# Big Data Technologies-CSP554

## Assignment-9

**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?**
As mentioned in 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 with what the data comes in.
On the other hand, 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?**
Storm and Samza are pure stream-oriented systems with very low latency and somewhat high per-item costs, on the other hand, 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.**
A 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. 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?**
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. It also automatically manages data flow and distribution.

**Exercise 2)**
**Starting an EMR Cluster**

```
aasth@LAPTOP-HJTR6HMR MINGW64 ~
$ ssh -i Downloads/emr-key-pair.pem hadoop@ec2-54-209-245-237.compute-1.amazonaw
s.com
The authenticity of host 'ec2-54-209-245-237.compute-1.amazonaws.com (54.209.245
.237)' can't be established.
ED25519 key fingerprint is SHA256:pexzlzUrwKO6WwcDn5+kxGL7ReZks9p/zwuUoDUj1/E.
This key is not known by any other names
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'ec2-54-209-245-237.compute-1.amazonaws.com' (ED25519
) to the list of known hosts.

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

https://aws.amazon.com/amazon-linux-2/
7 package(s) needed for security, out of 10 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEEE MMMMMMMM             MMMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::E M:::::::M           M:::::::M R::::::::::::::R
EE::::EEEEEEEEE::::E M::::::::M           M::::::::M R:::::RRRRRR:::::R
  E::::E       EEEEE M:::::::::M         M:::::::::M RR::::R      R::::R
  E::::E             M::::::M::M       M:::M::::::M   R:::R      R::::R
  E:::::EEEEEEEEEE    M:::::M M:::M   M:::M M:::::M   R:::RRRRRR:::::R
  E::::::::::::::E    M:::::M  M:::M:::M  M:::::M   R:::::::::::RR
  E:::::EEEEEEEEEE    M:::::M   M:::::M   M:::::M   R:::RRRRRR::::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
EEEEEEEEEEEEEEEEEEEEE MMMMMMM             MMMMMMM RRRRRRR      RRRRRR

[hadoop@ip-172-31-50-139 ~]$
```

**Moving kafka file to home/hadoop directory**

```
aasth@LAPTOP-HJTR6HMR MINGW64 ~
$ scp -i Downloads/emr-key-pair.pem Downloads/kafka_2.13-3.0.0.tgz hadoop@ec2-54
-209-245-237.compute-1.amazonaws.com:/home/hadoop
kafka_2.13-3.0.0.tgz                        100%   82MB 740.0KB/s   01:54

aasth@LAPTOP-HJTR6HMR MINGW64 ~
$
```

**Extracting the Kafka Package**

```
[hadoop@ip-172-31-50-139 ~]$ ls
[hadoop@ip-172-31-50-139 ~]$ ls
kafka_2.13-3.0.0.tgz
[hadoop@ip-172-31-50-139 ~]$ tar -xzf kafka_2.13-3.0.0.tgz
[hadoop@ip-172-31-50-139 ~]$ ls
kafka_2.13-3.0.0   kafka_2.13-3.0.0.tgz
```

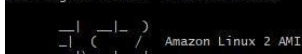**Installing Kafka python package**

```
[hadoop@ip-172-31-50-139 ~]$ python --version
Python 3.7.10
[hadoop@ip-172-31-50-139 ~]$ pip3 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 29.3 MB/s
Installing collected packages: kafka-python
Successfully installed kafka-python-2.0.2
[hadoop@ip-172-31-50-139 ~]$
```

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



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



## At Producer Terminal
## 2nd EMR Connection

**Topic Sample Created**

```
[hadoop@ip-172-31-50-139 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-172-31-50-139 kafka_2.13-3.0.0]$ |
```

**More Topics created**

```
[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$ bin/kafka-topics.sh --create --replication-factor 1 --partitions 1 --bootstrap-server localhost:9092 --topic aastha_dhir
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 aastha_dhir.
[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$ bin/kafka-topics.sh --create --replication-factor 1 --partitions 1 --bootstrap-server localhost:9092 --topic aasthadhir23
Created topic aasthadhir23.
[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$ |
```

**List of all topics created**

```
[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$ bin/kafka-topics.sh --list --bootstrap-server localhost:9092
aastha_dhir
aasthadhir23
sample
[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$ |
```

**Producer Terminal**
**vim put.py**

hadoop@ip-172-31-50-139:~/kafka_2.13-3.0.0

```python
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()
~
~
~
```

**phython put.py**

```
[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$ python3 put.py
Enter your ID: A20468022
Enter your name: Aastha Dhir
Enter your eye color: Brown

[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$
[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$ |
```

**At Consumer Terminal**
**3rd EMR Connection**

```
aasth@LAPTOP-HJTR6HMR MINGW64 ~
$ ssh -i Downloads/emr-key-pair.pem hadoop@ec2-54-209-245-237.compute-1.amazonaw
s.com
Last login: Tue Nov  8 03:51:39 2022

      __|  __|_  )
      _|  (     /   Amazon Linux 2 AMI
     ___|\___|___|

https://aws.amazon.com/amazon-linux-2/
7 package(s) needed for security, out of 10 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEEE MMMMMMMM           MMMMMMMM RRRRRRRRRRRRRRR
E::::::::::::::::::::E M::::::M           M::::::M R::::::::::::::R
EE:::::EEEEEEEEE:::E M:::::::M           M:::::::M R:::::RRRRRR:::::R
  E::::E       EEEEE M::::::::M         M::::::::M RR::::R      R::::R
  E::::E             M:::::::M:::M     M:::M:::::M   R:::R      R::::R
  E:::::EEEEEEEEEE    M:::::M M:::M   M:::M M:::::M   R:::RRRRRR:::::R
  E::::::::::::::E    M:::::M  M:::M:::M  M:::::M   R:::::::::::RR
  E:::::EEEEEEEEEE    M:::::M   M:::::M   M:::::M   R:::RRRRRR::::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
EEEEEEEEEEEEEEEEEEEEE MMMMMMM           MMMMMMM RRRRRRR      RRRRRR

[hadoop@ip-172-31-50-139 ~]$ |
```

**Consumer Terminal**
**vim get.py**

```
hadoop@ip-172-31-50-139:~/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_o
ffset_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()
~
~
~
~
```

**phython3 get.py**

```
[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$ vim get.py
[hadoop@ip-172-31-50-139 kafka_2.13-3.0.0]$ python3 get.py
key=MYID value=A20468022
key=MYNAME value=Aastha Dhir
key=MYEYECOLOR value=Brown
```

**Submitted By:-**
**Aastha Dhir**
**CWID- A20468022**
**adhir2@hawk.iit.edu**