

Assignment - 9

Q1) What is the Kappa architecture and how does it differ from the lambda architecture?

Sol) As per the article in the kappa architecture, precomputation of data does not happen on regular basis in the batch layer. It is the stream processing system in which all the computation is done. The re-computation is performed only when the business logic changes by replaying past data. For this, the Kappa Architecture uses powerful stream processor which can deal with data at faster speed than the data comes in.

Whereas Lambda architecture is a system consisting of three layers: Batch, Speed and Serving layer. It targets both Volume and Velocity challenge of big data at the same time. It has both batch-oriented system and real-time system.

Q2) What are the advantages and drawbacks of pure streaming versus micro-batch real-time processing systems?

Sol) Storm and Samza are pure stream-oriented systems with very low latency and somewhat high per-item costs, whereas 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 lengthening the time an individual item spends in the data flow. Storm Trident and Spark Streaming are two examples of this type.

Q3) In few sentences describe the data processing pipeline in Storm.

Sol) A topology in Storm refers to a data pipeline or application. Spouts are the nodes that intake 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 for cases in which a specific processing step must be performed on a specific node.

Storm, on the other hand, does not ensure the sequence in which tuples are processed; however, it does offer the option of at-least-once processing via an acknowledgement feature that maintains the processing status of every single tuple as it travels through the topology.

Q4) How does Spark streaming shift the Spark batch processing approach to work on real-time data streams?

Sol) By chunking the stream of incoming data items into small batches, translating 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.

Rishabh Jain
CSP-554
A20495530
Extra Credit exercise

Starting a EMR cluster

```

Desktop — hadoop@ip-172-31-10-170:~ — ssh -i ~/Desktop/new-key-pair-emr.pem hadoop@ec2-18-209-241-14.compute-1.amazonaws.com — 165x47
+ Desktop ssh -i ~/Desktop/new-key-pair-emr.pem hadoop@ec2-18-209-241-14.compute-1.amazonaws.com

  _ _  _ _  _
 _ _  _ _  _ /   Amazon Linux 2 AMI

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

EEEEEEEEEEEEEEEEEEEE MMMMMMMM      MMMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::E M::::::::M      M::::::::M R::::::::::::R
EE::::::::EEEEEEEE::::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 R::RRRRRR::::R
E::::E      M::::::::M M::::::::M M::::::::M R::RRRRRR::::R
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      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 RRRRRRR      RRRRRR

[hadoop@ip-172-31-10-170 ~]$

```

Extracting the Kafka Package

```
[hadoop@ip-172-31-10-170 ~]$ ls
kafka_2.13-3.0.0.tgz
[hadoop@ip-172-31-10-170 ~]$ tar -xzf kafka_2.13-3.0.0.tgz
[hadoop@ip-172-31-10-170 ~]$ ls
kafka_2.13-3.0.0  kafka_2.13-3.0.0.tgz
[hadoop@ip-172-31-10-170 ~]$
```

Installing the kafka-python package

```
[hadoop@ip-172-31-10-170 ~]$ python --version
Python 3.7.10
[hadoop@ip-172-31-10-170 ~]$ 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 16.6 MB/s
Installing collected packages: kafka-python
Successfully installed kafka-python-2.0.2
[hadoop@ip-172-31-10-170 ~]$
```

A20495530

```
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$ bin/zookeeper-server-start.sh config/zookeeper.properties &
[1] 7084
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$ [2022-03-29 01:57:50.068] INFO Reading configuration from: config/zookeeper.properties (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.074] WARN config/zookeeper.properties is relative. Prepend ./ to indicate that you're sure! (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.081] INFO clientPortAddress is 0.0.0.0:2181 (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.081] INFO secureClientPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.082] INFO observerMasterPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.082] INFO metricsProvider.className is org.apache.zookeeper.metrics.impl.DefaultMetricsProvider (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.089] INFO autopurge.snapRetainCount set to 3 (org.apache.zookeeper.server.DataDirCleanupManager)
[2022-03-29 01:57:50.089] INFO autopurge.purgeInterval set to 0 (org.apache.zookeeper.server.DataDirCleanupManager)
[2022-03-29 01:57:50.089] INFO Purge task is not scheduled. (org.apache.zookeeper.server.DataDirCleanupManager)
[2022-03-29 01:57:50.089] WARN Either no config or no quorum defined in config, running in standalone mode (org.apache.zookeeper.server.quorum.QuorumPeerMain)
[2022-03-29 01:57:50.101] INFO Log4j 1.2 jmx support found and enabled. (org.apache.zookeeper.jmx.ManagedUtil)
[2022-03-29 01:57:50.142] INFO Reading configuration from: config/zookeeper.properties (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.143] WARN config/zookeeper.properties is relative. Prepend ./ to indicate that you're sure! (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.144] INFO clientPortAddress is 0.0.0.0:2181 (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.144] INFO secureClientPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.144] INFO observerMasterPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.145] INFO metricsProvider.className is org.apache.zookeeper.metrics.impl.DefaultMetricsProvider (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2022-03-29 01:57:50.145] INFO Starting server (org.apache.zookeeper.server.ZooKeeperServerMain)
[2022-03-29 01:57:50.214] INFO ServerMetrics initialized with provider org.apache.zookeeper.metrics.impl.DefaultMetricsProvider@9f7f0c54 (org.apache.zookeeper.server.ServerMetrics)
[2022-03-29 01:57:50.240] INFO zookeeper.snapshot.trust.empty : false (org.apache.zookeeper.server.persistence.FileTxnSnapLog)
[2022-03-29 01:57:50.298] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.298] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.298] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.298] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.299] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.299] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.299] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.300] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.300] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.300] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2022-03-29 01:57:50.300] INFO (org.apache.zookeeper.server.ZooKeeperServer)
```

```
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$ bin/kafka-server-start.sh config/server.properties &
[2] 11922
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$ [2022-03-29 01:59:11,546] INFO Registered kafka:type=kafka.Log4jController MBean (kafka.utils.Log4jControllerRegistration$)
[2022-03-29 01:59:12,158] INFO Setting -Djdk.tls.rejectClientInitiatedRenegotiation=true to disable client-initiated TLS renegotiation (org.apache.zookeeper.common.XS509Util)
[2022-03-29 01:59:12,357] INFO Registered signal handlers for TERM, INT, HUP (org.apache.kafka.common.utils.LoggingSignalHandler)
[2022-03-29 01:59:12,383] INFO starting (kafka.server.KafkaServer)
[2022-03-29 01:59:12,384] INFO Connecting to zookeeper on localhost:2181 (kafka.server.KafkaServer)
[2022-03-29 01:59:12,423] INFO [ZooKeeperClient Kafka server] Initializing a new connection to localhost:2181. (kafka.zookeeper.ZooKeeperClient)
[2022-03-29 01:59:12,434] INFO Client environment:zookeeper.version=3.6.3--6401e4ad2087061bc6bf9f8dec2d69f2e3c8660a, built on 04/08/2021 16:35 GMT (org.apache.zookeeper.ZooKeeper)
[2022-03-29 01:59:12,434] INFO Client environment:host.name=ip-172-31-10-170.ec2.internal (org.apache.zookeeper.ZooKeeper)
[2022-03-29 01:59:12,435] INFO Client environment:java.version=1.8.0_322 (org.apache.zookeeper.ZooKeeper)
[2022-03-29 01:59:12,435] INFO Client environment:java.vendor=Amazon.com Inc. (org.apache.zookeeper.ZooKeeper)
[2022-03-29 01:59:12,435] INFO Client environment:java.home=/usr/lib/jvm/java-1.8.0-amazon-corretto.x86_64/jre (org.apache.zookeeper.ZooKeeper)
[2022-03-29 01:59:12,435] INFO Client environment:java.class.path=/home/hadoop/kafka_2.13-3.0.0/bin/../libs/activation-1.1.1.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/aopalliance-repackaged-2.6.1.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/argparse4j-0.7.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/audience-annotations-5.0.8.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/commons-clui-1.4.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/commons-lang3-3.8.1.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/connect-api-3.0.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/connect-basic-auth-extension-3.0.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/connect-file-3.0.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/connect-json-3.0.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/connect-mirror-3.0.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/connect-mirror-client-3.0.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/connect-runtime-3.0.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/connect-transforms-3.0.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/hk2-kafka-2.1.6.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/hk2-locator-2.6.1.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/hk2-utils-2.6.1.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jackson-annotations-2.12.3.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jackson-core-2.12.3.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jackson-databind-2.12.3.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jackson-dataformat-csv-2.12.3.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jackson-datatype-jdk8-2.12.3.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jackson-jaxrs-base-2.12.3.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jackson-jaxrs-provider-2.12.3.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jackson-module-jaxb-annotations-2.12.3.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jackson-module-scala-2.12-12.3.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jakarta.activation-api-1.2.1.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jakarta.annotation-api-1.3.5.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jakarta.inject-2.6.1.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jakarta.validation-api-2.0.2.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jakarta.ws.rs-api-2.1.6.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jakarta.xml.bind-api-2.3.2.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/javassist-3.27.0-GA.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/javax.servlet-api-3.1.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/javax.ws.rs-api-2.1.1.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jaxb-api-2.3.0.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jersey-client-2.34.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jersey-common-2.34.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jersey-container-servlet-2.34.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jersey-container-servlet-core-2.34.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jersey-hk2-2.34.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jetty-server-2.34.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jetty-client-9.4.43.v20210629.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jetty-continuation-9.4.43.v20210629.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jetty-http-9.4.43.v20210629.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jetty-io-9.4.43.v20210629.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jetty-security-9.4.43.v20210629.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jetty-server-9.4.43.v20210629.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jetty-servlet-9.4.43.v20210629.jar:/home/hadoop/kafka_2.13-3.0.0/bin/../libs/jetty-servlets-9.4.43.v20210629.jar:/home/hadoop/kafka
```

Rishabh Jain

CSP-554

A20495530

At Producer-terminal

Creating a topic named "sample"

Running command: `bin/kafka-topics.sh --create --replication-factor 1 --partitions 1 --bootstrap-server localhost:9092 --topic sample`

```
Desktop ssh -i ~/Desktop/new-key-pair-emr.pem hadoop@ec2-18-209-241-14.compute-1.amazonaws.com

  _| _|_ )
  _| ( _| /
  _|\_|\_|\

Amazon Linux 2 AMI

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

EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::E M::::::::M M::::::::M R:::::::::R
EE::::::::EEEEEEEE::E M::::::::M M::::::::M R::::RRRRR:::R
E:::E EEEEE M::::::::M M::::::::M RR:::R R:::R
E:::E M::::::::M M::M M::M M::M R:::R R:::R
E:::EEEEEEEEEE M:::M M::M M::M M::M R::RRRRR:::R
E::::::::EEEEEE M:::M M::M M::M M::M R:::RR:::RR
E:::EEEEEEEEEE M:::M 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::::::::EEEEEE M:::M M:::M RR:::R R:::R
EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRRR RRRRRR

[hadoop@ip-172-31-10-170 ~]$ cd kafka_2.13-3.0.0
[hadoop@ip-172-31-10-170 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-172-31-10-170 kafka_2.13-3.0.0]$
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$
```

More topics created

```
Desktop — hadoop@ip-172-31-10-170:~/kafka_2.13-3.0.0 — ssh -i ~/Desktop/new-key-pair-emr.pem hadoop@ec2-18-209-241-14....
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$ bin/kafka-topics.sh --create --replication-factor 1 --partitions 1 --bootstrap-server localhost:9092 --topic rishabh_jain
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 rishabh_jain.
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$ bin/kafka-topics.sh --create --replication-factor 1 --partitions 1 --bootstrap-server localhost:9092 --topic rishabh507
Created topic rishabh507.
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$
```

Rishabh Jain

CSP-554

A20495530

Listing all topics

```
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$ bin/kafka-topics.sh --list --bootstrap-server localhost:9092
rishabh507
rishabh_jain
sample
[hadoop@ip-172-31-10-170 kafka_2.13-3.0.0]$
```

On producer terminal

vim put.py (code file attached with the submission on blackboard).

```
1 from time import sleep
2 from json import dumps
3 from kafka import KafkaProducer
4
5 Producer = KafkaProducer(bootstrap_servers = ['localhost:9092'],value_serializer = lambda x: dumps(x).encode('utf-8'))
6
7 synmyid = 'MYID'
8 synmyname = 'MYNAME'
9 synmyeyecolor = 'MYEYECOLOR'
10
11 realid = input("Enter your ID: ")
12 realname = input("Enter your name: ")
13 realeyecolor = input("Enter your eye color: ")
14
15 my_dict = {}
16 my_dict[synmyid] = realid
17
18 my_dict1 = {}
19 my_dict1[synmyname] = realname
20
21 my_dict2 = {}
22 my_dict2[synmyeyecolor] = realeyecolor
23
24 myid = my_dict
25 Producer.send('sample',myid)
26 sleep(4)
27
28 myname = my_dict1
29 Producer.send('sample',myname)
30 sleep(4)
31
32 myeyecolor = my_dict2
33 Producer.send('sample',myeyecolor)
34 sleep(4)
35
36 Producer.close()
```

Python3 put.py

```
Desktop — hadoop@ip-172-31-15-220:~/kafka_2.13-3.0.0 — ssh -i ~/Desktop/new-key-pair-emr.pem hadoop@ec2-44-195-77
[hadoop@ip-172-31-15-220 kafka_2.13-3.0.0]$ python3 put.py
Enter your ID: A20495530
Enter your name: Rishabh Jain
Enter your eye color: Brown
[hadoop@ip-172-31-15-220 kafka_2.13-3.0.0]$
```

Rishabh Jain

CSP-554

A20495530

On consumer Terminal

vim get.py (code file attached with the submission on blackboard).

```
1 from ensurepip import bootstrap
2 from kafka import KafkaConsumer
3 from json import loads
4
5 Consumer = KafkaConsumer('sample', bootstrap_servers = ['localhost:9092'], auto_offset_reset='earliest', enable_auto_commit=True,
6 group_id='my-group', value_deserializer = lambda x:loads(x.decode('utf-8')))
7
8 for i in Consumer:
9     for key, value in i.value.items():
10         print("key=%s value=%s" % (key,value))
11
12
13 Consumer.close()
```

Python3 get.py

```
[hadoop@ip-172-31-15-220 kafka_2.13-3.0.0]$ python3 get.py
key=MYID value=A20495530
key=MYNAME value=Rishabh Jain
key=MYEYECOLOR value=Brown
[]
```