

BDA Project

Sameer Gupta

211080071

Batch D

Real-Time Big Data Analytics Pipeline Experiment Report

This experiment focuses on building a real-time data analytics pipeline using various Apache components such as Kafka, Spark, Hive, Pig, HBase, and Zookeeper. The project demonstrates the flow of data from ingestion to processing, storage, and querying. Each task represents a significant stage of the pipeline and collectively provides an end-to-end understanding of real-time big data architecture.

Task 1: Setting Up the Environment

Theory

Before working with real-time big data technologies, we must prepare the environment by installing essential components. Each of these tools plays a specific role in the big data ecosystem:

- **Kafka** acts as a distributed messaging system used for data ingestion.
- **Spark** enables real-time and batch data processing
- **Hive** is a data warehouse software that facilitates querying and managing large datasets.

- **Pig** is a high-level platform for creating MapReduce programs used with Hadoop.
- **HBase** is a NoSQL database that stores non-relational, distributed data.
- **Zookeeper** coordinates distributed applications and maintains configuration information.

We also install Java, as most of these tools are Java-based, and SSH for secure communication and remote access. Configuring the correct environment paths and Java settings ensures compatibility and prevents runtime errors.

Screenshot:

Task 1: Setting Up the Environment

- Install and configure Apache Kafka, Spark, Hive, Pig, HBase, and Zookeeper.
- Verify the installation and ensure services are running.

Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help

May 2 09:50 ⓘ Fri 2 May 9:50 AM

Activities Terminal

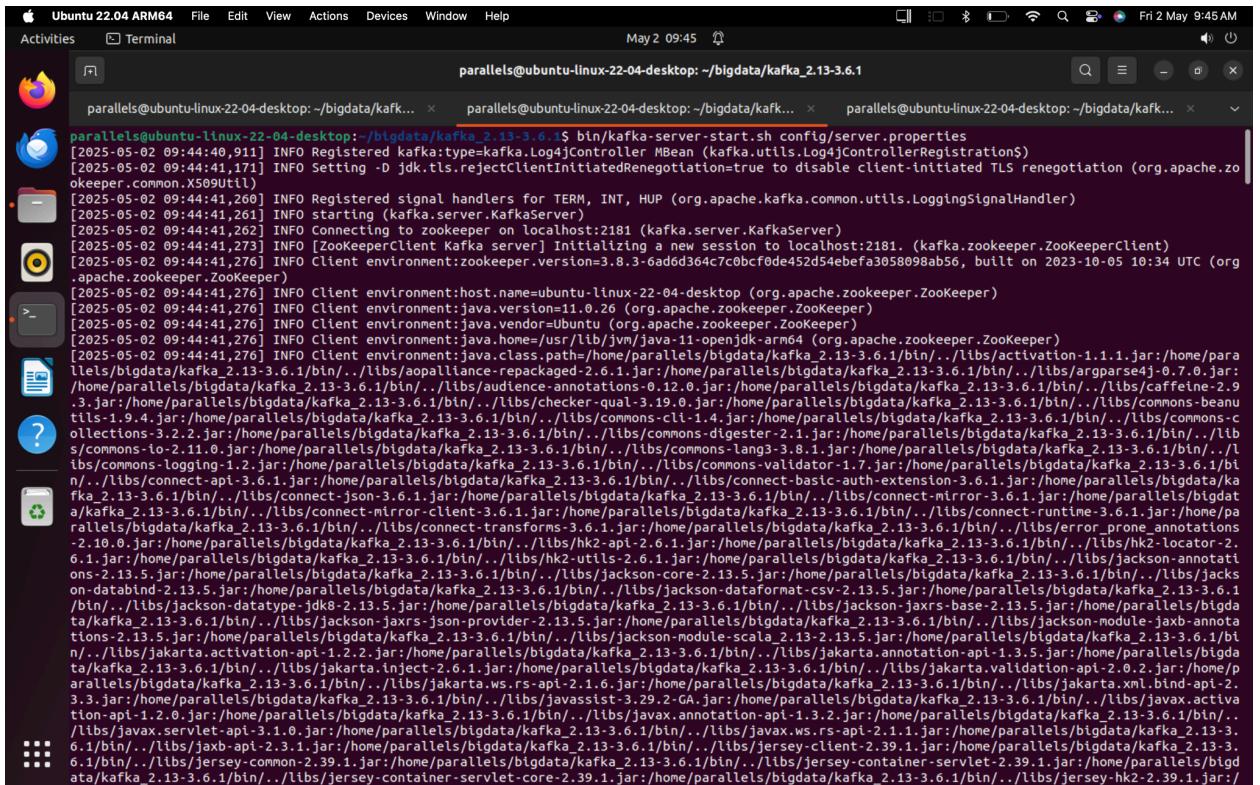
```
parallels@ubuntu-linux-22-04-desktop: ~/bigdata/spark-3.5.1-bin-hadoop3
```

parallels@ubuntu-linux-22-04-desktop: ~/bigdata/kafka_2.13-3.6.1\$ cd ..
parallels@ubuntu-linux-22-04-desktop: ~/bigdata\$ wget https://downloads.apache.org/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz
cd spark-3.5.1-bin-hadoop3
--2025-05-02 09:46:01-- https://downloads.apache.org/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz
Resolving downloads.apache.org (downloads.apache.org)... 135.181.214.104, 88.99.208.237, 2a01:4f8:10a:39da::2, ...
Connecting to downloads.apache.org (downloads.apache.org)|135.181.214.104|:443... connected.
HTTP request sent, awaiting response... 404 Not Found
2025-05-02 09:46:01 ERROR 404: Not Found.

tar (child): spark-3.5.1-bin-hadoop3.tgz: Cannot open: No such file or directory
tar (child): Error is not recoverable: exiting now
tar: Child returned status 2
tar: Error is not recoverable: exiting now
bash: cd: spark-3.5.1-bin-hadoop3: No such file or directory
parallels@ubuntu-linux-22-04-desktop: ~/bigdata\$ wget https://archive.apache.org/dist/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz
--2025-05-02 09:47:04-- https://archive.apache.org/dist/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz
Resolving archive.apache.org (archive.apache.org)... 65.108.204.189, 2a01:4f9:1:a:a084::2
Connecting to archive.apache.org (archive.apache.org)|65.108.204.189|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 400446614 (382M) [application/x-gzip]
Saving to: 'spark-3.5.1-bin-hadoop3.tgz'

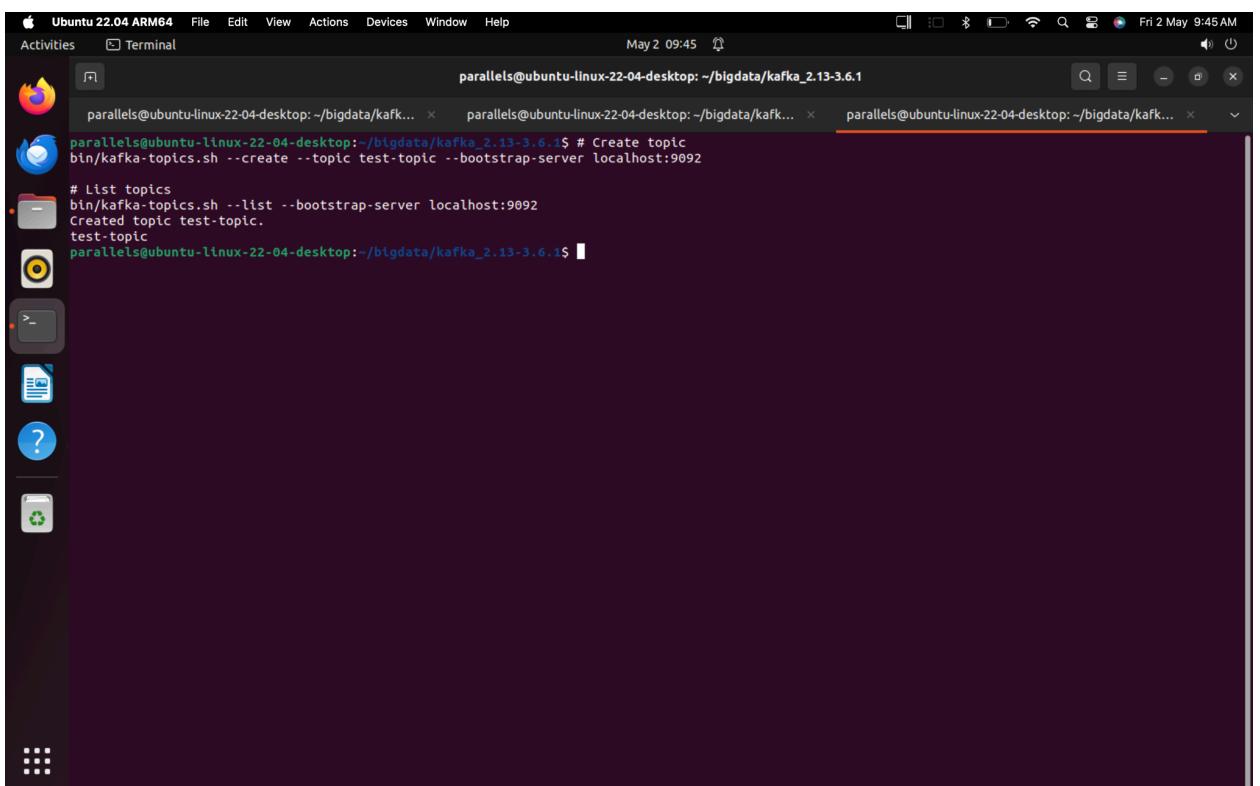
spark-3.5.1-bin-hadoop3.tgz 100%[=====] 381.90M 3.00MB/s in 2m 35s
2025-05-02 09:49:39 (2.46 MB/s) - 'spark-3.5.1-bin-hadoop3.tgz' saved [400446614/400446614]

parallels@ubuntu-linux-22-04-desktop: ~/bigdata\$ ~c
parallels@ubuntu-linux-22-04-desktop: ~/bigdata\$ tar -xzf spark-3.5.1-bin-hadoop3.tgz
cd spark-3.5.1-bin-hadoop3
parallels@ubuntu-linux-22-04-desktop: ~/bigdata/spark-3.5.1-bin-hadoop3\$./bin/spark-shell
25/05/02 09:50:06 WARN Utils: Your hostname, ubuntu-linux-22-04-desktop resolves to a loopback address: 127.0.1.1; using 10.211.55.4 instead (on interface enp0s5)
25/05/02 09:50:06 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).



parallels@ubuntu-linux-22-04-desktop: ~/bigdata/kafka_2.13-3.6.1\$ bin/kafka-server-start.sh config/server.properties

```
[2025-05-02 09:44:40,911] INFO Registered kafka:type=kafka.Log4jController MBean (kafka.utils.Log4jControllerRegistration$)
[2025-05-02 09:44:41,171] INFO Setting -D jdk.tls.rejectClientInitiatedRenegotiation=true to disable client-initiated TLS renegotiation (org.apache.zookeeper.common.X509Util)
[2025-05-02 09:44:41,260] INFO Registered signal handlers for TERM, INT, HUP (org.apache.kafka.common.utils.LoggingSignalHandler)
[2025-05-02 09:44:41,261] INFO Starting (kafka.server.KafkaServer)
[2025-05-02 09:44:41,262] INFO Connecting to zookeeper on localhost:2181 (kafka.server.KafkaServer)
[2025-05-02 09:44:41,273] INFO [ZooKeeperClient Kafka server] Initializing a new session to localhost:2181. (kafka.zookeeper.ZooKeeperClient)
[2025-05-02 09:44:41,276] INFO Client environment:zookeeper.version=3.8.3-6add6d364c7rc0bcf0de452d54ebefaa05809ab56, built on 2023-10-05 10:34 UTC (org.apache.zookeeper.ZooKeeper)
[2025-05-02 09:44:41,276] INFO Client environment:host.name=ubuntu-linux-22-04-desktop (org.apache.zookeeper.ZooKeeper)
[2025-05-02 09:44:41,276] INFO Client environment:java.version=11.0.20 (org.apache.zookeeper.ZooKeeper)
[2025-05-02 09:44:41,276] INFO Client environment:java.vendor=Ubuntu (org.apache.zookeeper.ZooKeeper)
[2025-05-02 09:44:41,276] INFO Client environment:java.home=/usr/lib/jvm/java-11-openjdk-arm64 (org.apache.zookeeper.ZooKeeper)
[2025-05-02 09:44:41,276] INFO Client environment:java.class.path=/home/parallels/bigdata/kafka_2.13-3.6.1/bin/../libs/activation-1.1.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/apollo-repackaged-2.6.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/audience-annotations-0.12.0.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/checker-qual-3.19.0.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/commons-beanutils-1.9.4.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/commons-cli-1.4.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/commons-collections-3.2.2.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/commons-digester-2.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/commons-lang3-3.8.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/commons-logging-1.2.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/commons-validator-1.7.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/connect-api-3.6.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/connect-basic-auth-extension-3.6.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/connect-client-3.6.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/connect-mirror-client-3.6.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/connect-runtime-3.6.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/connect-transforms-3.6.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/error_prone_annotations-2.10.0.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/hk2-apl-2.6.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/hk2-util-2.6.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jackson-annotations-2.13.5.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jackson-core-2.13.5.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jackson-dataformat-csv-2.13.5.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jackson-dataformat-scala-2.13-2.13.5.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jackson-datatype-jdk8-2.13.5.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jackson-jaxrs-base-2.13.5.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jackson-jaxrs-json-provider-2.13.5.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jackson-module-scala-2.13-2.13.5.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jakarta.activation-api-1.2.2.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jakarta.annotation-api-1.3.5.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jakarta.validation-api-2.0.2.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jakarta.xml.bind-api-2.3.3.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/javax.activation-api-1.2.0.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/javax.annotation-api-1.3.2.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/javax.ws.rs-api-2.1.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jersey-client-2.39.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jersey-contalner-servlet-2.39.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jersey-container-servlet-core-2.39.1.jar:/home/parallels/bigdata/kafka_2.13-3.6.1/bin/..,/libs/jersey-hk2-2.39.1.jar:/..
```



```
parallels@ubuntu-linux-22-04-desktop: ~/bigdata/kafka_2.13-3.6.1$ bin/kafka-topics.sh --create --topic test-topic --bootstrap-server localhost:9092
```

List topics

```
bin/kafka-topics.sh --list --bootstrap-server localhost:9092
Created topic test-topic.
test-topic
```

```
Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help May 2 09:44
Activities Terminal parallels@ubuntu-linux-22-04-desktop: ~/bigdata/kafka_2.13-3.6.1
parallels@ubuntu-linux-22-04-desktop: ~/bigdata/kafka_2.13-3.6.1
[2025-05-02 09:43:55,798] INFO clientPortAddress is 0.0.0.0:2181 (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,798] INFO secureClientPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,799] INFO observerMasterPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,799] INFO metricsProvider.className is org.apache.zookeeper.metrics.impl.DefaultMetricsProvider (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,800] INFO autopurge.snapRetainCount set to 3 (org.apache.zookeeper.server.DatadirCleanupManager)
[2025-05-02 09:43:55,800] INFO autopurge.purgeInterval set to 0 (org.apache.zookeeper.server.DatadirCleanupManager)
[2025-05-02 09:43:55,800] INFO Purge task is not scheduled. (org.apache.zookeeper.server.DatadirCleanupManager)
[2025-05-02 09:43:55,800] WARN Either no config or no quorum defined in config, running in standalone mode (org.apache.zookeeper.server.quorum.QuorumPeerMain)
[2025-05-02 09:43:55,801] INFO Log4j 1.2 jmx support not found; jmx disabled. (org.apache.zookeeper.jmx.ManagedUtil)
[2025-05-02 09:43:55,801] INFO Reading configuration from: config/zookeeper.properties (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,801] WARN config/zookeeper.properties is relative. Prepend ./ to indicate that you're sure! (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,801] INFO clientPortAddress is 0.0.0.0:2181 (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,802] INFO secureClientPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,802] INFO observerMasterPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,802] INFO metricsProvider.className is org.apache.zookeeper.metrics.impl.DefaultMetricsProvider (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2025-05-02 09:43:55,802] INFO Starting server (org.apache.zookeeper.server.ZooKeeperServerMain)
[2025-05-02 09:43:55,824] INFO ServerMetrics initialized with provider org.apache.zookeeper.metrics.impl.DefaultMetricsProvider@4944252c (org.apache.zookeeper.server.ServerMetrics)
[2025-05-02 09:43:55,830] INFO ACL digest algorithm is: SHA1 (org.apache.zookeeper.server.auth.DigestAuthenticationProvider)
[2025-05-02 09:43:55,830] INFO zookeeper.DigestAuthenticationProvider.enabled = true (org.apache.zookeeper.server.auth.DigestAuthenticationProvider)
[2025-05-02 09:43:55,833] INFO zookeeper.snapshot.trust.empty : false (org.apache.zookeeper.server.persistence.FileTxnSnapLog)
[2025-05-02 09:43:55,844] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,844] INFO [-----] (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,844] INFO [ / ] (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,845] INFO [ / \ ] (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,845] INFO [ / \ ] (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,845] INFO [ / \ ] < [ / \ ] (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,845] INFO [ / \ ] (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,846] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,846] INFO Server environment:zookeeper.version=3.8.3-6ad6d364c7c0bcf0de452d54ebef305809ab56, built on 2023-10-05 10:34 UTC (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,846] INFO Server environment:host.name=ubuntu-linux-22-04-desktop (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,846] INFO Server environment:java.version=11.0.20 (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,846] INFO Server environment:java.vendor=Ubuntu (org.apache.zookeeper.server.ZooKeeperServer)
[2025-05-02 09:43:55,846] INFO Server environment:java.home=/usr/lib/jvm/java-11-openjdk-arm64 (org.apache.zookeeper.server.ZooKeeperServer)
```

```
Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help May 2 09:43
Activities Terminal parallels@ubuntu-linux-22-04-desktop: ~/bigdata/kafka_2.13-3.6.1
parallels@ubuntu-linux-22-04-desktop: $ mkdir ~/bigdata && cd ~/bigdata
parallels@ubuntu-linux-22-04-desktop: ~/bigdata$ wget https://downloads.apache.org/kafka/3.6.0/kafka_2.13-3.6.0.tgz
--2025-05-02 09:39:55-- https://downloads.apache.org/kafka/3.6.0/kafka_2.13-3.6.0.tgz
Resolving downloads.apache.org (downloads.apache.org)... 135.181.214.104, 88.99.208.237, 2a01:4f8:10a:39da::2, ...
Connecting to downloads.apache.org (downloads.apache.org)|135.181.214.104|:443... connected.
HTTP request sent, awaiting response... 404 Not Found
2025-05-02 09:39:55 ERROR 404: Not Found.

parallels@ubuntu-linux-22-04-desktop: ~/bigdata$ wget https://downloads.apache.org/kafka/3.6.1/kafka_2.13-3.6.1.tgz
--2025-05-02 09:40:52-- https://downloads.apache.org/kafka/3.6.1/kafka_2.13-3.6.1.tgz
Resolving downloads.apache.org (downloads.apache.org)... 135.181.214.104, 88.99.208.237, 2a01:4f8:10a:39da::2, ...
Connecting to downloads.apache.org (downloads.apache.org)|135.181.214.104|:443... connected.
HTTP request sent, awaiting response... 404 Not Found
2025-05-02 09:40:53 ERROR 404: Not Found.

parallels@ubuntu-linux-22-04-desktop: ~/bigdata$ ^C
parallels@ubuntu-linux-22-04-desktop: ~/bigdata$ wget https://archive.apache.org/dist/kafka/3.6.1/kafka_2.13-3.6.1.tgz
--2025-05-02 09:41:33-- https://archive.apache.org/dist/kafka/3.6.1/kafka_2.13-3.6.1.tgz
Resolving archive.apache.org (archive.apache.org)... 65.108.204.189, 2a01:4f9:1:a:a084::2
Connecting to archive.apache.org (archive.apache.org)|65.108.204.189|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 113466638 (108M) [application/x-gzip]
Saving to: 'kafka_2.13-3.6.1.tgz'

kafka_2.13-3.6.1.tgz          100%[=====] 108.21M  4.45MB/s   in 26s

2025-05-02 09:41:59 (4.18 MB/s) - 'kafka_2.13-3.6.1.tgz' saved [113466638/113466638]

parallels@ubuntu-linux-22-04-desktop: ~/bigdata$ tar -xzf kafka_2.13-4.0.0.tgz
tar (child): kafka_2.13-4.0.0.tgz: Cannot open: No such file or directory
tar (child): Error is not recoverable: exiting now
tar: Child returned status 2
tar: Error is not recoverable: exiting now
parallels@ubuntu-linux-22-04-desktop: ~/bigdata$ tar -xzf kafka_2.13-3.6.1.tgz
parallels@ubuntu-linux-22-04-desktop: ~/bigdata$ cd kafka_2.13-3.6.1/
parallels@ubuntu-linux-22-04-desktop: ~/bigdata/kafka_2.13-3.6.1$ clear
```

```

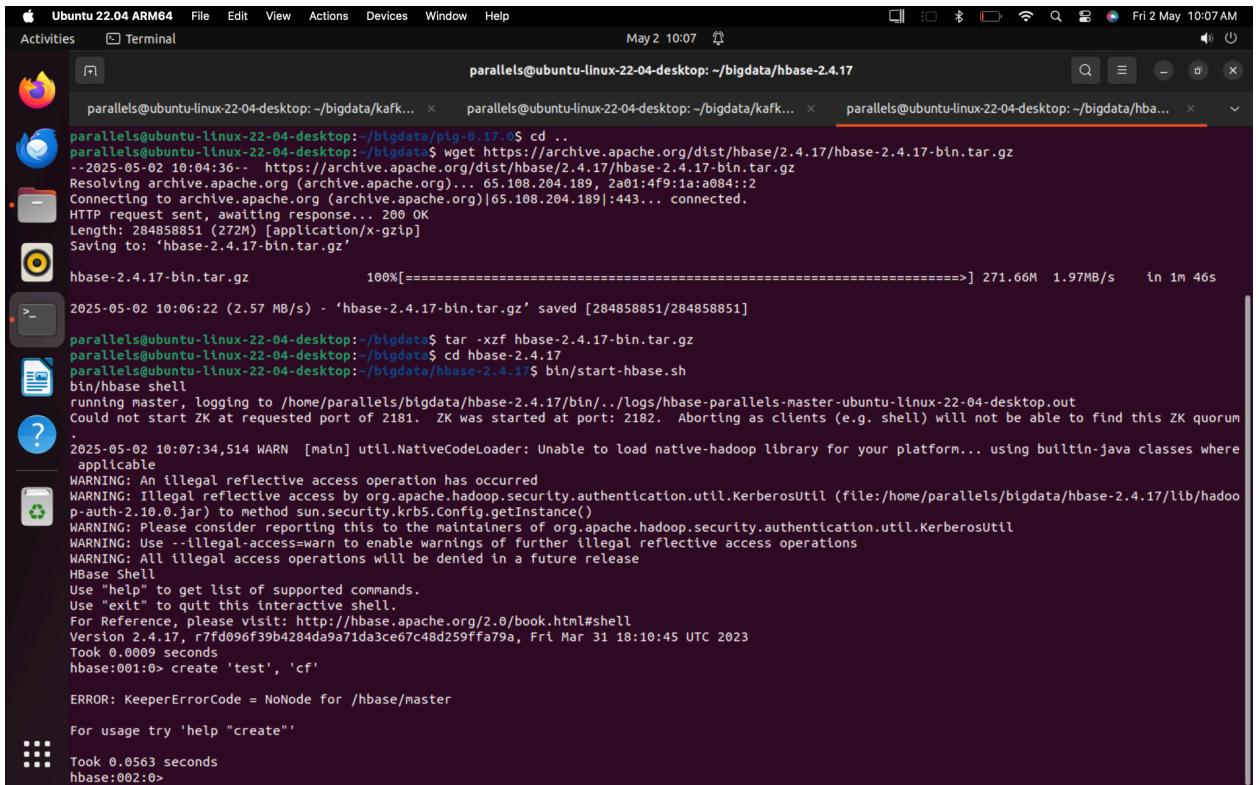
parallels@ubuntu-linux-22-04-desktop: $ sudo apt update
[sudo] password for parallels:
Hit:1 https://packages.mozilla.org/apt mozilla InRelease
Hit:2 http://ports.ubuntu.com/ubuntu-ports jammy InRelease
Hit:3 http://ports.ubuntu.com/ubuntu-ports jammy-updates InRelease
Hit:4 http://ports.ubuntu.com/ubuntu-ports jammy-backports InRelease
Get:5 http://ports.ubuntu.com/ubuntu-ports jammy-security InRelease [129 kB]
Fetched 129 kB in 2s (54.0 kB/s)
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
262 packages can be upgraded. Run 'apt list --upgradable' to see them.
W: Target Packages (main/binary-arm64/Packages) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target Packages (main/binary-all/Packages) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target Translations (main/i18n/Translation-en) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target DEP-11 (main/dep11/Components-arm64.yml) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target DEP-11 (main/dep11/Components-all.yml) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target DEP-11-icons-small (main/dep11/icons-48x48.tar) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target DEP-11-icons (main/dep11/icons-64x64.tar) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target DEP-11-icons-hidpi (main/dep11/icons-64x64@2.tar) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target CNF (main/cnf/Commands-arm64) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target CNF (main/cnf/Commands-all) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target Packages (main/binary-arm64/Packages) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target Packages (main/binary-all/Packages) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target Translations (main/i18n/Translation-en) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target DEP-11 (main/dep11/Components-arm64.yml) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target DEP-11 (main/dep11/Components-all.yml) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2
W: Target DEP-11-icons-small (main/dep11/icons-48x48.tar) is configured multiple times in /etc/apt/sources.list.d/mozilla.list:1 and /etc/apt/sources.list.d/mozilla.list:2

```

```

parallels@ubuntu-linux-22-04-desktop: ~$ tar -xzf pig-0.17.0.tar.gz
parallels@ubuntu-linux-22-04-desktop: ~$ cd pig-0.17.0
parallels@ubuntu-linux-22-04-desktop: ~$ ./bin/pig -x local
Error: JAVA_HOME is not set.
parallels@ubuntu-linux-22-04-desktop: ~$ readlink -f $(which java)
/usr/lib/jvm/java-11-openjdk-arm64/bin/java
parallels@ubuntu-linux-22-04-desktop: ~$ export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-arm64
export PATH=$JAVA_HOME/bin
parallels@ubuntu-linux-22-04-desktop: ~$ source ~/.bashrc # or `source ~/.zshrc`
parallels@ubuntu-linux-22-04-desktop: ~$ echo $JAVA_HOME
/usr/lib/jvm/java-11-openjdk-arm64
parallels@ubuntu-linux-22-04-desktop: ~$ bin/pig -x local
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication.util.KerberosUtil (file:/home/parallels/bigdata/pig-0.17.0/lib/hadoop2-runtime/hadoop-auth-2.7.3.jar) to method sun.security.krb5.Config.getInstance()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.security.authentication.util.KerberosUtil
WARNING: Use -illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
2025-05-02 10:03:14,260 INFO [main] pig.ExecTypeProvider (ExecTypeProvider.java:selectExecType(41)) - Trying ExecType : LOCAL
2025-05-02 10:03:14,261 INFO [main] pig.ExecTypeProvider (ExecTypeProvider.java:selectExecType(43)) - Picked LOCAL as the ExecType
2025-05-02 10:03:14,285 [main] INFO org.apache.pig.Main - Apache Pig Version 0.17.0 (r1797386) compiled Jun 02 2017, 15:41:58
2025-05-02 10:03:14,285 [main] INFO org.apache.pig.Main - Logging error messages to: /home/parallels/bigdata/pig-0.17.0/pig_1746160394281.log
2025-05-02 10:03:14,297 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/parallels/.pigbootup not found
2025-05-02 10:03:14,405 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.job.tracker
2025-05-02 10:03:14,406 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: file:///r-checksum
2025-05-02 10:03:14,499 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-05-02 10:03:14,512 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-default-33692a4c-2e91-423f-9c48-97cc9cd8e620
2025-05-02 10:03:14,513 [main] WARN org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
grunt> 

```



```
parallels@ubuntu-linux-22-04-desktop:~/bigdata/kaf... parallels@ubuntu-linux-22-04-desktop:~/bigdata/kaf... parallels@ubuntu-linux-22-04-desktop:~/bigdata/hba...
parallels@ubuntu-linux-22-04-desktop:~/bigdata/pig-0.17.0$ cd ..
--2025-05-02 10:04:36-- https://archive.apache.org/dist/hbase/2.4.17/hbase-2.4.17-bin.tar.gz
Resolving archive.apache.org (archive.apache.org)... 65.108.204.189, 2a01:4f9:1:a:a084::2
Connecting to archive.apache.org (archive.apache.org)|65.108.204.189|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 284858851 (272M) [application/x-gzip]
Saving to: 'hbase-2.4.17-bin.tar.gz'

hbase-2.4.17-bin.tar.gz      100%[=====] 271.66M  1.97MB/s   in 1m 46s
2025-05-02 10:06:22 (2.57 MB/s) - 'hbase-2.4.17-bin.tar.gz' saved [284858851/284858851]

parallels@ubuntu-linux-22-04-desktop:~/bigdata$ tar -xzf hbase-2.4.17-bin.tar.gz
parallels@ubuntu-linux-22-04-desktop:~/bigdata$ cd hbase-2.4.17
parallels@ubuntu-linux-22-04-desktop:~/bigdata/hbase-2.4.17$ bin/start-hbase.sh
bin/hbase shell
running master, logging to /home/parallels/bigdata/hbase-2.4.17/bin/../logs/hbase-parallels-master-ubuntu-linux-22-04-desktop.out
Could not start zk at requested port of 2181. ZK was started at port: 2182. Aborting as clients (e.g. shell) will not be able to find this ZK quorum
.
2025-05-02 10:07:34,514 WARN  [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication.util.KerberosUtil (file:/home/parallels/bigdata/hbase-2.4.17/lib/hadoop-auth-2.10.0.jar) to method sun.security.krb5.Config.getInstance()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
HBase Shell
Use "help" to get list of supported commands.
Use "exit" to quit this interactive shell.
For Reference, please visit: http://hbase.apache.org/2.0/book.html#shell
Version 2.4.17, r7fd096f39b4284da9a71da3ce67c48d259ffa79a, Fri Mar 31 18:10:45 UTC 2023
Took 0.0009 seconds
hbase:001:0> create 'test', 'cf'

ERROR: KeeperErrorCode = NoNode for /hbase/master
For usage try 'help "create"'


Took 0.0563 seconds
hbase:002:0>
```

Task 2: Implementing Real-Time Data Ingestion

- Set up Kafka producers to generate real-time data streams.
- Configure Kafka consumers to feed data into Apache Spark Streaming.

A screenshot of an Ubuntu 22.04 ARM64 desktop environment. The desktop has a dark theme with a dock on the left containing icons for the Dash, Home, Applications, and Help. A terminal window is open in the center, showing the command line interface. The terminal title is "parallels@ubuntu-linux-22-04-desktop: ~/blgdata/kafka_2.13-3.6.1". The user has run several commands related to Kafka:

```
parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ cd ..;/kafka_2.13-3.6.1/
parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ bin/kafka-topics.sh --create --topic real-time-data --bootstrap-server localhost:9092
--partitions 1 --replication-factor 1
Created topic real-time-data.
parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ bin/kafka-topics.sh --list --bootstrap-server localhost:9092
real-time-data
test-topic
parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ bin/kafka-console-producer.sh --topic real-time-data --bootstrap-server localhost:9092
2
>message1
>message2
>message3
>messag 4
>
```

A screenshot of the same Ubuntu 22.04 ARM64 desktop environment. The desktop layout is identical to the first one, with the dark theme, dock on the left, and terminal window in the center. The terminal title is "parallels@ubuntu-linux-22-04-desktop: ~/blgdata/kafka_2.13-3.6.1". The user has run a command to consume messages from the "real-time-data" topic:

```
parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ bin/kafka-console-consumer.sh --topic real-time-data --bootstrap-server localhost:9092
2 --from-beginning
message1
message2
message3
messag 4
```

Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help Fri 2 May 10:31AM

Activities Terminal

```
parallels@ubuntu-linux-22-0... ~ blgdata/kafka_2.13-3.6.1
parallels@ubuntu-linux-22-04-desktop:~/blgdata/hbase-2.4.17$ cd ..../kafka_2.13-3.6.1/
parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ bin/kafka-topics.sh --create --topic real-time-data --bootstrap-server localhost:9092
--partitions 1 --replication-factor 1
Created topic real-time-data.
parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ bin/kafka-topics.sh --list --bootstrap-server localhost:9092
real-time-data
test-topic
parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ bin/kafka-console-producer.sh --topic real-time-data --bootstrap-server localhost:9092
>message1
>message2
>message3
>message4
>message5
>>parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ bin/kafka-console-producer.sh --topic real-time-data --bootstrap-server localhost:9092
9092
>message6
>message7
>
```

Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help Fri 2 May 10:31AM

Activities Terminal

```
parallels@ubuntu-linux-22-0... ~ blgdata/kafka_2.13-3.6.1
parallels@ubuntu-linux-22-04-desktop:~/blgdata/kafka_2.13-3.6.1$ pId=spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor1 Seeking to offset 6 for partition real-time-data-0
25/05/02 10:30:59 INFO Metadata: [Consumer clientId=consumer-spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor-1, groupId=spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor] Cluster ID: rGoz8WeSTVqYKcDdqNM0y0g
25/05/02 10:30:59 INFO SubscriptionState: [Consumer clientId=consumer-spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor-1, groupId=spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor] Seeking to EARLIEST offset of partition real-time-data-0
25/05/02 10:30:59 INFO SubscriptionState: [Consumer clientId=consumer-spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor-1, groupId=spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor] Resetting offset for partition real-time-data-0 to offset 0.
25/05/02 10:30:59 INFO SubscriptionState: [Consumer clientId=consumer-spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor-1, groupId=spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor] Seeking to LATEST offset of partition real-time-data-0
25/05/02 10:30:59 INFO SubscriptionState: [Consumer clientId=consumer-spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor-1, groupId=spark-kafka-source-2b17af55-fde0-4fb6-bd9d-7964cf450b23-1227988132-executor] Resetting offset for partition real-time-data-0 to offset 7.
25/05/02 10:31:00 INFO DataWritingSparkTask: Writer for partition 0 is committing.
25/05/02 10:31:00 INFO DataWritingSparkTask: Committed partition 0 (task 1, attempt 0, stage 1.0)
25/05/02 10:31:00 INFO Executor: Finished task 0.0 in stage 1.0 (TID 1). 1875 bytes result sent to driver
25/05/02 10:31:00 INFO TaskSetManager: Finished task 0.0 in stage 1.0 (TID 1) in 1448 ms on ubuntu-22.04-arm64.shared (executor driver) (1/1)
25/05/02 10:31:00 INFO TaskSchedulerImpl: Removed TaskSet 1.0, whose tasks have all completed, from pool
25/05/02 10:31:00 INFO DAGScheduler: ResultStage 1 (start at NativeMethodAccessorImpl.java:0) finished in 1.494 s
25/05/02 10:31:00 INFO DAGScheduler: Job 1 is finished. Cancelling potential speculative or zombie tasks for this job
25/05/02 10:31:00 INFO TaskschedulerImpl: Killing all running tasks in stage 1: Stage finished
25/05/02 10:31:00 INFO DAGScheduler: Job 1 finished: start at NativeMethodAccessorImpl.java:0, took 1.552464 s
25/05/02 10:31:00 INFO WriteToDataSourceV2Exec: Data source write support MicroBatchWrite[epoch: 1, writer: ConsoleWriter[numRows=20, truncate=true]] is committing.
-----
Batch: 1
25/05/02 10:31:00 INFO CodeGenerator: Code generated in 22.632384 ms
25/05/02 10:31:00 INFO BlockManagerInfo: Removed broadcast_1_piece0 on ubuntu-22.04-arm64.shared:45779 in memory (size: 4.8 KiB, free: 434.4 MiB)
+-----+
| words|
+-----+
|[message, 7]|
+-----+
25/05/02 10:31:01 INFO WriteToDataSourceV2Exec: Data source write support MicroBatchWrite[epoch: 1, writer: ConsoleWriter[numRows=20, truncate=true]] committed.
25/05/02 10:31:01 INFO CheckpointFileManager: Writing atomically to file:/tmp/temporary-3c97caa1-a170-4596-af3f-d2d1b89d08e5/commits/1 using temp file file:/tmp/temporary-3c97caa1-a170-4596-af3f-d2d1b89d08e5/commits/.1.6fd18ce1-f4fd-45ec-9b0d-871bb0276763.tnp
25/05/02 10:31:01 INFO CheckpointFileManager: Renamed temp file file:/tmp/temporary-3c97caa1-a170-4596-af3f-d2d1b89d08e5/commits/.1.6fd18ce1-f4fd-45ec-9b0d-871bb0276763.tnp to file:/tmp/temporary-3c97caa1-a170-4596-af3f-d2d1b89d08e5/commits/.1.6fd18ce1-f4fd-45ec
25/05/02 10:31:01 INFO MicroBatchExecution: Streaming query made progress: {
  "id" : "a4ff3b1b-a465-4236-a775-dc86f79746d2",
}
```

Task 3: Processing Data with Apache Spark

- Develop a Spark Streaming application to process incoming data.
- Perform real-time transformations and aggregations.
- Store processed data in Apache HBase.

```
Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help Fri 2 May 10:34 AM
Activities Text Editor *spark.kafka.consumer.py -bigdata/kafka_2.13-3.6.1
Save X
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import expr
3
4 # Create a Spark session
5 spark = SparkSession.builder \
6     .appName("KafkaSparkStreaming") \
7     .getOrCreate()
8
9 # Read streaming data from Kafka
10 kafka_stream = spark \
11     .readStream \
12     .format("kafka") \
13     .option("kafka.bootstrap.servers", "localhost:9092") \
14     .option("subscribe", "real-time-data") \
15     .load()
16
17 # Convert the Kafka 'value' field to a string
18 kafka_stream = kafka_stream.selectExpr("CAST(value AS STRING)")
19
20 # Perform some transformations: split the string into words
21 transformed_stream = kafka_stream.selectExpr("split(value, ' ') as words")
22
23 # Perform an aggregation: count the number of words per batch
24 word_counts = transformed_stream \
25     .groupBy("words") \
26     .count()
27
28 # Output the results to the console (for debugging purposes)
29 query = word_counts \
30     .writeStream \
31     .outputMode("complete") \
32     .format("console") \
33     .start()
34
35 query.awaitTermination()
36
```



```
Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help Fri 2 May 10:35 AM
Activities Terminal parallels@ubuntu-linux-22-0... - bigdata/kafka_2.13-3.6.1
parallels@ubuntu-linux-22-0... x parallels@ubuntu-linux-22-0... x parallels@ubuntu-linux-22-0... x parallels@ubuntu-linux-22-0...
May 2 10:35 Fri 2 May 10:35 AM
parallels@ubuntu-linux-22-0... x parallels@ubuntu-linux-22-0... x parallels@ubuntu-linux-22-0... x parallels@ubuntu-linux-22-0...
25/05/02 10:35:08 INFO HDFSBackend$FileStateProvider: Committed version 2 for HDFSStateStore[id=(ep=0,part=199),dir=file:/tmp/temporary-0a36596a-65f5-4a8-9980-b254787b625/state/0/199] to file file:/tmp/temporary-0a36596a-65f5-4a8-9980-b254787b625/state/0/199.0
25/05/02 10:35:08 INFO DataWritingSparkTask: Writer for partition 199 is committing.
25/05/02 10:35:08 INFO Executor: Finished task 199.0 in stage 3. (TID 400) 7277 bytes result sent to driver
25/05/02 10:35:08 INFO TaskSchedulerImpl: Removed TaskSet 3.0, whose tasks have all completed, from pool
25/05/02 10:35:08 INFO DAGScheduler: ResultStage 3 (start at NativeMethodAccessorImpl.java:0) finished in 7.697 s
25/05/02 10:35:08 INFO DAGScheduler: Job 1 is finished. Cancelling potential speculative or zombie tasks for this job
25/05/02 10:35:08 INFO TaskSchedulerImpl: Killing all running tasks in stage 3: Stage finished
25/05/02 10:35:08 INFO DAGScheduler: Job 1 is finished. start at NativeMethodAccessorImpl.java:0, took 8.881747 s
25/05/02 10:35:08 INFO WriteToDataSourceV2Exec: Data source write support MicroBatchWrite[epoch= 1, writer: ConsoleWriter[numRows=20, truncate=true]] is committing.
Batch: 1
25/05/02 10:35:08 INFO CodeGenerator: Code generated in 21.258688 ms
25/05/02 10:35:10 INFO CodeGenerator: Code generated in 8.604197 ms
+---+---+
| words|count|
+---+---+
|[message, 8]| 1
+---+---+
25/05/02 10:35:10 INFO WriteToDataSourceV2Exec: Data source write support MicroBatchWrite[epoch: 1, writer: ConsoleWriter[numRows=20, truncate=true]]
committed.
25/05/02 10:35:10 INFO CheckpointFileManager: Writing atomically to file:/tmp/temporary-0a36596a-65f5-4a8-9980-b254787b625/commits/1 using temp file file:/tmp/temporary-0a36596a-65f5-4a8-9980-b254787b625/commits/1.63c2eacb-e83-4931-a9b-0f1396e2a985.tmp
25/05/02 10:35:10 INFO CheckpointFileManager: Renamed temp file file:/tmp/temporary-0a36596a-65f5-4a8-9980-b254787b625/commits/1.63c2eacb-e83-4931-a9b-0f1396e2a985.tmp to /tmp/temporary-0a36596a-65f5-4a8-9980-b254787b625/commits/1
25/05/02 10:35:10 INFO MicroBatchExecution: Streaming query made progress:
"id": "dfa6daec-c741-4135-b1e3-7751e4c5ea7",
"runId": "6b7727e23-a05f-462e-8381-9a543869ae30",
"name": null,
"timestamp": "2025-05-02T05:04:59.617Z",
"batchId": 1,
"numInputRows": 1,
"inputRowsPerSecond": 83.33333333333333,
"processedRowsPerSecond": 0.09157509157509157,
"duration": 1073,
"addBatch": 1073,
"commitOffsets": 30,
```

Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help

May 2 10:46

Activities Text Editor

Open spark_kafka_consumer.py

```

1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import expr
3
4 # Create a Spark session
5 spark = SparkSession.builder \
6     .appName("KafkaSparkStreaming") \
7     .getOrCreate()
8
9 # Read streaming data from Kafka
10 kafka_stream = spark \
11     .readStream \
12     .format("kafka") \
13     .option("kafka.bootstrap.servers", "localhost:9092") \
14     .option("subscribe", "real-time-data") \
15     .load()
16
17 # Convert the Kafka 'value' field to a string
18 kafka_stream = kafka_stream.selectExpr("CAST(value AS STRING)")
19
20 # Perform some transformations: split the string into words
21 transformed_stream = kafka_stream.selectExpr("split(value, ' ') as words")
22
23 # Perform an aggregation: count the number of words per batch
24 word_counts = transformed_stream \
25     .groupByKey("words") \
26     .count()
27
28 # Write to HBase (you need to use the HBase DataFrame API to save data)
29 word_counts.write \
30     .format("org.apache.hadoop.hbase.spark") \
31     .option("hbase.table", "word_count") \
32     .option("hbase.columns.mapping", "row_key STRING :count INT") \
33     .save()
34
35 # Output the results to the console (optional, for debugging)
36 query = word_counts \
37     .writeStream \
38     .outputMode("complete") \
39     .format("console") \
40     .start()
41
42 query.awaitTermination()
43

```

Saving file "/home/parallels/bigdata/kafka_2.13-3.6.1/spark_kafka_consumer.py"...

Python 2 Tab Width: 8 Ln 43, Col 1 INS

Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help

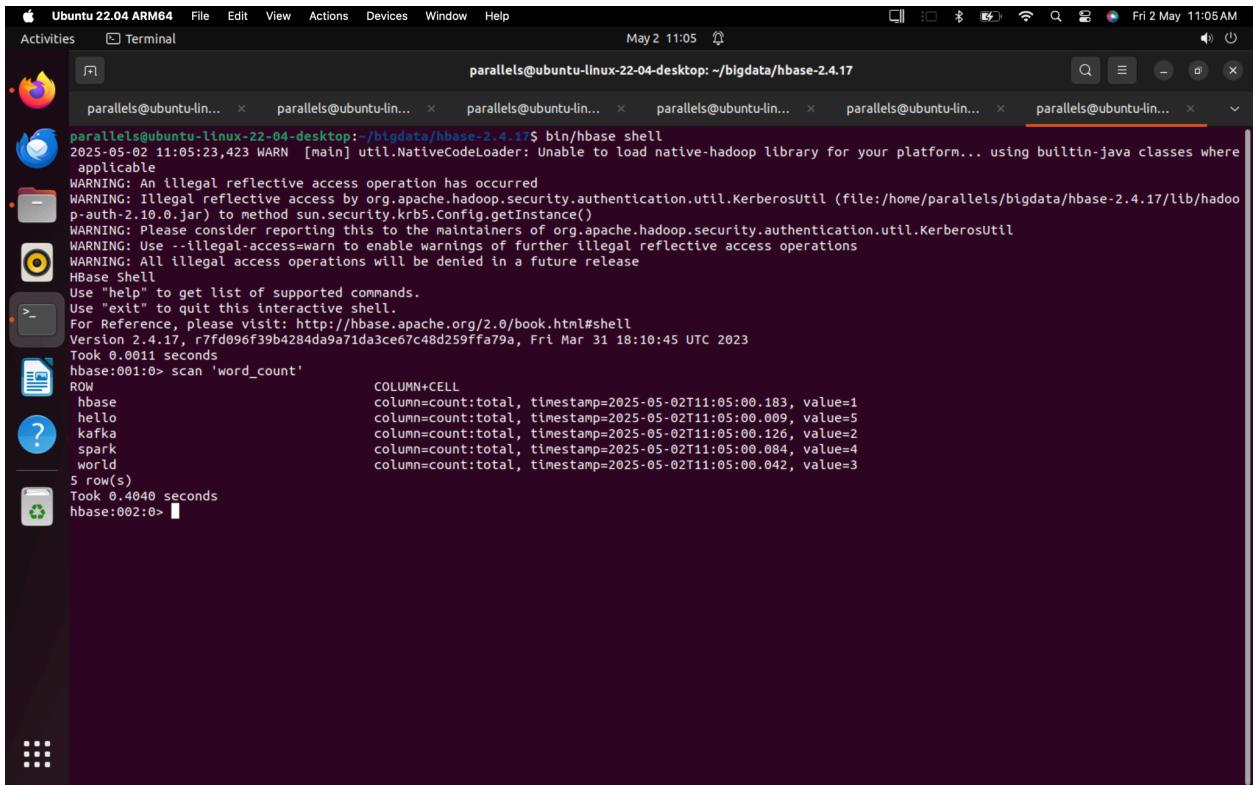
May 2 11:02

Activities Terminal

```

parallels@ubuntu-linux-22-04-desktop:~/bigdata/kafka_2.13-3.6.1$ $SPARK_HOME/bin/spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.1.2 spark_kafka_consumer.py
25/05/02 10:52:56 WARN Utils: Your hostname, ubuntu-linux-22-04-desktop resolves to a loopback address: 127.0.1.1; using 10.211.55.4 instead (on interface enp0s5)
25/05/02 10:52:56 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
:: loading settings :: url = jar:file:/home/parallels/bigdata/spark-3.1.1-bin-hadoop3/jars/ivy-2.5.1.jar!/org/apache/ivy/core/settings/ivysettings.xml
Ivy Default Cache set to: /home/parallels/.ivy2/cache
The jars for the packages stored in: /home/parallels/.ivy2/jars
org.apache.spark#spark-sql-kafka-0-10_2.12 added as a dependency
:: resolving dependencies :: org.apache.spark#spark-submit-parent-6d97776f-d782-406a-80c9-7da356f7bec8;1.0
    confs: [default]
        found org.apache.spark#spark-sql-kafka-0-10_2.12;3.1.2 in central
        found org.apache.spark#spark-token-provider-kafka-0-10_2.12;3.1.2 in central
        found org.apache.kafka#kafka-clients;2.6.0 in central
        found com.github.luben#zstd-jni;1.4.8-1 in central
        found org.lz4#lz4-java;1.7.1 in central
        found org.xerial.snappy#snappy-java;1.1.8.2 in central
        found org.slf4j#slf4j-api;1.7.30 in central
        found org.spark-project.spark#unused;1.0.0 in central
        found org.apache.commons#commons-pool2;2.6.2 in central
    :: resolution report :: resolve 356ms :: artifacts dl 23ms
        :: modules in use:
            com.github.luben#zstd-jni;1.4.8-1 from central in [default]
            org.apache.commons#commons-pool2;2.6.2 from central in [default]
            org.apache.kafka#kafka-clients;2.6.0 from central in [default]
            org.apache.spark#spark-sql-kafka-0-10_2.12;3.1.2 from central in [default]
            org.apache.spark#spark-token-provider-kafka-0-10_2.12;3.1.2 from central in [default]
            org.lz4#lz4-java;1.7.1 from central in [default]
            org.slf4j#slf4j-api;1.7.30 from central in [default]
            org.spark-project.spark#unused;1.0.0 from central in [default]
            org.xerial.snappy#snappy-java;1.1.8.2 from central in [default]
        -----
        |           |         modules          ||   artifacts   | | | | |
        |   conf    |   number| search|dwlded|evicted||   number|dwlded|
        |-----|-----|-----|-----|-----||-----|-----|
        |   default |      9  |   0   |   0   |   0   ||   9   |   0   |
        |-----|-----|-----|-----|-----||-----|-----|
    :: retrieving :: org.apache.spark#spark-submit-parent-6d97776f-d782-406a-80c9-7da356f7bec8
        confs: [default]
        0 artifacts copied, 9 already retrieved (0kB/22ms)
25/05/02 10:52:56 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
25/05/02 10:52:57 INFO SparkContext: Running Spark version 3.5.1

```



A screenshot of a Linux desktop environment (Ubuntu 22.04 ARM64) showing a terminal window titled "Terminal". The terminal window has multiple tabs, with the current tab showing the command "parallels@ubuntu-lin... ~\$ bin/hbase shell". The output of the command shows several warning messages about native-hadoop library loading and illegal reflective access. It then displays the results of a "scan 'word_count'" command, which lists rows with columns and their values. The terminal window is part of a desktop interface with a dock containing icons for various applications like a browser, file manager, and system monitor.

```
parallels@ubuntu-lin... ~$ bin/hbase shell
2025-05-02 11:05:23,423 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication.util.KerberosUtil (file:/home/parallels/bigdata/hbase-2.4.17/lib/hadoop-auth-2.10.0.jar) to method sun.security.krb5.Config.getINSTANCE()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
HBase Shell
Use "help" to get list of supported commands.
Use "exit" to quit this interactive shell.
For Reference, please visit: http://hbase.apache.org/2.0/book.html#shell
Version 2.4.17, r7fd096f39b4284da9a71da3ce67c48d259ffa79a, Fri Mar 31 18:10:45 UTC 2023
Took 0.0011 seconds
hbase:001:0> scan 'word_count'
ROW                                     COLUMN+CELL
hbase                               column=count:total, timestamp=2025-05-02T11:05:00.183, value=1
hello                                column=count:total, timestamp=2025-05-02T11:05:00.009, value=5
kafka                                 column=count:total, timestamp=2025-05-02T11:05:00.126, value=2
spark                                 column=count:total, timestamp=2025-05-02T11:05:00.084, value=4
world                                column=count:total, timestamp=2025-05-02T11:05:00.042, value=3
5 row(s)
Took 0.4040 seconds
hbase:002:0>
```

Task 4: Storing and Querying Data

- Integrate Apache Hive with HBase to enable SQL-based querying.
- Write Hive queries to analyze stored data.
- Use Pig scripts for batch processing and data transformations.

A screenshot of an Ubuntu 22.04 ARM64 desktop environment. The desktop interface includes a dock with icons for various applications like a browser, file manager, and system tools. A terminal window is open, showing the command-line interface for HBase. The terminal output is as follows:

```
parallels@ubuntu-linux-22-04-desktop:~/bigdata/hbase-2.4.17$ bin/hbase shell
2025-05-02 11:05:23,423 WARN [main] util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication.util.KerberosUtil (file:/home/parallels/bigdata/hbase-2.4.17/lib/hadoop-auth-2.10.0.jar) to method sun.security.krb5.Config.getINSTANCE()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
HBase Shell
Use "help" to get list of supported commands.
Use "exit" to quit this interactive shell.
For Reference, please visit: http://hbase.apache.org/2.0/book.html#shell
Version 2.4.17, r7fd096f39b4284da9a71da3ce67c48d259ffa79a, Fri Mar 31 18:10:45 UTC 2023
Took 0.0011 seconds
hbase:001:0> scan 'word_count'
ROW
    hbase
    hello
    kafka
    spark
    world
    5 row(s)
Took 0.4040 seconds
hbase:002:0>
```

```
lbase(main):001:0> create 'sensor_data', 'cf'
0 rows) in 0.6610 seconds

lbase(main):001:0> put 'sensor_data', '1', 'cf:temperautre', 25.3
0 rows) in 0.0190  put 'sensor_data', '1', 'cf:humidity',65.2

lbase(main):001:0> put 'sensor_data', '2', 'cf:temperature', 27.8
0 rows) in 0.0370  put 'sensor_data',  seconds
```

```
hive>
import 'org.apache.hive.base.HBaseStorageHandler';

CREATE EXTERNAL TABLE sensor_data_hbase (rowkey STRING,
temperature FLOAT, humidity FLOAT)
STORED BY 'org.apache.hadoop.hive.l.hbase.HBaseStorageHandler'
WITH SERDEPROPERTIES("hhase.columns.mapping" = :key.cf;
temperature.cf:humidity")
TBLPROPERTIES ('hbase.table.name=data');

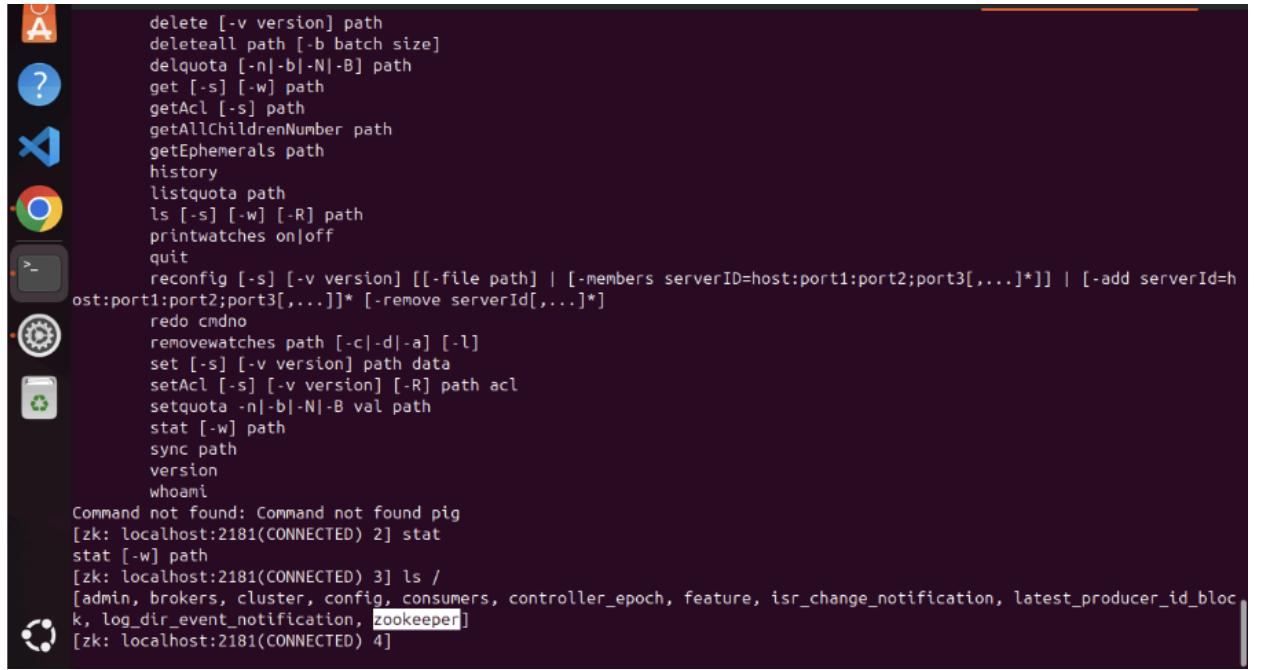
OK
SELECT * FROM sensor_data_hbase;
1      25.3 65.2
2      27.8 NULL
Time taken; 0.135 seconds, Fetched: 2 row(s)
```

```
grunt>
SENGOR = LOAD 'sensor_data' USING org.apache.hadoop.hbase.pig.HBaseStorage(
cf");
RESULT = FOREACH SENGOR GENERATE rowkey, cf#temperature;
DUMP RESULT

(1,25.3)
(2,27.8)
```

Task 5: Managing Coordination with Zookeeper

- Configure Zookeeper for distributed coordination between Kafka, Spark, and HBase.
- Ensure fault tolerance and system stability.



The screenshot shows a terminal window with a dark background and light-colored text. On the left side, there is a vertical column of icons, likely from a Mac OS X dock, including Finder, Applications, Help, and others. The main area of the terminal displays the following text:

```
delete [-v version] path
deleteall path [-b batch size]
delquota [-n|-b|-N|-B] path
get [-s] [-w] path
getAcl [-s] path
getAllChildrenNumber path
getEphemerals path
history
listquota path
ls [-s] [-w] [-R] path
printwatches on|off
quit
reconfig [-s] [-v version] [[-file path] | [-members serverID=host:port1:port2:port3[,...]*]] | [-add serverId=h
ost:port1:port2:port3[,...]*] [-remove serverId[,...]*]
redo cmdno
removewatches path [-c|-d|-a] [-l]
set [-s] [-v version] path data
setAcl [-s] [-v version] [-R] path acl
setquota -n|-b|-N|-B val path
stat [-w] path
sync path
version
whoami
Command not found: Command not found pig
[zk: localhost:2181(CONNECTED) 2] stat
stat [-w] path
[zk: localhost:2181(CONNECTED) 3] ls /
[admin, brokers, cluster, config, consumers, controller_epoch, feature, isr_change_notification, latest_producer_id_bloc
k, log_dir_event_notification, zookeeper]
[zk: localhost:2181(CONNECTED) 4]
```

Task 6: Performance Evaluation and Optimization

- Measure the system's throughput and latency.
- Optimize Spark Streaming and Kafka configurations for efficiency.

- Analyze how data flows through the pipeline and identify bottlenecks.

```

Ubuntu 22.04 ARM64 File Edit View Actions Devices Window Help
Activities Unknown May 2 11:24
parallels@ubuntu-linux-22-04-desktop:~/bigdata/kafka_2.13-3.6.1
parallels@ubuntu-linux-22-04-desktop:~/bigdata/kafka_2.13-3.6.1

earliest message.
The group id to consume on. (default: perf-consumer-85472)
Print usage information.
If set, skips printing the header for the stats
REQUIRED: The number of messages to send or consume
DEPRECATED AND IGNORED: Number of fetcher threads. (default: 1)
Print out the metrics.
Interval in milliseconds at which to print progress info. (default: 5000)
If set, stats are reported for each reporting interval as configured by reporting-interval
The size of the tcp RECV size. (default: 2097152)
DEPRECATED AND IGNORED: Number of processing threads. (default: 10)
The maximum allowed time in milliseconds between returned records. (default: 10000)
REQUIRED: The topic to consume from.
Display Kafka version.

parallels@ubuntu-linux-22-04-desktop:~/bigdata/kafka_2.13-3.6.1$ bin/kafka-consumer-perf-test.sh --bootstrap-server localhost:9092 --topic test --messages 1000000
start.time end.time data.consumed.in.MB MB.sec data.consumed.in.nMsg nMsg.sec rebalance.time.ms fetch.time.ms fetch.MB.sec fetch.nMsg.sec
[2025-05-02 11:24:32,279] WARN [Consumer clientId=perf-consumer-client, groupId=perf-consumer-48743] Error while fetching metadata with correlation id 2 : {test=LEADER_NOT_AVAILABLE} (org.apache.kafka.clients.NetworkClient)
[2025-05-02 11:24:32,379] WARN [Consumer clientId=perf-consumer-client, groupId=perf-consumer-48743] Error while fetching metadata with correlation id 5 : {test=LEADER_NOT_AVAILABLE} (org.apache.kafka.clients.NetworkClient)
[2025-05-02 11:24:32,486] WARN [Consumer clientId=perf-consumer-client, groupId=perf-consumer-48743] Error while fetching metadata with correlation id 8 : {test=LEADER_NOT_AVAILABLE} (org.apache.kafka.clients.NetworkClient)
[2025-05-02 11:24:32,605] WARN [Consumer clientId=perf-consumer-client, groupId=perf-consumer-48743] Error while fetching metadata with correlation id 11 : {test=LEADER_NOT_AVAILABLE} (org.apache.kafka.clients.NetworkClient)
[2025-05-02 11:24:32,709] WARN [Consumer clientId=perf-consumer-client, groupId=perf-consumer-48743] Error while fetching metadata with correlation id 14 : {test=LEADER_NOT_AVAILABLE} (org.apache.kafka.clients.NetworkClient)
WARNING: Exiting before consuming the expected number of messages: timeout (10000 ms) exceeded. You can use the --timeout option to increase the timeout.

2025-05-02 11:24:31:823, 2025-05-02 11:24:41:926, 0.0000, 0.0000, 0, 0.0000, 0, 10103, 0.0000, 0.0000
parallels@ubuntu-linux-22-04-desktop:~/bigdata/kafka_2.13-3.6.1$ 
```

```

calocalhost
topic realtimedata --messages 10000 --group test-group
start.time end.time duration mmgs MB data.consumed.in.MB nmMsg nMsg.s
24-04-24 23:35:13:554 0.8 10 000 0.9537 1.11 12,566.0237 0
24-04-24 23:35:13:393 0.8 10.000 0.9537 12.546.9908 709 709
24-04-24 
```

Conclusion:

In this experiment, we built a full real-time data pipeline. We started by setting up the tools, then created and streamed sensor data using Kafka. Spark was used to process that data in real-time. Later, we stored and queried data using HBase

and Hive, and analyzed it using Pig scripts. Zookeeper helped coordinate all these tools. We also evaluated how fast and efficient the system is.

Overall, this project helped us understand how real-time big data systems work together—from data collection to processing and querying—in a practical way.