Processing on tweets

## TEAM #: 5

**PES1UG20CS286:** Pragathya Baskar
**PES1UG20CS045:** Ananya Prakash
**PES1UG21CS810:** Divya M
**PES1UG20CS265:** Nidhi R Jois

# Class of Prof.  K S Srinivas

Pragathya
Ananya
Divya
Nidhi

# Table of Contents

Pragathya
Ananya
Divya
Nidhi

## 1. Introduction

- The task involves the use of various technologies and frameworks to process streaming data and run batch queries on the same data. The aim is to compare the performance and accuracy of processing the data in both streaming and batch modes.

- Apache Spark Streaming and Spark SQL will be used to execute multiple workloads on the input data. These workloads will include Spark SQL queries to perform actions, transformations, and aggregations on the input data.

- Apache Kafka Streaming will be used to publish and subscribe to the results or produce and consume from three or more topics. The data will be stored in a DBMS of choice such as Postgres or MySQL.

- The input data for the computation examples will be a streaming data source such as Twitter feed (tweets).

## 2.Installation of Software [include version #s and URLs]
➔ Streaming Tools Used
- Apache Spark Streaming:
  - Version # : Spark 3.4.0
  - URL: https://spark.apache.org/downloads.html
- Apache Kafka Streaming:
  - Version #: Kafka 3.4.0
  - URL: https://kafka.apache.org/downloads
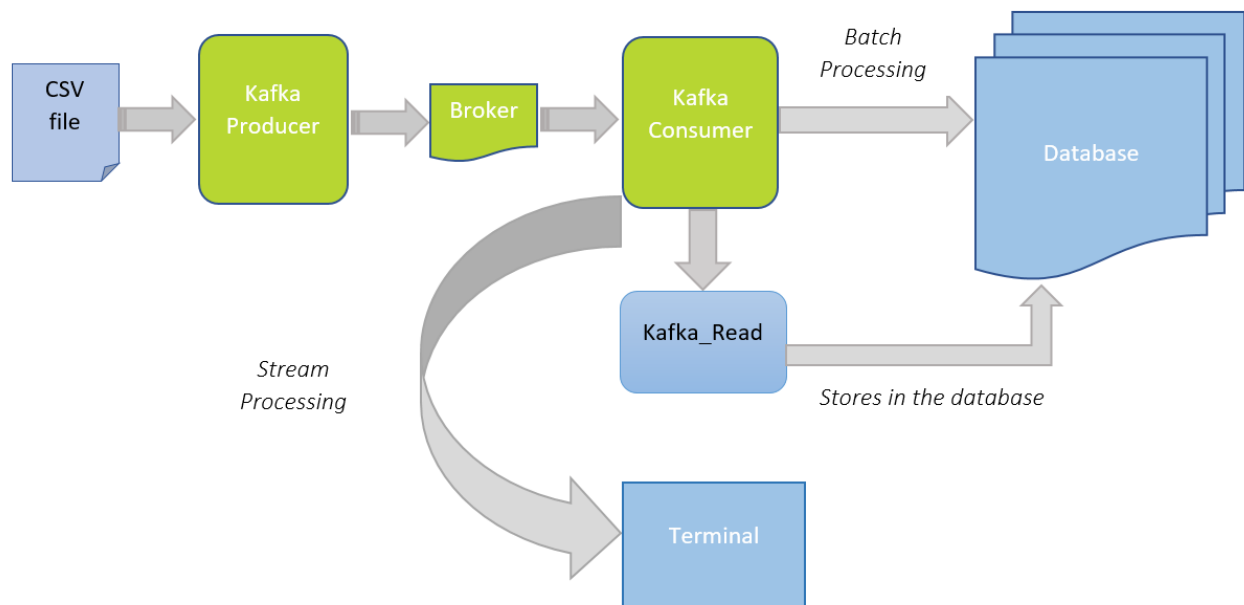
Pragathya
Ananya
Divya
Nidhi

➔DBMS Used:
◆MySQL database

# 3.Problem Description

- Apache Spark Streaming and Spark SQL is used to execute multiple workloads on the input data.
- These workloads will include Spark SQL queries to perform actions, transformations, and aggregations on the input data.

- Apache Kafka Streaming will be used to publish and subscribe to the results or produce and consume from three or more topics.

- The data will be stored in the MySQL database.

# 4.Architecture Diagram

Pragathya
Ananya
Divya
Nidhi

# 5. Input Data
- Source
  - Kaggle dataset
- Description
  - The dataset containing tweets is taken from a producer which publishes it to kafka as topic and value.
  - Then kafka will store the data in MySQL database after doing some preprocessing.
  - Then the consumer will subscribe to a topic and perform stream and batch processing.

# 6. Streaming Mode Experiment

Pragathya
Ananya
Divya
Nidhi

- Windows: Ubuntu
- The type of window used here: tumbling window
- Workloads:
    - For Spark, data processing tasks such as batch processing and streaming processing.
- Code like SQL Scripts : Spark SQL and Pyspark
    - Spark SQL code used in consumer_batch.py and consumer_stream.py
    - positive_df = spark.sql("SELECT * FROM tweets WHERE topic = 'positive'")

- Inputs and Corresponding Results
    - Input is a twitter dataset from which the producer read and publishes the topics to the kafka broker
    - The result is selecting all tweets of topic 'positive' and then counting the number of hashtags in the tweet

Input screenShot:

Dataset.xlsx

Pragathya
Ananya
Divya
Nidhi

| File Name | tweet | sentiment |
|-----------|-------|-----------|
| 1.txt | How I feel today #legday #jelly #aching #gym | negative |
| 10.txt | @ArrivaTW absolute disgrace two carriages from Bangor half way there standing room only #disgraced | negative |
| 100.txt | This is my Valentine's from 1 of my nephews. I am elated; sometimes the little things are the biggest & best things! | positive |
| 1000.txt | betterfeelingfilms: RT via Instagram: First day of filming #powerless back in 2011. Can't ¡ | neutral |
| 1001.txt | Zoe's first love #Rattled @JohnnyHarper15 | positive |
| 1002.txt | Chaotic Love - giclee print ?65 at #art #love #chaotic #abstract #blue #silver #prints #buy | positive |
| 1003.txt | They gna be mad when I reach that goal though. #Rejected the wrong girl ? just getting started & already turn heads.? | negative |
| 1004.txt | On day 9.. It's now in my daily routine.. Feeling guuuurdddd! ? #Aching #PainNoGain #FeelingGood | negative |
| 1005.txt | #ANIMALABUSE #TORONTO #PUPPY #TORTURE WE OFFER $1K #REWARD puppy #beaten #bound #burned | neutral |
| 1006.txt | Mike will not accept this plastic rose. @wfaamike @wfaachannel8 @wfaagmt #rejected | negative |
| 1007.txt | Just ate four cookies. #remorse | negative |
| 1008.txt | It's shocking what is acceptable in kids TV shows these days #shocking #shocked | negative |
| 1009.txt | We are so #excited to announce that we have launched our #affiliate program please visit us at | positive |
| 101.txt | Just when you thought you'd seen everything! .... #ParkingSpace #Lanzarote #ItsNOTtaxed #zimmer #oldpeople #alarmed ? | negative |
| 1010.txt | RT @MissGem: So this @parcelforce van thinks it's ok to nearly run people off the road?! #disgusting #shocked #disgraceful | negative |



| File Name | tweet | sentiment |
|-----------|-------|-----------|
| 1010.txt | RT @MissGem: So this @parcelforce van thinks it's ok to nearly run people off the road?! #disgusting #shocked #disgraceful | negative |
| 1011.txt | Today I #StepBackInTime !!! @PWLHitFactory @kylieminogue #BetterTheDevilYouKnow #WhatDoIHaveToDo #Shocked | neutral |
| 1012.txt | Photos: #Photographer got a rumble in the #jungle as he was #beaten by 30-stone #gorilla | neutral |
| 1013.txt | ¡°@Dreenayz: Eto day yung mukha ng stolen shot ?? @kineyerrrr feeling #shocked | negative |
| 1014.txt | @ThatKidRalph: I was going to be a child and violate butttt I'm grown #shook | negative |
| 1015.txt | @ThatKidRalph: I was going to be a child and violate butttt I'm grown #shook | negative |
| 1016.txt | Just sat down on the plane! #shook #a380 | neutral |
| 1017.txt | in #Spain #paris #entertainment #Portug¡ | neutral |
| 1018.txt | nephews ????????? #boys #brothers #caring #instafamous #f4f #l4¡ | positive |
| 1019.txt | followed by a forfeit... #shook #forfeit #QuakerGrit @PennWrestling @MikeSteltenkamp | neutral |
| 102.txt | RT @IDA_SINGAPORE: It's @TimDraper w @ChannelNewsAsia Tim hosted @steveleonardSG @DraperHeroCity & we r elated to host Tim @comebash http:/¡ | positive |
| 1020.txt | RT @GooleAFC: Tonight we recieved a donation of ?200 from our friends @GooleUnitedAFC towards the VPG defibrillator! #speechless | positive |
| 1021.txt | #SPEECHLESS. Deah & Yusor were just married in December. #ChapelHillShooting | neutral |
| | I'm #Speechless ... #ChapelHillShooting #MuslimLivesMatter | |

Pragathya
Ananya
Divya
Nidhi

Output Screenshot:

Pragathya
Ananya
Divya
Nidhi

Pragathya
Ananya
Divya
Nidhi

# The output of kafka_read.py



 it stores the dataset in the database which will be further used by batch processing

Pragathya
Ananya
Divya
Nidhi

# 7. Batch Mode Experiment

- Description:
    - Input to the consumer is through the database which has stored the processed tweets which were received from the producer.
- Data Size
    - 375 KB
- Results
    - The result is selecting all tweets of topic 'positive' and then counting the number of hashtags in the tweet
- Output ScreenShot

Pragathya
Ananya
Divya
Nidhi

- ○
- ○ the command to run consumer_batch.py
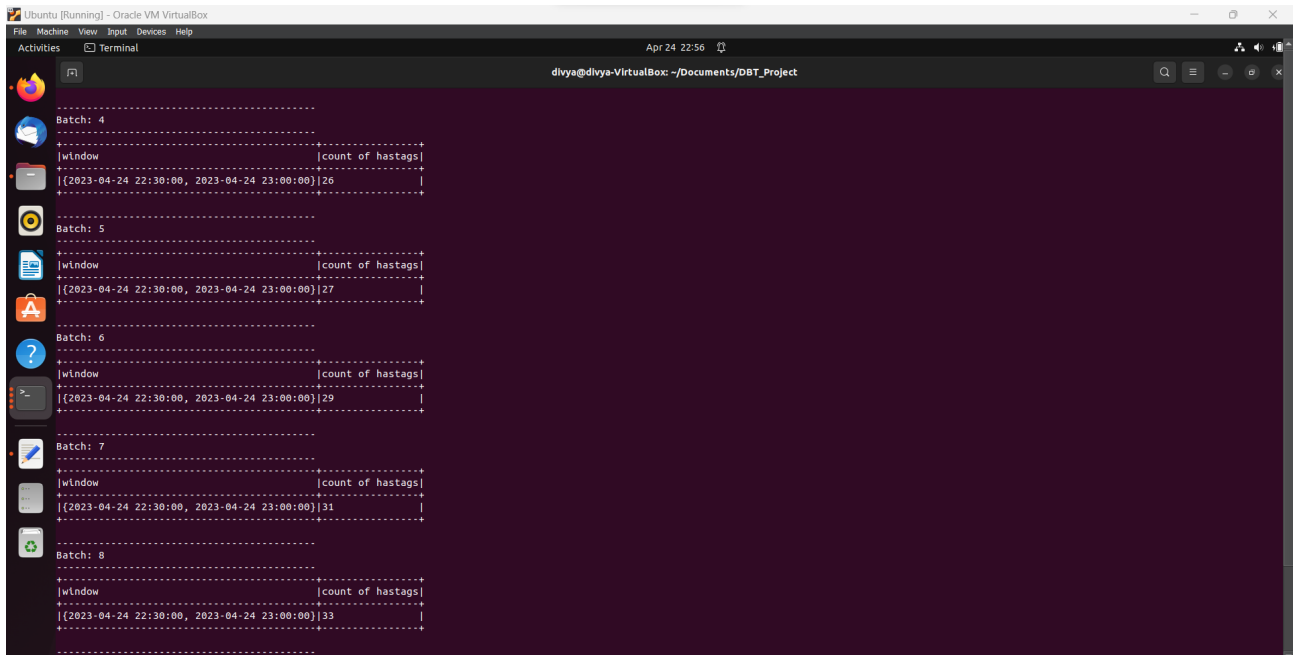


# 8. Comparison of Streaming & Batch Modes

Stream processing is much faster than batch processing as the number of tweets is counted for every 30 min where batch processing it is not real time and it reads data from the dataset

# 9. Results and Discussion

Pragathya
Ananya
Divya
Nidhi

## Output for Stream processing



## Output for batch processing

Pragathya
Ananya
Divya
Nidhi

```
23/04/24 23:14:46 INFO DAGScheduler: Job 1 finished: showString
23/04/24 23:14:46 INFO CodeGenerator: Code generated in 32.32541
+--------------------+--------------------+-----+
|              window|            hashtags|count|
+--------------------+--------------------+-----+
|{2023-04-24 22:30...|Cant wait for Feb...|    1|
|{2023-04-24 22:30...|in a few days exc...|    1|
|{2023-04-24 22:30...|later finally uau...|    1|
|{2023-04-24 22:30...|purepleasure simp...|    1|
|{2023-04-24 22:30...|what it takes But...|    1|
|{2023-04-24 22:30...|So proud of these...|    1|
|{2023-04-24 22:30...|Strike Spinnerbai...|    1|
|{2023-04-24 22:30...|NotJustForGirls P...|    1|
|{2023-04-24 22:30...|illustration art ...|    1|
|{2023-04-24 22:30...|uaufa  uaufb and ...|    1|
|{2023-04-24 22:30...|sustainable compa...|    1|
|{2023-04-24 22:30...|feel after I lift...|    1|
|{2023-04-24 22:30...|What does everybo...|    1|
|{2023-04-24 22:30...|respect the man s...|    1|
|{2023-04-24 22:30...|birthday balloons...|    1|
|{2023-04-24 22:30...|for free online o...|    1|
|{2023-04-24 22:30...|Wow we got an ama...|    1|
|{2023-04-24 22:30...|Compassion is a m...|    1|
|{2023-04-24 22:30...|Last  days ago jo...|    1|
|{2023-04-24 22:30...|know this couch c...|    1|
+--------------------+--------------------+-----+
only showing top 20 rows
```

## 10. Conclusion

Stream processing is much faster than batch processing in certain conditions but here stream processing is faster.

## 11. References

- [Streaming Spark Programming guide](#)
- [Kafka Streaming overview](#)

Pragathya
Ananya
Divya
Nidhi