



Big Data Technologies

Assignment #9

Ashmita Gupta (A20512498)



Exercise 1) 5 points

Read the article “Real-time stream processing for Big Data” available on the blackboard in the ‘Articles’ section and then answer the following questions:

- a) (1.25 points) What is the Kappa architecture and how does it differ from the lambda architecture?

Ans.

- 1) **Kappa:** As per the article on the kappa architecture, the precomputation of data does not happen on regular basis in the batch layer. All the computation is done in the stream processing system. The re-computation is performed only when the business logic changes by replaying past data. For this, the Kappa Architecture uses a powerful stream processor which can deal with data at a faster speed than the speed with what the data comes in.
 - 2) **Lambda:** Lambda architecture is a system consisting of three layers: Batch, Speed, and Serving layer. It targets both the volume and velocity challenge of big data at the same time. It has both a batch-oriented system and a real-time system.
- b) (1.25 points) What are the advantages and drawbacks of pure streaming versus micro-batch real-time processing systems?

Ans.

- 1) Storm and Samza are pure stream-oriented systems with very low latency and somewhat high per-item costs.
- 2) Batch-oriented systems offer unmatched resource efficiency at the tradeoff of unreasonably high latency for real-time applications. Data is buffered and processed in batches in micro-batch real-time processing systems. It improves efficiency while also increasing the time an individual item spends in the data flow. Storm Trident and Spark Streaming are two examples of this type.

- c) (1.25 points) In few sentences describe the data processing pipeline in Storm.

Ans.

- 1) The topology in Storm refers to a data pipeline or application. Spouts are the nodes that take data and so start the data flow in the topology. Spouts output tuples to bolts, which execute processing, write data to external storage and may transmit tuples further downstream.
- 2) Data flow between nodes is controlled by storm groupings. Storm distributes spouts and bolts in a round-robin fashion by default, but the scheduler can be modified to accommodate cases in which a specific processing step must be performed on a specific node. Storm offers the option of at least one processing via an acknowledgment feature that maintains the processing status of every

single tuple as it travels through the topology. It does not ensure the sequence in which tuples are processed.

- d) (1.25 points) How does Spark streaming shift the Spark batch processing approach to work on real-time data streams?

Ans.

It does so by chunking the stream of incoming data items into small batches, changing them into DDs, and processing them as usual, Spark streaming shifts the batch-processing method towards real-time requirements. Further, it automatically manages data flow and distribution.

Exercise 2) 5 points (extra credit; if you don't want to try or if you try and can't get things to work, this won't impact your score negatively)

Refer to the `python-Kafka` Documentation from the Free Books and Chapters section of our blackboard site

Step A – Start an EMR cluster

Starting an EMR Cluster

[illegible]

Step B – Copy the Kafka software to the EMR master node

Moving kafka file to home/hadoop directory

```

jashm@ashwin: ~/OneDrive/Desktop/Big Data$ scp -i amr-key-pair-2.pem kafka_2.13-3.0.0.tgz hadoop@ec2-3-235-107-43.compute-1.amazonaws.com:/home/hadoop
kafka_2.13-3.0.0.tgz
100% 82MB 1.3MB/s 01:01

```

Step C – Install the Kafka software and start it

Installing Kafka python package

```
[hadoop@ip-172-31-9-164 ~]$ tar -xzf kafka_2.13-3.0.0.tgz
[hadoop@ip-172-31-9-164 ~]$ ls
kafka_2.13-3.0.0  kafka_2.13-3.0.0.tgz
[hadoop@ip-172-31-9-164 ~]$ pip install kafka-python
Defaulting to user installation because normal site-packages is not writeable
Collecting kafka-python
  Downloading kafka_python-2.0.2-py2.py3-none-any.whl (246 kB)
    |████████████████████████████████████████| 246 kB 43.0 MB/s
Installing collected packages: kafka-python
Successfully installed kafka-python-2.0.2
```

Running Command- bin/zookeeper-server-start.sh config/zookeeper.properties &

[illegible]

Running Command- bin/kafka-server-start.sh config/server.properties &

```
hadoop@ip-172-31-9-164 kafka-2.13-3.0.0$ bin/kafka-server-start.sh config/server.properties &
[2] 20496
[2023-11-07 06:11:55.121] INFO Setting up log4j, rejectClientInitializetokennegotiationtrue to disable client-initiated token negotiation (org.apache.zookeeper.ZooKeeper.XS09U11)
[2023-11-07 06:11:55.217] INFO Registered signal handlers for TERM, INT, SIG (org.apache.kafka.common.utils.LoggingSignalHandler)
[2023-11-07 06:11:55.243] INFO starting (kafka.server.KafkaServer)
[2023-11-07 06:11:55.245] INFO [ZookeeperClient Kafka server] Initializing a new session to localhost:2181 (kafka.server.KafkaServer)
[2023-11-07 06:11:55.275] INFO [ZookeeperClient Kafka server] Initializetokennegotiationtrue to disable client-initiated token negotiation (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.282] INFO Client environment:zookeeper.version=3.6.3-6401e4ad2087061bcb9f80dec2d69f2a3c8660a, built on 04/08/2021 16:35 GMT (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.283] INFO Client environment:host.name=ip-172-31-9-164.ec2.internal (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.283] INFO Client environment:java.version=1.8.0_382 (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.283] INFO Client environment:java.vendor=Amazon.com Inc. (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.283] INFO Client environment:java.class.path=/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/activation-1.1.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/aopalliance-repackaged-2.6.1.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/args4j-2.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/audience-annotations-0.5.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/commons-cli-1.4.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/commons-lang3-3.8.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/connect-api-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/connect-basic-auth-extensions-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/connect-file-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/connect-jdbc-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/connect-mirror-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/connect-mirror-client-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/connect-runtime-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/connect-transtormers-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/hk2-api-2.6.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/hk2-locator-2.6.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/hk2-utils-2.6.1.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jackson-annotations-2.12.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jackson-core-2.12.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jackson-databind-2.12.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jackson-datatype-csv-2.12.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jackson-datatype-jdk8-2.12.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jackson-javrs-base-2.12.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jackson-jaxrs-json-provider-2.12.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jackson-module-jaxb-annotations-2.12.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jackson-module-scala-2.12.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jakarta.activation-api-1.2.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jakarta.annotation-api-2.1.3.5.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jakarta.inject-2.6.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jakarta.validation-api-2.0.2.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jakarta.ws.rs-api-2.1.6.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jakarta.xml.bind-api-2.3.2.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/javassist-3.27.0-GA.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/javax.servlet-api-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/javax.ws.rs-api-2.1.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jaxb-api-2.3.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jersey-client-2.34.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jersey-common-2.34.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jersey-container-servlet-2.34.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jersey-container-servlet-core-2.34.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jersey-hk2-2.34.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jersey-server-2.34.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jetty-client-9.4.43.v20210629.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jetty-continuation-9.4.43.v20210629.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jetty-http-9.4.43.v20210629.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jetty-io-9.4.43.v20210629.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jetty-security-9.4.43.v20210629.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jetty-server-9.4.43.v20210629.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jetty-util-9.4.43.v20210629.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jetty-util-ajax-9.4.43.v20210629.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jline-3.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/jost-time-5.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-2.13-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-clients-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-log4j-appender-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-metadata-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-raft-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-server-common-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-shell-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-storage-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-streams-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-streams-test-utils-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/kafka-tools-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/log4j-2.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/log4j-1.7.30.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/maven-artifact-3.8.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-codec-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-common-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-buffer-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-codec-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-common-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-handler-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-resolver-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-transport-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-transport-native-epoll-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/netty-transport-native-unix-common-4.1.12.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/osgi-resource-locator-1.0.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/paramanet-2.8.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/plexus-utils-3.2.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/reflections-0.9.12.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/rocksdbjni-6.15.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/scalacollection-compat-2.13-2.4.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/scala-java8-compat-2.13-1.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/scala-library-2.13.6.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/scala-logging-2.13-3.9.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/scala-reflect-2.13.6.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/slf4j-api-1.7.30.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/slf4j-log4j12-1.7.30.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/snappy-java-1.1.6.1.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/trogdor-3.0.0.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/zookeeper-3.6.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/zookeeper-jute-3.6.3.jar:/home/hadoop/kafka-2.13-3.0.0/bin/./:/lib/zstd-jni-1.5.0-2.jar (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.283] INFO Client environment:java.library.path=/usr/java/packages/lib/amd64:/usr/lib64:/lib:/lib64:/usr/lib (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.284] INFO Client environment:java.compiler=NA (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.284] INFO Client environment:os.name=Linux (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.284] INFO Client environment:os.arch=amd64 (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.284] INFO Client environment:os.version=4.14.322-244.539.amzn2.x86_64 (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.284] INFO Client environment:user.name=hadoop (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.284] INFO Client environment:user.dir=/home/hadoop (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.284] INFO Client environment:os.memory.free=1011MB (org.apache.zookeeper.ZooKeeper)
[2023-11-07 06:11:55.284] INFO Client environment:os.memory.max=1024MB (org.apache.zookeeper.ZooKeeper)
```

Step D – Prepare to run Kafka producers and consumers

2nd EMR Connection (Producer)

```
gashm@Ashmita MINGW64 ~/OneDrive/Desktop/Big Data
$ ssh -i emr-key-pair-2.pem hadoop@ec2-3-235-107-43.compute-1.amazonaws.com
Last login: Tue Nov 7 06:07:19 2023 from 097-069-218-126.res.spectrum.com
```

```
┌─┴─┐
└─┴─┘ Amazon Linux 2 AMI
```

```
https://aws.amazon.com/amazon-linux-2/
44 package(s) needed for security, out of 71 available
Run "sudo yum update" to apply all updates.
```

```
EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRRRRRRRRRRR
E:::::EEEEEEEEEEEE M:::::M M:::::M R:::::R
EE::::EEEEEEEEEEEE M:::::M M:::::M R:::::R
E::::E EEEEE M:::::M M:::::M R:::::R
E::::E M:::::M M:::::M R:::::R
E::::EEEEEEEEEEEE M:::::M M:::::M R:::::R
E:::::EEEEEEEEEEEE M:::::M M:::::M R:::::R
E::::E M:::::M M:::::M R:::::R
E::::E EEEEE M:::::M M:::::M R:::::R
EE::::EEEEEEEEEEEE M:::::M M:::::M R:::::R
E:::::EEEEEEEEEEEE M:::::M M:::::M R:::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRRR RRRRRR
```


3rd EMR Connection (Consumer)

```
gashm@Ashmita MINGW64 ~/OneDrive/Desktop/Big Data
$ ssh -i emr-key-pair-2.pem hadoop@ec2-3-235-107-43.compute-1.amazonaws.com
Last login: Tue Nov  7 06:07:19 2023 from 097-069-218-126.res.spectrum.com
```

```
  _ | _ | _ )
 _ | ( _ | /
 _ | \ _ | _ |
      Amazon Linux 2 AMI
```

```
https://aws.amazon.com/amazon-linux-2/
44 package(s) needed for security, out of 71 available
Run "sudo yum update" to apply all updates.
```

```
EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRRRRRRRRRR
E::::::::::::::::::::E M::::::::M M::::::::M R:::::::::R
EE::::::::EEEEEEEE::E M::::::::M M::::::::M R:::::::::R
E::::E EEEEE M::::::::M M::::::::M RR::::R R::::R
E::::E M::::::::M M::M::::M R::R R::R
E::::EEEEEEEE M:::M M::M M::M M::M R::RRRRR::::R
E::::::::::::E M:::M M::M::::M M::M R:::::::::RR
E::::EEEEEEEE M:::M M::M M::M R::RRRRR::::R
E::::E M:::M M::M M::M R::R R::R
E::::E EEEEE M:::M MMM M::M R::R R::R
EE::::::::EEEEEEEE::E M:::M M::M R::R R::R
E::::::::::::E M:::M M::M RR::::R R::R
EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRR RRRRR
```

Step E – Creating a Kafka topics

```
[hadoop@ip-172-31-9-164 ~]$ cd kafka_2.13-3.0.0
[hadoop@ip-172-31-9-164 kafka_2.13-3.0.0]$ bin/kafka-topics.sh --create --replication-factor 1 --partitions 1 --bootstrap-server localhost:9092 --topic sample
Created topic sample.
[hadoop@ip-172-31-9-164 kafka_2.13-3.0.0]$ bin/kafka-topics.sh --create --replication-factor 1 --partitions 1 --bootstrap-server localhost:9092 --topic ashmita_gupta
WARNING: Due to limitations in metric names, topics with a period('.') or underscore('_') could collide. To avoid issues it is best to use either, but not both.
Created topic ashmita_gupta.
[hadoop@ip-172-31-9-164 kafka_2.13-3.0.0]$ bin/kafka-topics.sh --list --bootstrap-server localhost:9092
ashmita_gupta
sample
```

Producer Terminal

Key	Value
'MYID'	Your student id
'MYNAME'	Your name
'MYEYECOLOR'	Your eye color (make it up if you can't remember)

vi put.py

```

hadoop@ip-172-31-9-164:~/kafka_2.13-3.0.0
from time import sleep
from json import dumps
from kafka import KafkaProducer

Producer = KafkaProducer(bootstrap_servers = ['localhost:9092'], value_serializer = lambda x: dumps(x).encode('utf-8'))

synmyid = 'MYID'
synmyname = 'MYNAME'
synmyeyecolor = "MYEYECOLOR"

realid = input("Enter your ID: ")
realname = input("Enter your name: ")
realeyecolor = input("Enter your eye color: ")

my_dict = {}
my_dict[synmyid] = realid

my_dict1 = {}
my_dict1[synmyname] = realname

my_dict2 = {}
my_dict2[synmyeyecolor] = realeyecolor

myid = my_dict
Producer.send('sample', myid)
sleep(4)

myname = my_dict1
Producer.send('sample', myname)
sleep(4)

myeyecolor = my_dict2
Producer.send('sample', myeyecolor)
sleep(4)

Producer.close()

"put.py" 37L, 729B

```

```
python put.py
```

```
[hadoop@ip-172-31-9-164 kafka_2.13-3.0.0]$ python put.py
Enter your ID: A20512498
Enter your name: Ashmita Gupta
Enter your eye color: Black
```

Consumer Terminal vi get.py

```
hadoop@ip-172-31-9-164:~/kafka_2.13-3.0.0
from ensurepip import bootstrap
from kafka import KafkaConsumer
from json import loads

Consumer = KafkaConsumer('sample', bootstrap_servers = ['localhost:9092'], auto_offset_reset= 'earliest', enable_auto_commit=True, group_id='my-group', value_deserializer = lambda x:loads(x.decode('utf-8')))

for i in Consumer:
    for key, value in i.value.items():
        print ("key=%s value=%s" % (key, value))
Consumer.close()
```

"get.py" 11L, 422B

5,24 A11

python get.py

```
[hadoop@ip-172-31-9-164 kafka_2.13-3.0.0]$ python get.py
key=MYID value=A20512498
key=MYNAME value=Ashmita Gupta
key=MYEYECOLOR value=Black
```