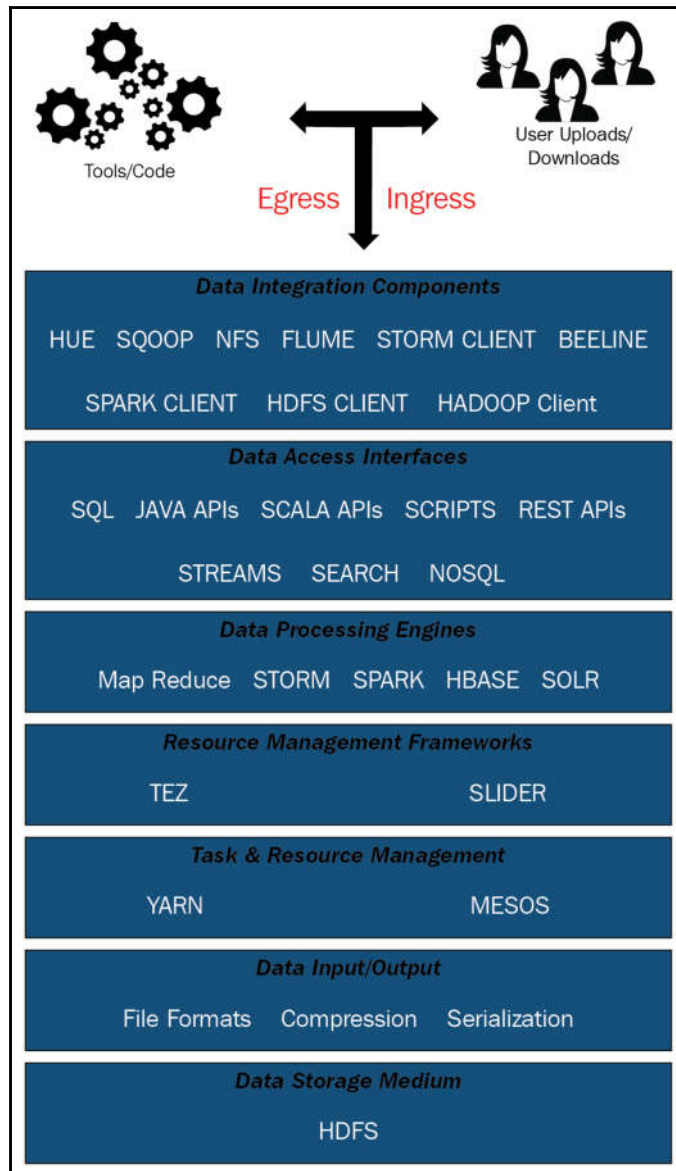
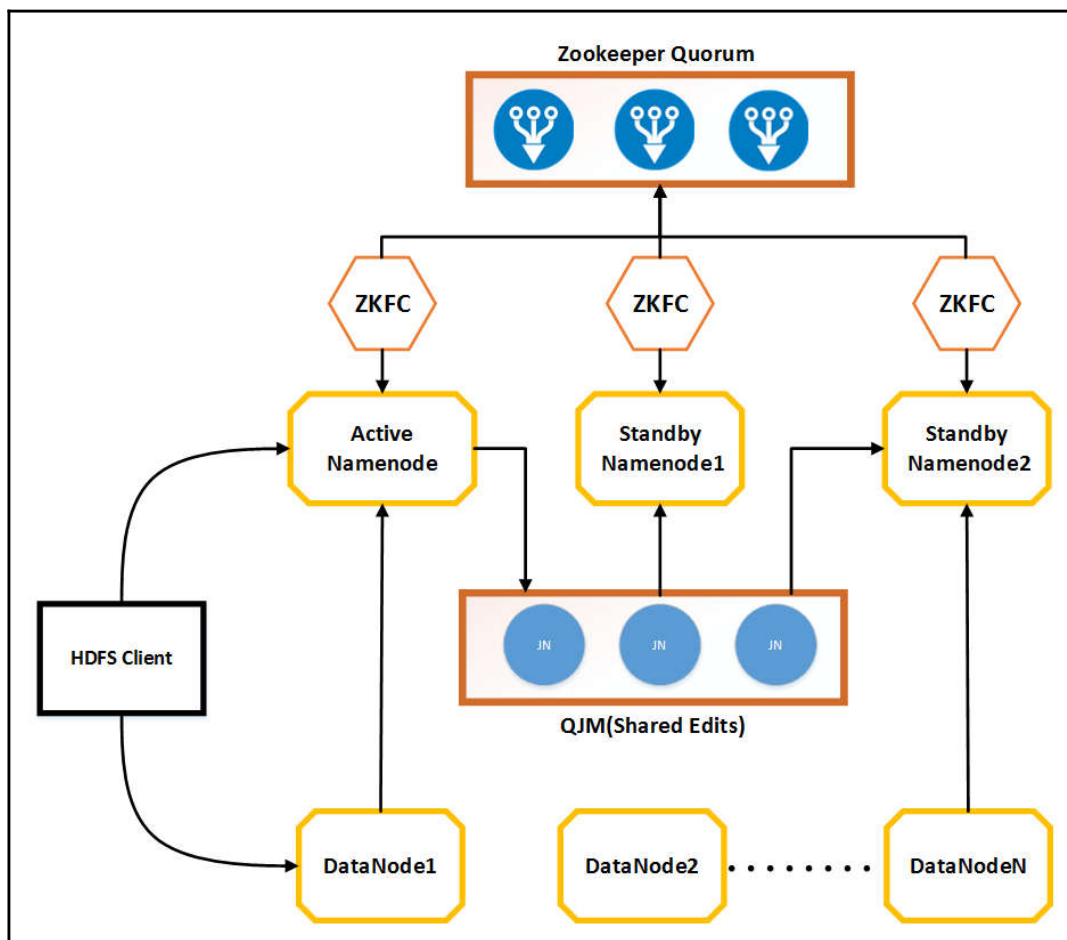
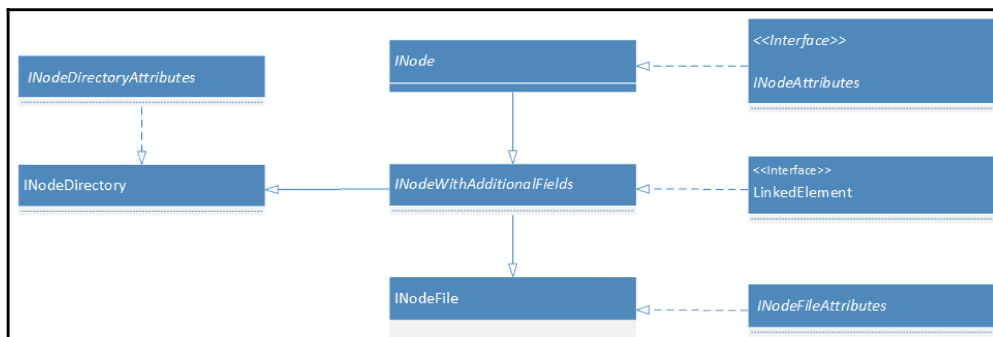
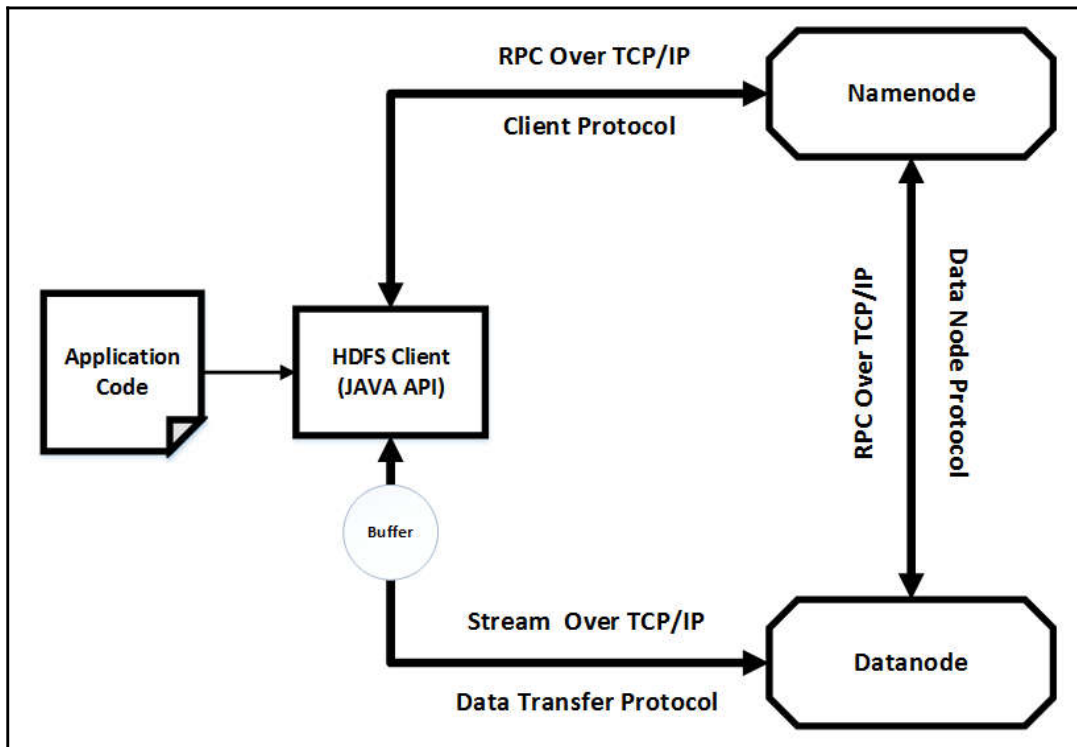


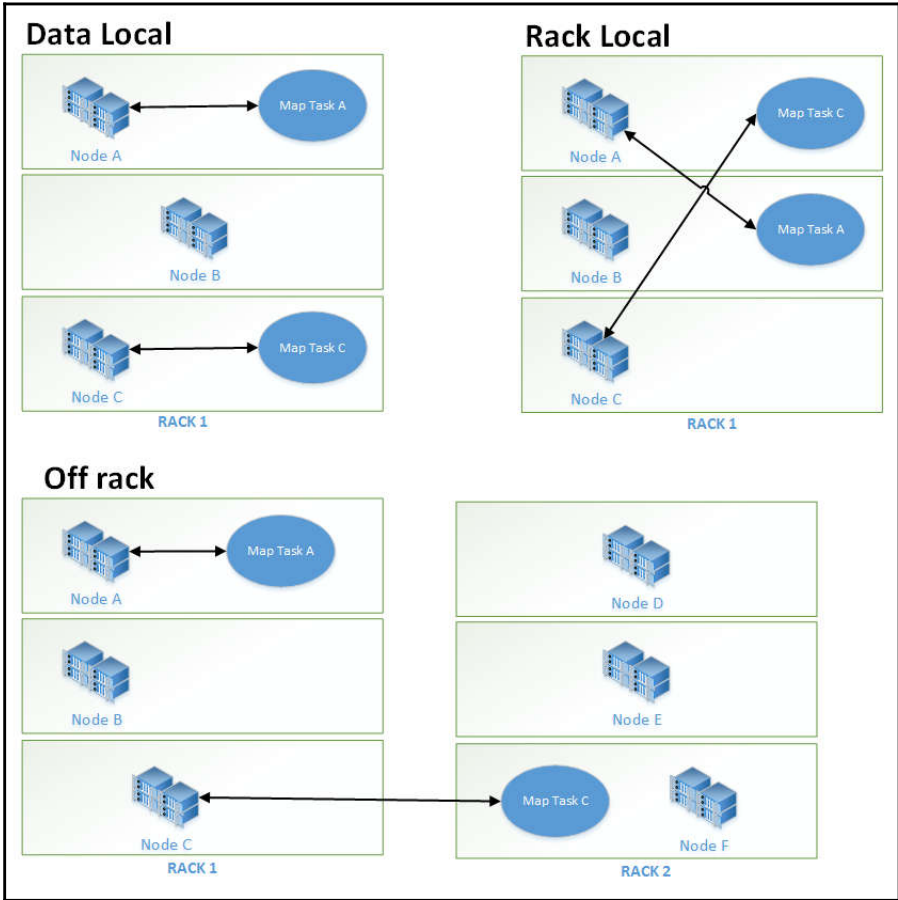
Chapter 1: Journey to Hadoop 3



Chapter 2: Deep Dive into the Hadoop Distributed File System







```
chanchal@chanchal-Lenovo-ideapad-510-15IKB:~$ hdfs oiv --help
```

```
Error parsing command-line options:
```

```
Usage: bin/hdfs oiv [OPTIONS] -i INPUTFILE -o OUTPUTFILE
```

```
Offline Image Viewer
```

```
View a Hadoop fsimage INPUTFILE using the specified PROCESSOR,  
saving the results in OUTPUTFILE.
```

The oiv utility will attempt to parse correctly formed image files and will abort fail with mal-formed image files.

The tool works offline and does not require a running cluster in order to process an image file.

The following image processors are available:

- * XML: This processor creates an XML document with all elements of the fsimage enumerated, suitable for further analysis by XML tools.
- * FileDistribution: This processor analyzes the file size distribution in the image.
 - maxSize specifies the range [0, maxSize] of file sizes to be analyzed (128GB by default).
 - step defines the granularity of the distribution. (2MB by default)
- * Web: Run a viewer to expose read-only WebHDFS API.
 - addr specifies the address to listen. (localhost:5978 by default)
- * Delimited (experimental): Generate a text file with all of the elements common to both inodes and inodes-under-construction, separated by a delimiter. The default delimiter is \t, though this may be changed via the -delimiter argument.

Required command line arguments:

```
-i,--inputFile <arg> FSImage file to process.
```

Optional command line arguments:

```
-o,--outputFile <arg> Name of output file. If the specified  
file exists, it will be overwritten.  
(output to stdout by default)
```

```
-p,--processor <arg> Select which type of processor to apply  
against image file. (XML|FileDistribution|Web|Delimited)  
(Web by default)
```

```
-delimiter <arg> Delimiting string to use with Delimited processor.
```

```
-t,--temp <arg> Use temporary dir to cache intermediate result to generate  
Delimited outputs. If not set, Delimited processor constructs  
the namespace in memory before outputting text.
```

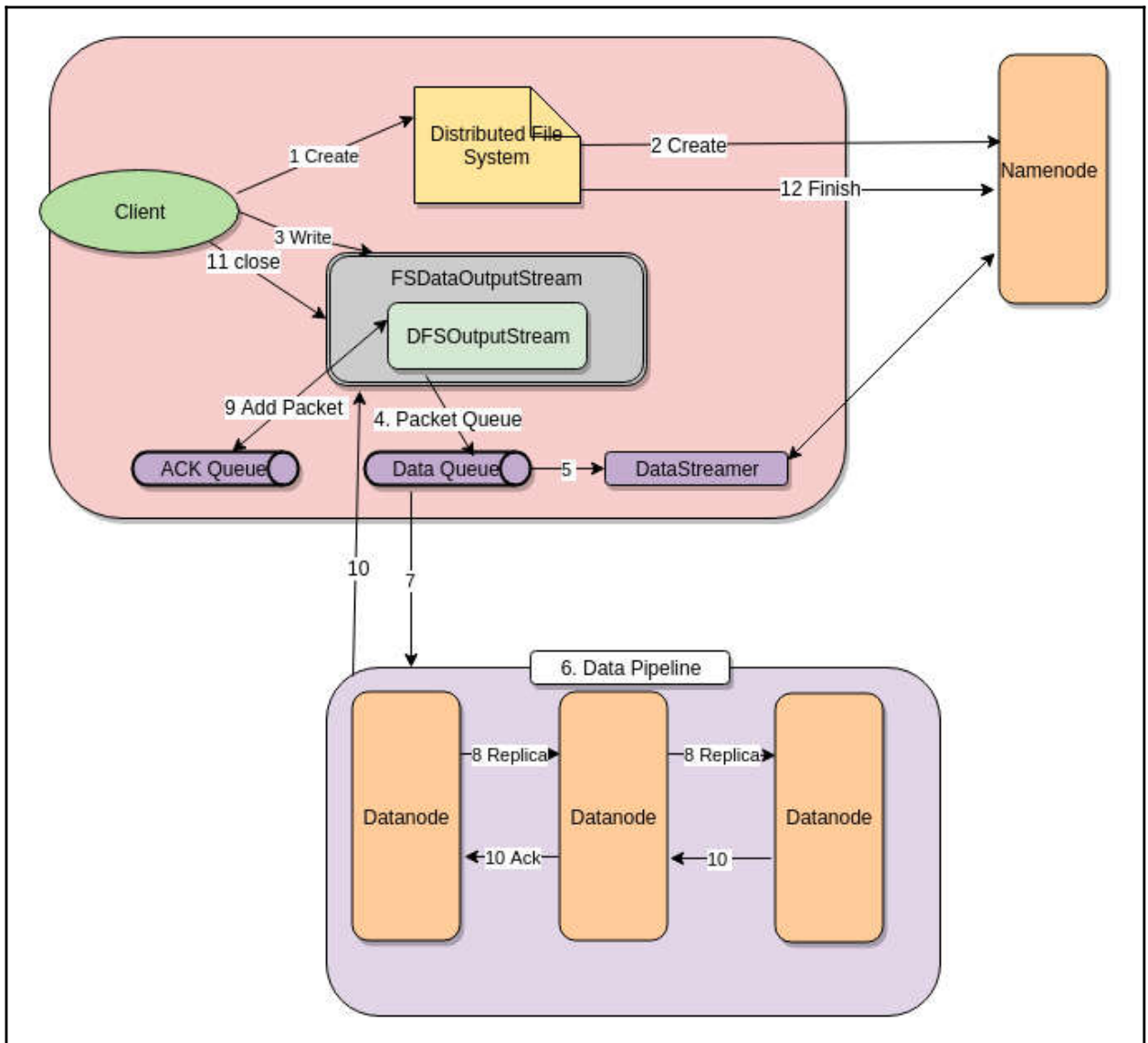
```
<?xml version="1.0" encoding="UTF-8"?>
<EDITS>
  <EDITS_VERSION>-63</EDITS_VERSION>
  <RECORD>
    <OPCODE>OP_START_LOG_SEGMENT</OPCODE>
    <DATA>
      <TXID>488053</TXID>
    </DATA>
  </RECORD>
  <RECORD>
    <OPCODE>OP_MKDIR</OPCODE>
    <DATA>
      <TXID>488054</TXID>
      <LENGTH>0</LENGTH>
      <INODEID>190335</INODEID>
      <PATH>/tmp/hive/hive/124dd7e2-d4d3-413e-838e-3dbbbd185a69</PATH>
      <TIMESTAMP>1509663411129</TIMESTAMP>
      <PERMISSION_STATUS>
        <USERNAME>hive</USERNAME>
        <GROUPNAME>hdfs</GROUPNAME>
        <MODE>448</MODE>
      </PERMISSION_STATUS>
    </DATA>
  </RECORD>
  <RECORD>
    <OPCODE>OP_ADD</OPCODE>
    <DATA>
      <TXID>488055</TXID>
      <LENGTH>0</LENGTH>
      <INODEID>190336</INODEID>
      <PATH>/tmp/hive/hive/124dd7e2-d4d3-413e-838e-3dbbbd185a69/inuse.info</PATH>
      <REPLICATION>3</REPLICATION>
      <MTIME>1509663411169</MTIME>
      <ATIME>1509663411169</ATIME>
      <BLOCKSIZE>134217728</BLOCKSIZE>
      <CLIENT_NAME>DFSClient_NONMAPREDUCE_1006023362_1</CLIENT_NAME>
      <CLIENT_MACHINE>10.1.2.26</CLIENT_MACHINE>
      <OVERWRITE>>true</OVERWRITE>
      <PERMISSION_STATUS>
        <USERNAME>hive</USERNAME>
        <GROUPNAME>hdfs</GROUPNAME>
        <MODE>420</MODE>
      </PERMISSION_STATUS>
    </DATA>
  </RECORD>
</EDITS>
```

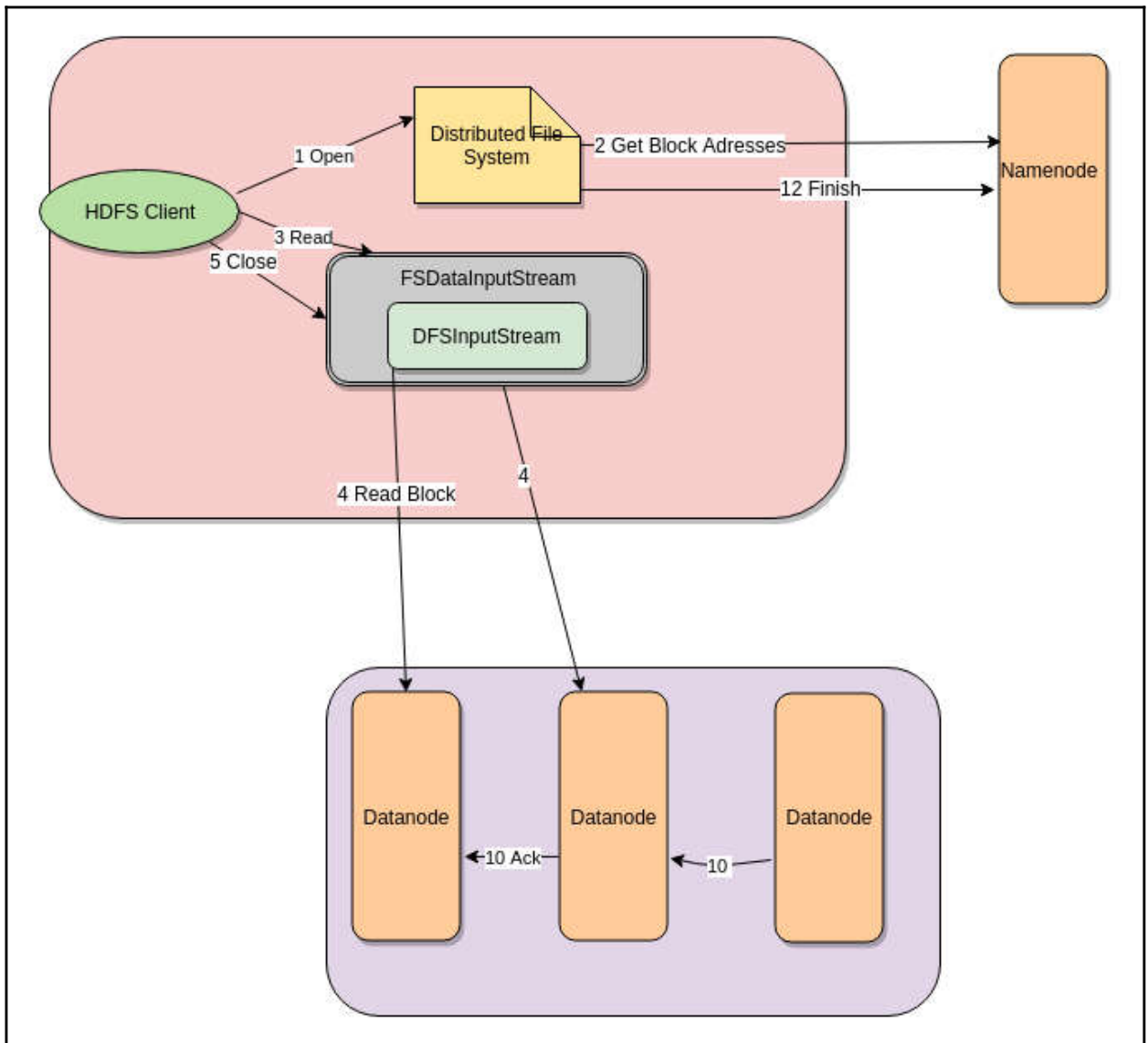


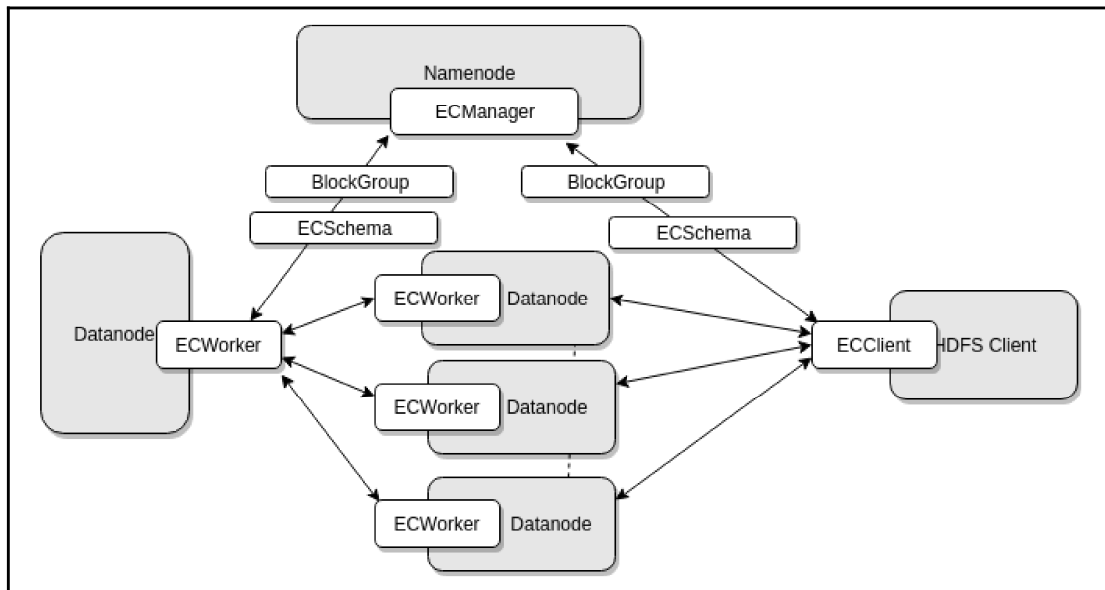
```
sshuser@hn0-apdc-s:~$ hdfs balancer --help
Usage: hdfs balancer
    [-policy <policy>]          the balancing policy: datanode or blockpool
    [-threshold <threshold>]    Percentage of disk capacity
    [-exclude [-f <hosts-file> | <comma-separated list of hosts>]] Excludes the specified datanodes.
    [-include [-f <hosts-file> | <comma-separated list of hosts>]] Includes only the specified datanodes.
    [-source [-f <hosts-file> | <comma-separated list of hosts>]] Pick only the specified datanodes as source nodes.
    [-idleiterations <idleiterations>] Number of consecutive idle iterations (-1 for Infinite) before exit.
    [-runDuringUpgrade]        Whether to run the balancer during an ongoing HDFS upgrade. This is usually not desired since it will not
                                used space on over-utilized machines.

Generic options supported are
    -conf <configuration file>    specify an application configuration file
    -D <property=value>           use value for given property
    -fs <local|namenode:port>     specify a namenode
    -jt <local|resourcemanager:port> specify a ResourceManager
    -files <comma separated list of files> specify comma separated files to be copied to the map reduce cluster
    -libjars <comma separated list of jars> specify comma separated jar files to include in the classpath.
    -archives <comma separated list of archives> specify comma separated archives to be unarchived on the compute machines.

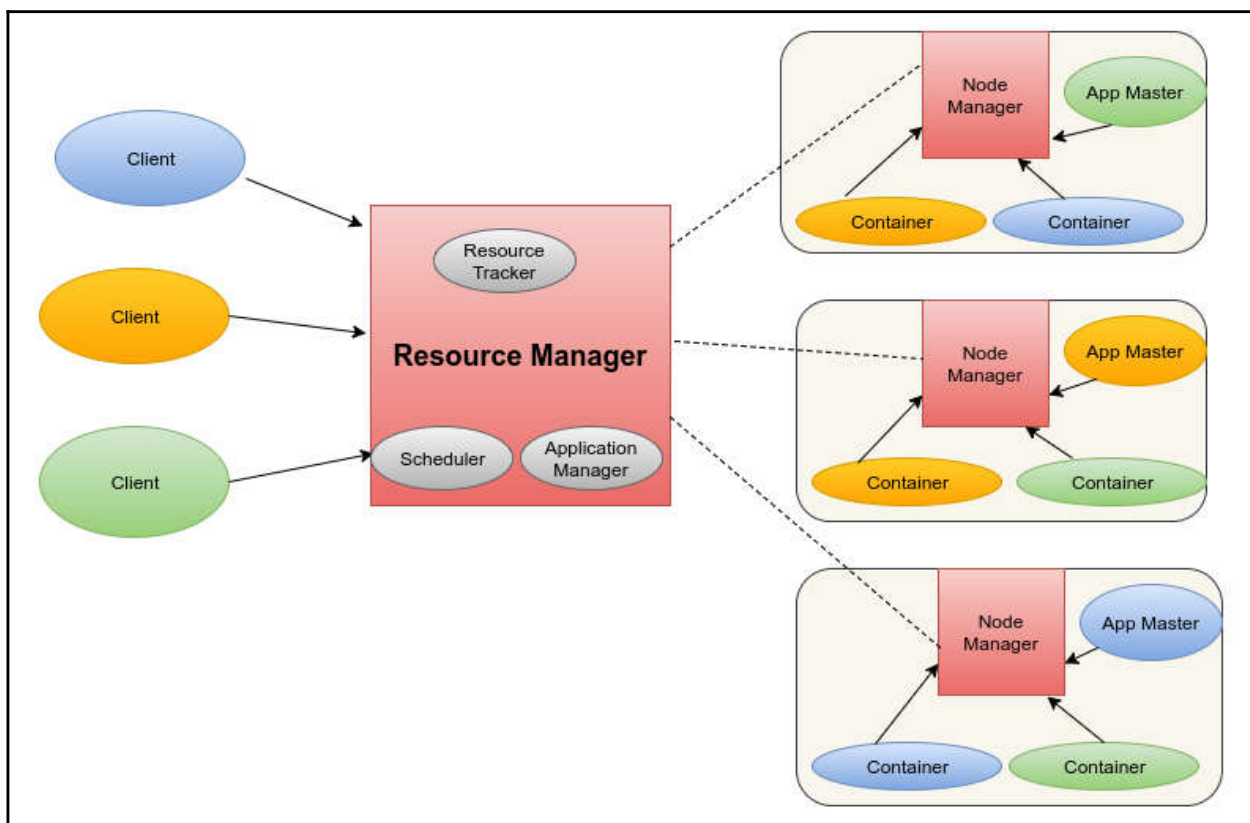
The general command line syntax is
bin/hadoop command [genericOptions] [commandOptions]
```

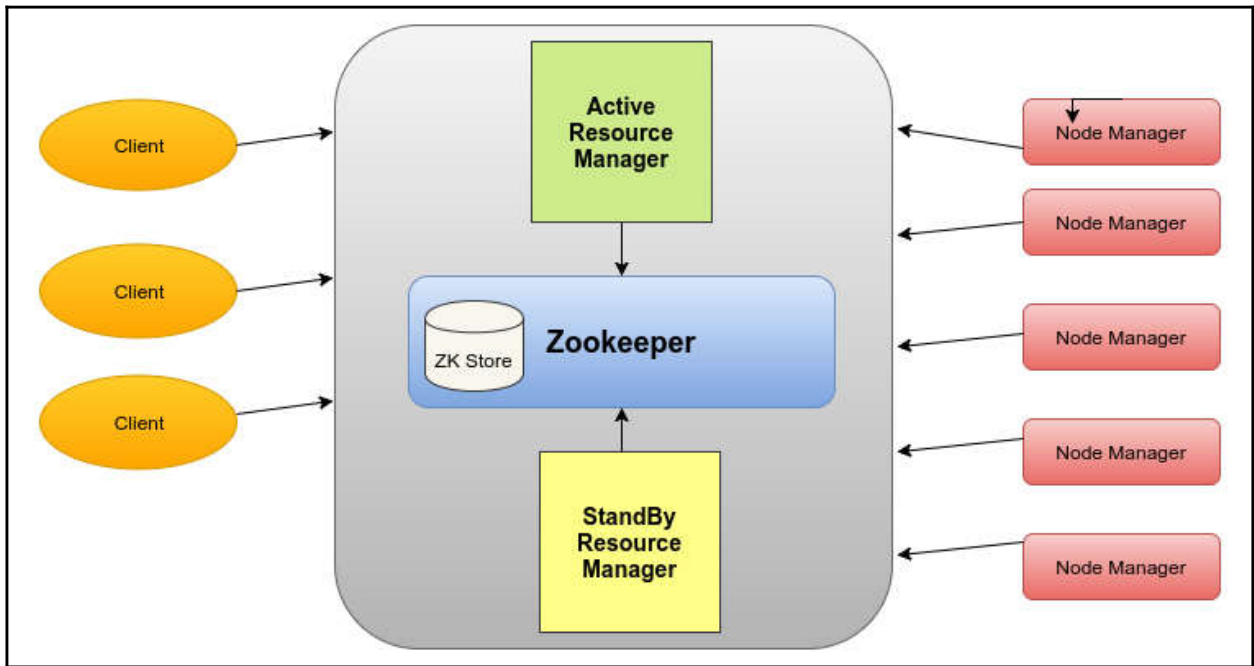
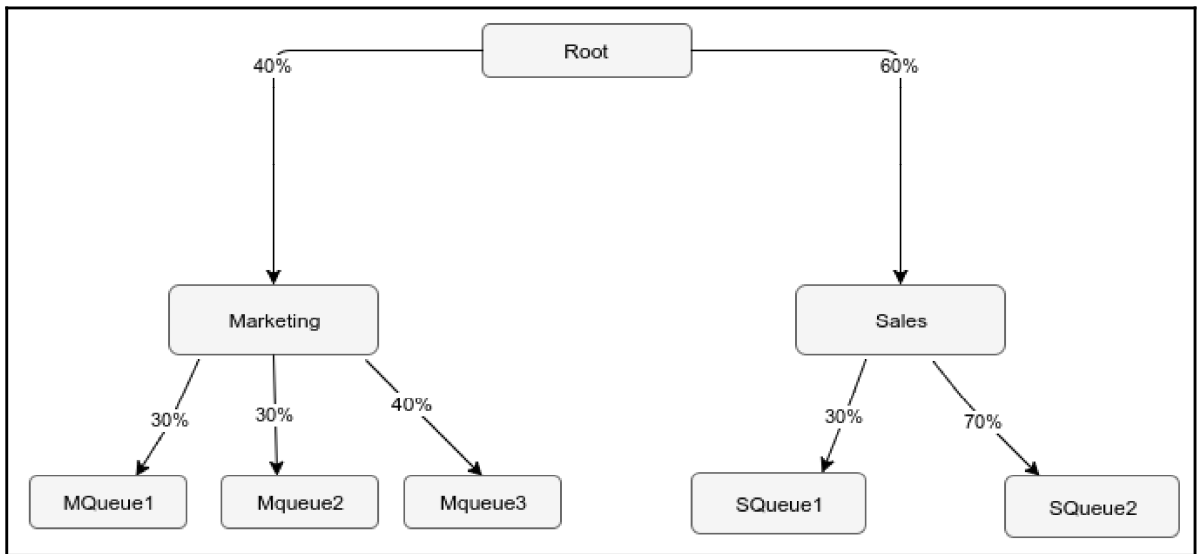






Chapter 3: YARN Resource Management in Hadoop





```

[root@ip-10-254-0-45 hadoop]# yarn application -list -appStates failed
18/01/18 08:08:03 INFO Impl.TimelineClientImpl: Timeline service address: http://ip-10-254-0-45.ap-south-1.compute.internal:8188/ws/v1/timeline/
18/01/18 08:08:03 INFO client.RMProxy: Connecting to ResourceManager at ip-10-254-0-45.ap-south-1.compute.internal/10.254.0.45:8032
Total number of applications (application-types: [] and states: [FAILED]):5
Application-Id      Application-Name      Application-Type      User      Queue      State
inal-State      Progress      Tracking-URL
application_1513582536692_0011  HIVE-dcbc8001-3339-470e-a303-b367d17fa83c      TEZ      hadoop      default
FAILED      FAILED      0% http://ip-10-254-0-45.ap-south-1.compute.internal:8088/cluster/app/application_1513582536692_0011
application_1513582536692_0009  select count(*) from tests3...service_region(Stage-1)      MAPREDUCE      hadoop      default
FAILED      FAILED      0% http://ip-10-254-0-45.ap-south-1.compute.internal:8088/cluster/app/application_1513582536692_0009
application_1513582536692_0010  HIVE-9a9368ef-6e76-4ddb-8d72-e8b0dcca534      TEZ      hadoop      default
FAILED      FAILED      0% http://ip-10-254-0-45.ap-south-1.compute.internal:8088/cluster/app/application_1513582536692_0010
application_1513582536692_0007  select sum(transaction_amou...service_region(Stage-1)      MAPREDUCE      hadoop      default
FAILED      FAILED      0% http://ip-10-254-0-45.ap-south-1.compute.internal:8088/cluster/app/application_1513582536692_0007
application_1513582536692_0008  select sum(transaction_amount) from tests3(Stage-1)      MAPREDUCE      hadoop      default
FAILED      FAILED      0% http://ip-10-254-0-45.ap-south-1.compute.internal:8088/cluster/app/application_1513582536692_0008
[root@ip-10-254-0-45 hadoop]#

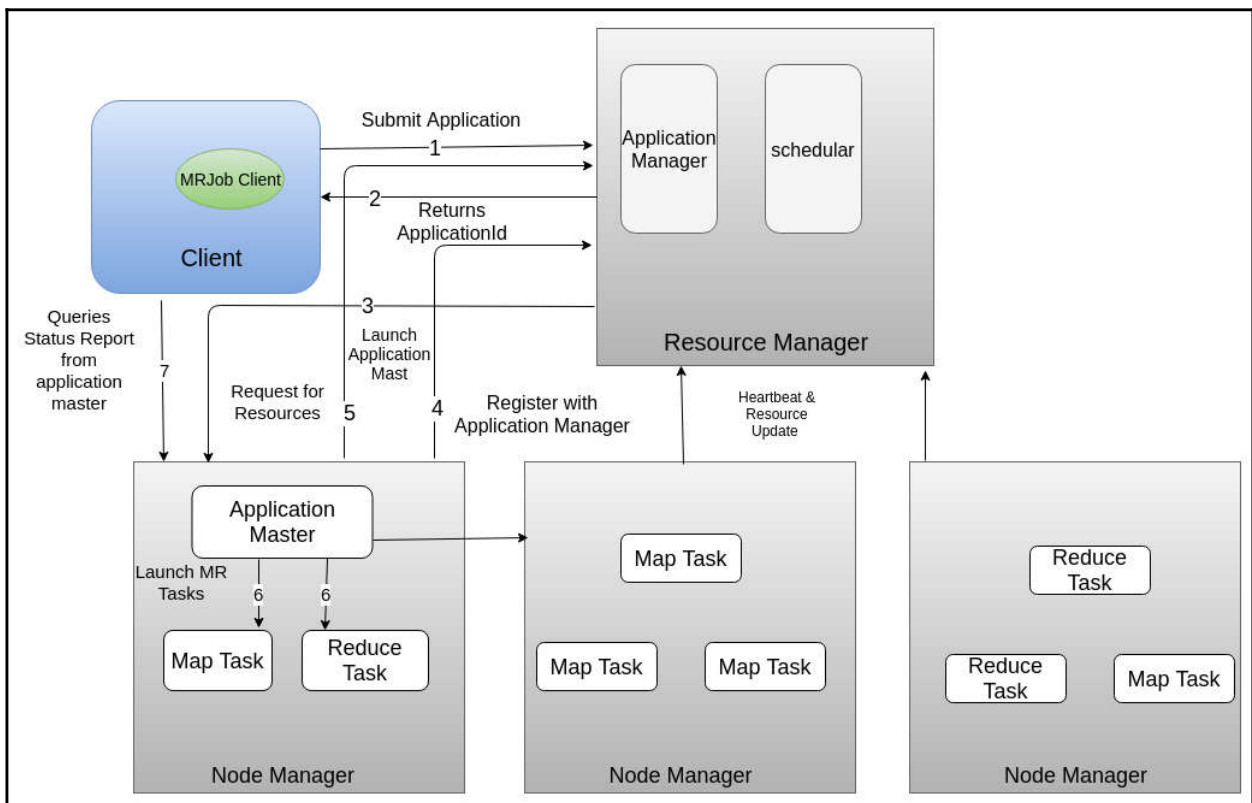
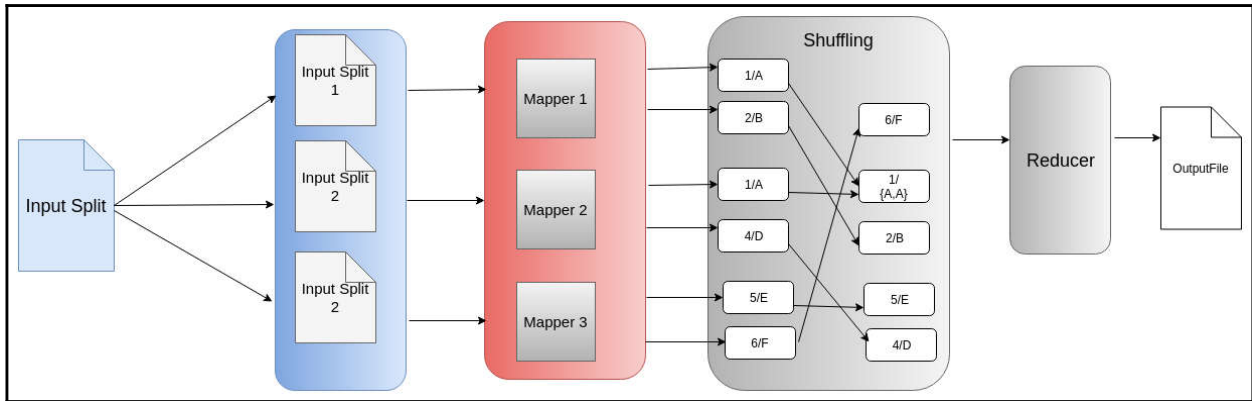
```

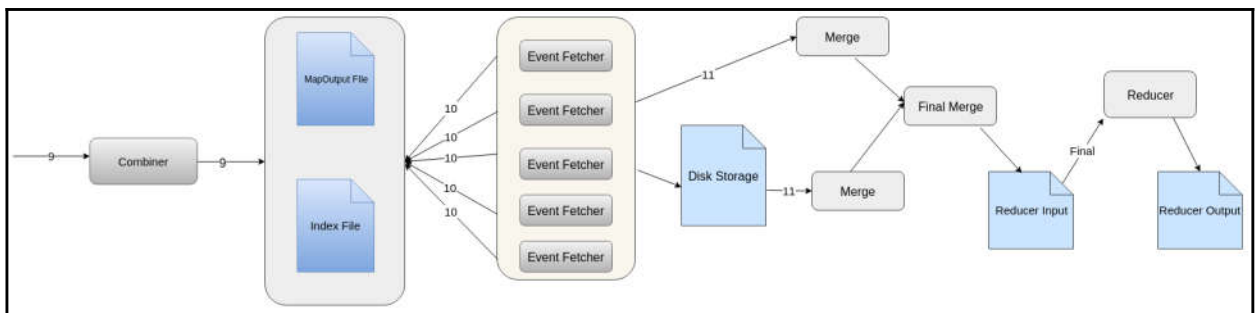
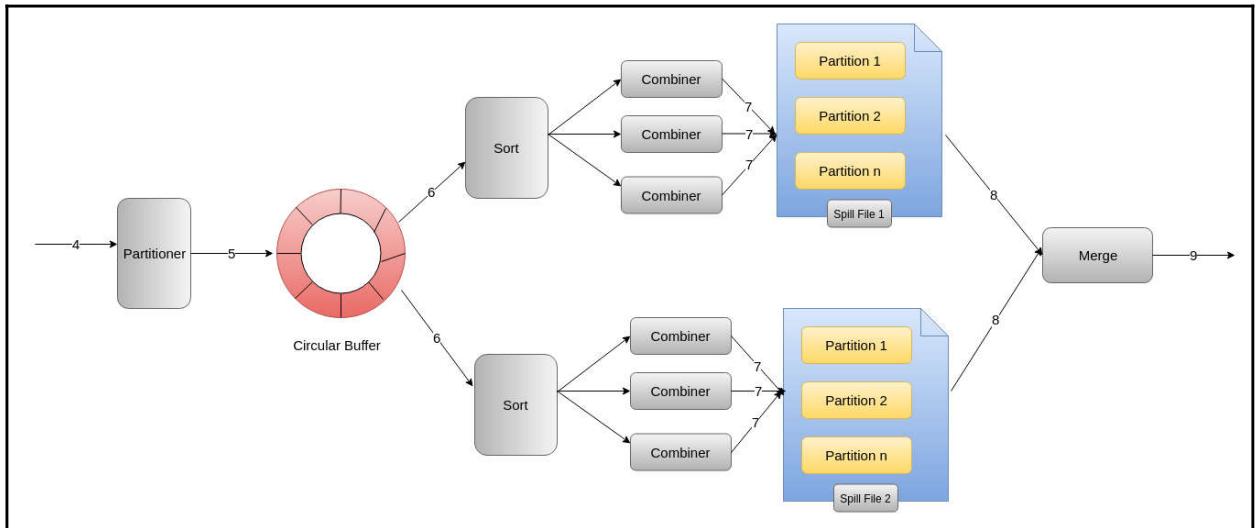
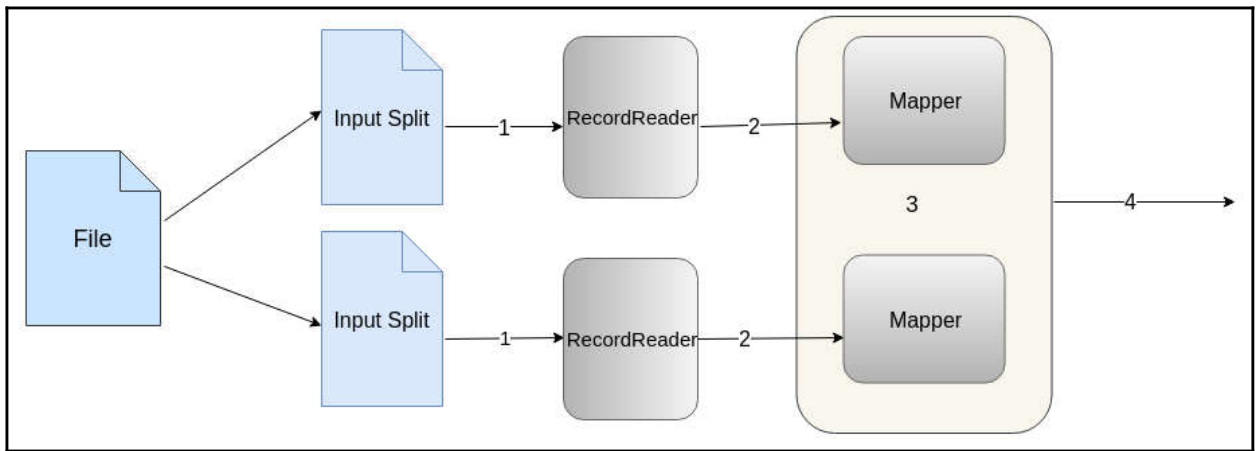
```

[root@ip-10-254-0-45 hadoop]# yarn application -status application_1513582536692_0119
18/01/18 08:20:26 INFO Impl.TimelineClientImpl: Timeline service address: http://ip-10-254-0-45.ap-south-1.compute.internal:8188/ws/v1/timeline/
18/01/18 08:20:26 INFO client.RMProxy: Connecting to ResourceManager at ip-10-254-0-45.ap-south-1.compute.internal/10.254.0.45:8032
Application Report :
  Application-Id : application_1513582536692_0119
  Application-Name : HIVE-a34d1844-cdb8-4691-a289-8a53a55f9dec
  Application-Type : TEZ
  User : root
  Queue : default
  Start-Time : 1515995486809
  Finish-Time : 1515995488561
  Progress : 100%
  State : KILLED
  Final-State : KILLED
  Tracking-URL : http://ip-10-254-0-45.ap-south-1.compute.internal:8088/cluster/app/application_1513582536692_0119
  RPC Port : -1
  AM Host : N/A
  Aggregate Resource Allocation : 1337 MB-seconds, 1 vcore-seconds
  Diagnostics : Application killed by user.
[root@ip-10-254-0-45 hadoop]#

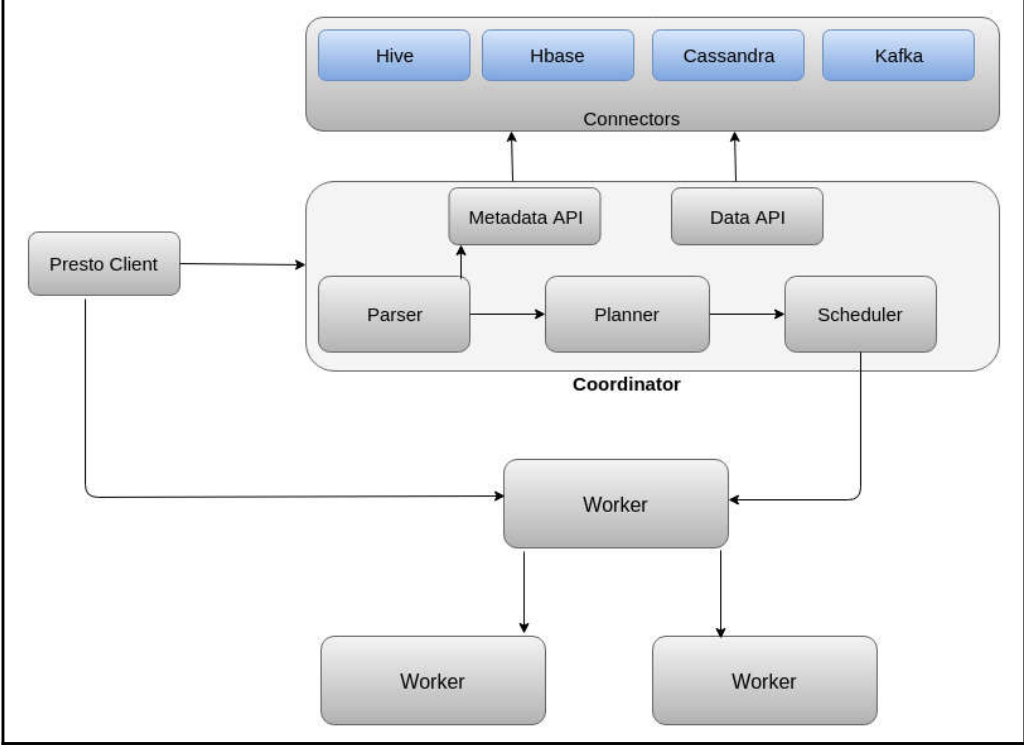
```

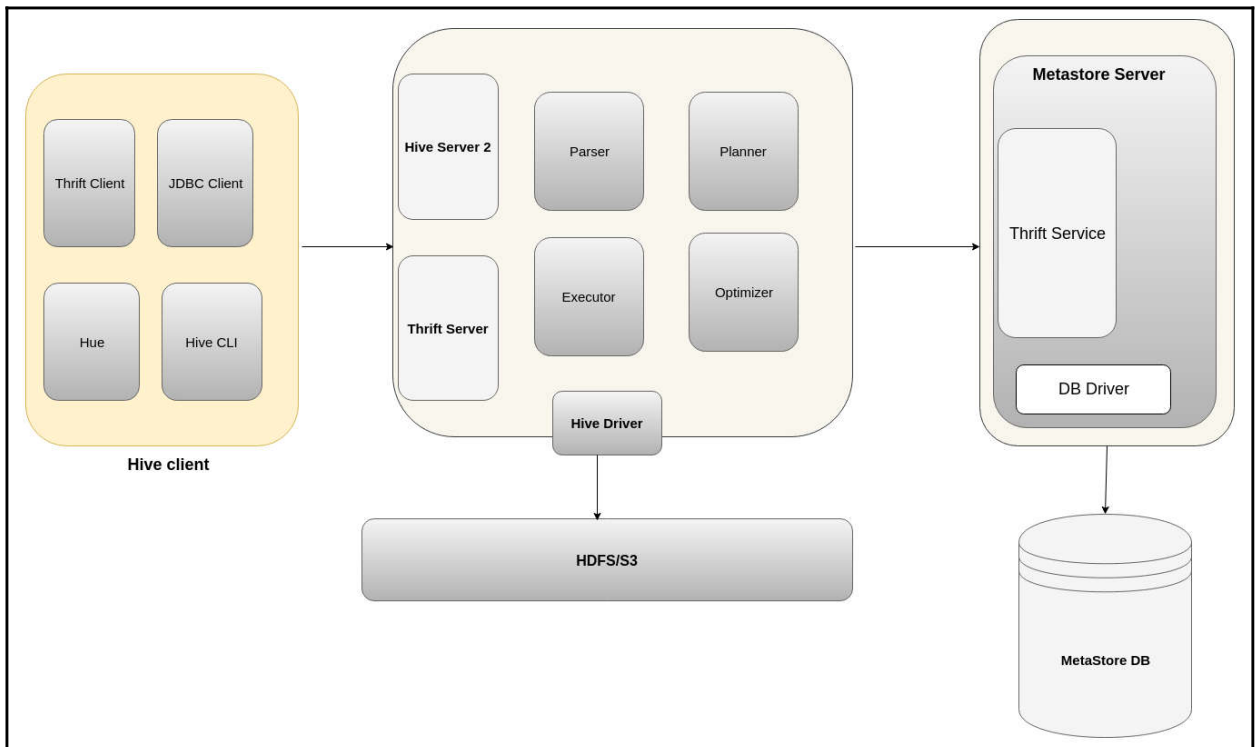
Chapter 4: Internals of MapReduce





Chapter 5: SQL on Hadoop





```
hive> describe formatted product;
OK
# col_name          data_type          comment
product_id         int
product_name       string
product_price      double
manufacturer       string

# Detailed Table Information
Database:          default
Owner:            hadoop
CreateTime:       Mon May 21 08:40:53 UTC 2018
LastAccessTime:   UNKNOWN
Retention:        0
Location:         hdfs://ip-10-254-0-45.ap-south-1.compute.internal:8020/user/hive/warehouse/product
Table Type:       EXTERNAL_TABLE
Table Parameters:
  EXTERNAL          TRUE
  numFiles          1
  numRows           1
  rawDataSize       22
  totalSize         27
  transient_lastDdlTime 1526894412

# Storage Information
SerDe Library:    org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
InputFormat:      org.apache.hadoop.mapred.TextInputFormat
OutputFormat:     org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
Compressed:       No
Num Buckets:      -1
Bucket Columns:   []
Sort Columns:     []
Storage Desc Params:
  field.delim      \t
  line.delim       \n
  serialization.format \t
```

```

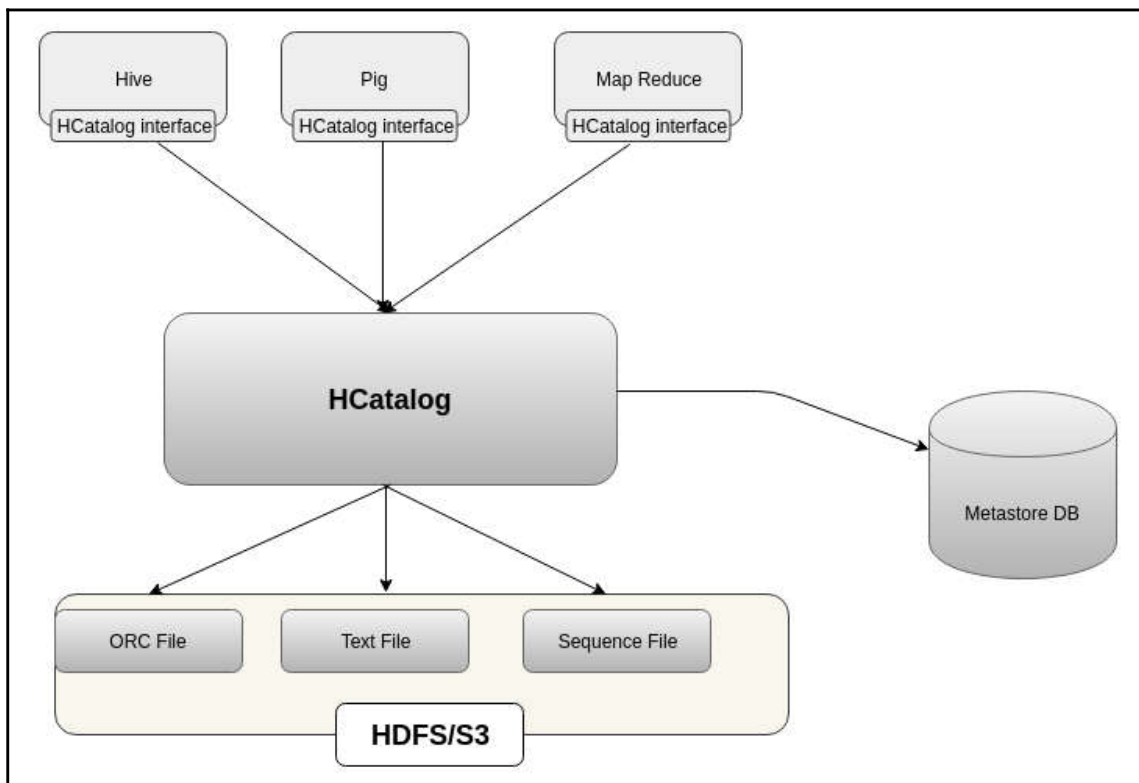
hive> explain Select count(distinct product_id) , manufacturer from product group by manufacturer;
OK
Plan optimized by CBO.

Vertex dependency in root stage
Reducer 2 <- Map 1 (SIMPLE_EDGE)

Stage-0
  Fetch Operator
    limit:-1
  Stage-1
    Reducer 2
      File Output Operator [FS_11]
      Select Operator [SEL_10] (rows=1 width=22)
        Output:["_col0", "_col1"]
      Group By Operator [GBY_9] (rows=1 width=22)
        Output:["_col0", "_col1"],aggregations:["count(_col0)"],keys:_col1
      Select Operator [SEL_5] (rows=1 width=22)
        Output:["_col0", "_col1"]
      Group By Operator [GBY_4] (rows=1 width=22)
        Output:["_col0", "_col1"],keys:KEY._col0, KEY._col1
    <-Map 1 [SIMPLE_EDGE]
      SHUFFLE [RS_3]
        PartitionCols:_col0
        Group By Operator [GBY_2] (rows=1 width=22)
          Output:["_col0", "_col1"],keys:manufacturer, product_id
        Select Operator [SEL_1] (rows=1 width=22)
          Output:["manufacturer", "product_id"]
        TableScan [TS_0] (rows=1 width=22)
          default@product,product,Tbl:COMPLETE,Col:NONE,Output:["product_id", "manufacturer"]

Time taken: 0.108 seconds, Fetched: 29 row(s)

```



```

ravishankarnair — bash — 80x7
ravion:~ ravishankarnair$ hive --version
Hive 2.3.2
Git git://stakiar-MBP.local/Users/stakiar/Desktop/scratch-space/apache-hive -r 8
57a9fd8ad725a53bd95c1b2d6612f9b1155f44d
Compiled by stakiar on Thu Nov 9 09:11:39 PST 2017
From source with checksum dc38920061a4eb32c4d15ebd5429ac8a
ravion:~ ravishankarnair$

```

```

ravishankarnair — java -Xmx256m -Djava.library.path=/Users/ravishankarnair/cdh545/hadoop/lib/native -Djava.net.preferIPv4Stack=true -Dhadoop.log.dir=/Users/ravishankar...
hive> set hive.support.concurrency = true;
hive> set hive.enforce.bucketing = true;
hive> set hive.exec.dynamic.partition.mode = nonstrict;
hive> set hive.txn.manager = org.apache.hadoop.hive.ql.lockmgr.DbTxnManager;
hive> set hive.compactor.initiator.on = true;
hive> set hive.compactor.worker.threads = 1;
hive>

```

```
ravishankarnair — java -Xmx256m -Djava.library.path=/Users/ravishankarnair/cdh545/hadoop/lib/native -Djava.net.preferIPv4Stack=true -Dhadoop.log.dir=/Users/ravishankar...  
hive> set hive.enforce.bucketing;  
hive.enforce.bucketing=true  
hive>
```

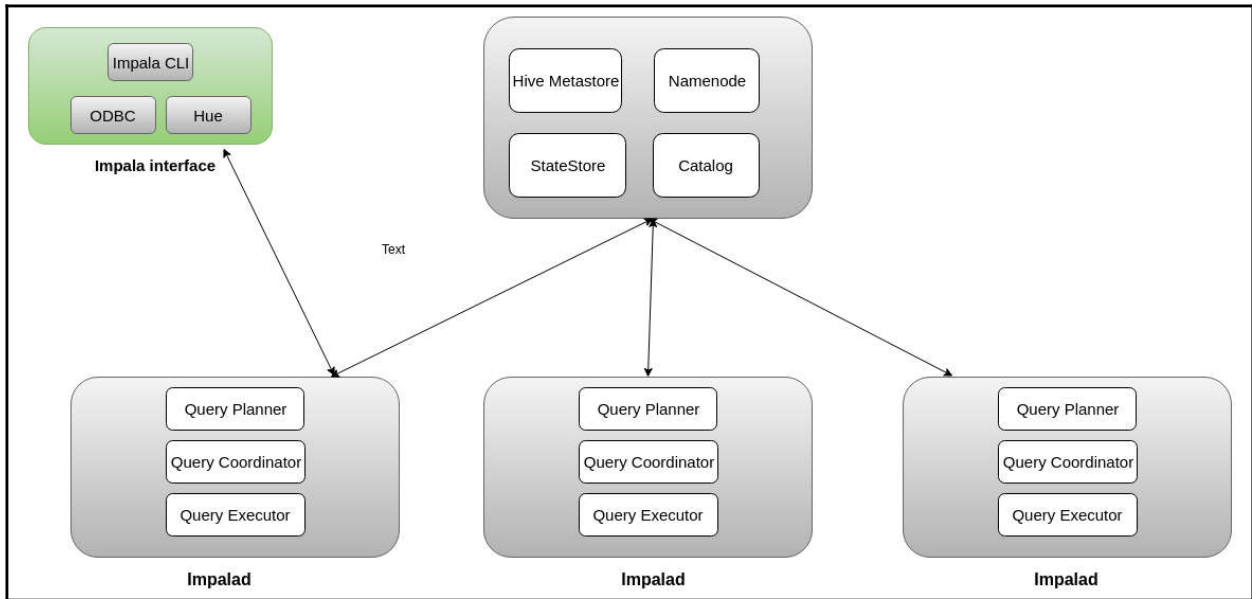
```
ravishankarnair — java -Xmx256m -Djava.library.path=/Users/ravishankarnair/cdh545/hadoop/lib/native -Djava.net.preferIPv4Stack=true -Dhadoop.log.dir=/Users/ravishankar...  
hive> create table tschools(school_id int, school_name string, school_loc string  
) clustered by (school_id) into 5 buckets stored as ORC TBLPROPERTIES ('transactional' = 'true');  
OK  
Time taken: 0.231 seconds  
hive>
```

```
ravishankarnair — java -Xmx256m -Djava.library.path=/Users/ravishankarnair/cdh545/hadoop/lib/native -Djava.net.preferIPv4Stack=true -Dhadoop.log.dir=/Users/ravishankar...  
hive> select * from tschools;  
OK  
5      fgh      fhg  
5      fgh      fhg  
1      abc      acb  
1      abc      acb  
2      bcd      bdc  
2      bcd      bdc  
3      cde      ced  
3      cde      ced  
4      efg      egf  
4      efg      egf  
Time taken: 0.298 seconds, Fetched: 10 row(s)  
hive>
```

```
ravishankarnair — java -Xmx256m -Djava.library.path=/Users/ravishankarnair/cdh545/hadoop/lib/native -Djava.net.preferIPv4Stack=true -Dhadoop.log.dir=/Users/ravishankar...  
hive> update tschools set school_id = 10 where school_id = 5;  
FAILED: SemanticException [Error 10302]: Updating values of bucketing columns is  
not supported. Column school_id.  
hive>
```

```
ravishankarnair — java -Xmx256m -Djava.library.path=/Users/ravishankarnair/cdh545/hadoop/lib/native -Djava.net.preferIPv4Stack=true -Dhadoop.log.dir=/Users/ravishankar...
hive> select * from tschools;
OK
5      MIT      fhg
5      MIT      fhg
1      abc      acb
1      abc      acb
2      bcd      bdc
2      bcd      bdc
3      cde      ced
3      cde      ced
4      efg      egf
4      efg      egf
Time taken: 0.171 seconds, Fetched: 10 row(s)
hive> 
```

```
ravishankarnair — java -Xmx256m -Djava.library.path=/Users/ravishankarnair/cdh545/hadoop/lib/native -Djava.net.preferIPv4Stack=true -Dhadoop.log.dir=/Users/ravishankar...
hive> select * from tschools;
OK
1      abc      acb
1      abc      acb
2      bcd      bdc
2      bcd      bdc
3      cde      ced
3      cde      ced
4      efg      egf
4      efg      egf
Time taken: 0.108 seconds, Fetched: 8 row(s)
hive> 
```

```

cloudera@quickstart:~$ impala-shell -i localhost --quiet
Starting Impala Shell without Kerberos authentication
*****
***
Welcome to the Impala shell. Copyright (c) 2015 Cloudera, Inc. All rights reserved.
(Impala Shell v2.6.0-cdh5.8.0 (5464d17) built on Thu Jun 16 12:35:00 PDT 2016)

To see a summary of a query's progress that updates in real-time, run 'set
LIVE_PROGRESS=1;'.
*****
***
[localhost:21000] >
  
```

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
[localhost:21000] > version();  
Shell version: Impala Shell v2.6.0-cdh5.8.0 (5464d17) built on Thu Jun 16 12:35:  
00 PDT 2016  
Server version: impalad version 2.6.0-cdh5.8.0 RELEASE (build 5464d1750381b40a7e  
7163b12b09f11b891b4de3)  
[localhost:21000] > show databases;  
+-----+  
| name          | comment                                     |  
+-----+  
| _impala_builtins | System database for Impala builtin functions |  
| default        | Default Hive database                     |  
| passion        |                                             |  
+-----+  
[localhost:21000] > █
```

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
[localhost:21000] > select current_database();  
+-----+  
| current_database() |  
+-----+  
| default            |  
+-----+  
[localhost:21000] > show tables in default;  
+-----+  
| name          |  
+-----+  
| customers     |  
| departments   |  
| employee      |  
| orders        |  
| orders1       |  
| schools       |  
+-----+  
[localhost:21000] > █
```

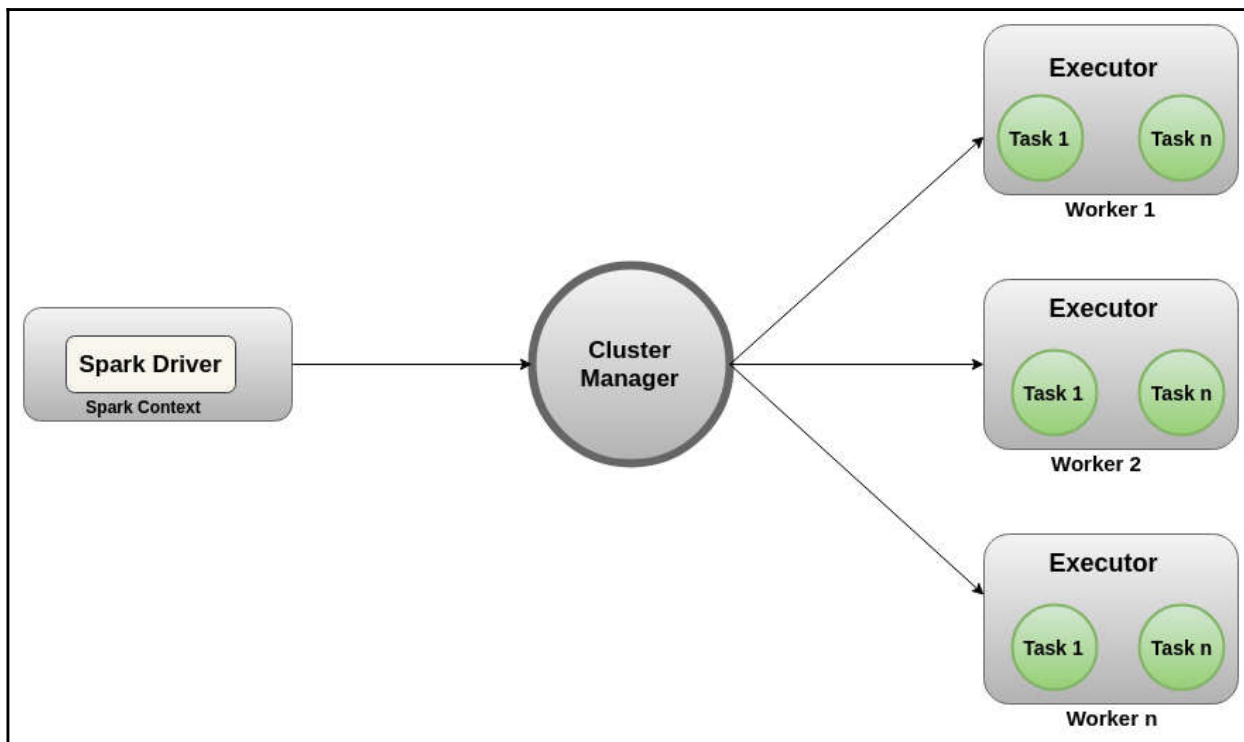
```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
[localhost:21000] > select * from departments;  
+-----+-----+  
| department_id | department_name |  
+-----+-----+  
| 4              | Apparel         |  
| 5              | Golf            |  
| 2              | Fitness         |  
| 3              | Footwear        |  
| 6              | Outdoors        |  
| 7              | Fan Shop        |  
+-----+-----+  
[localhost:21000] > █
```

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
[localhost:21000] > describe departments;  
+-----+-----+-----+  
| name          | type  | comment |  
+-----+-----+-----+  
| department_id | int   |         |  
| department_name | string |         |  
+-----+-----+-----+  
[localhost:21000] > █
```

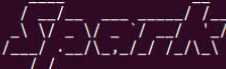
```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
[cloudera@quickstart ~]$ hadoop fs -mkdir -p passion/sample1  
[cloudera@quickstart ~]$ hadoop fs -mkdir -p passion/sample2  
[cloudera@quickstart ~]$ hadoop fs -put data1.csv passion/sample1  
[cloudera@quickstart ~]$ hadoop fs -put data2.csv passion/sample2  
[cloudera@quickstart ~]$  
[cloudera@quickstart ~]$  
[cloudera@quickstart ~]$
```

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
[localhost:21000] > select * from table1;  
+-----+-----+-----+-----+  
| id | col_1 | col_2 | col_3 |  
+-----+-----+-----+-----+  
| 1 | true | 123.123 | 2012-10-24 08:55:00 |  
| 2 | false | 1243.5 | 2012-10-25 13:40:00 |  
| 3 | false | 24453.325 | 2008-08-22 09:33:21.123000000 |  
| 4 | false | 243423.325 | 2007-05-12 22:32:21.334540000 |  
| 5 | true | 243.325 | 1953-04-22 09:11:33 |  
+-----+-----+-----+-----+  
[localhost:21000] > select * from table2;  
+-----+-----+-----+  
| id | col_1 | col_2 |  
+-----+-----+-----+  
| 1 | true | 12789.123 |  
| 2 | false | 1243.5 |  
| 3 | false | 24453.325 |  
| 4 | false | 2423.3254 |  
| 5 | true | 243.325 |  
| 60 | false | 243565423.325 |  
| 70 | true | 243.325 |  
| 80 | false | 243423.325 |
```

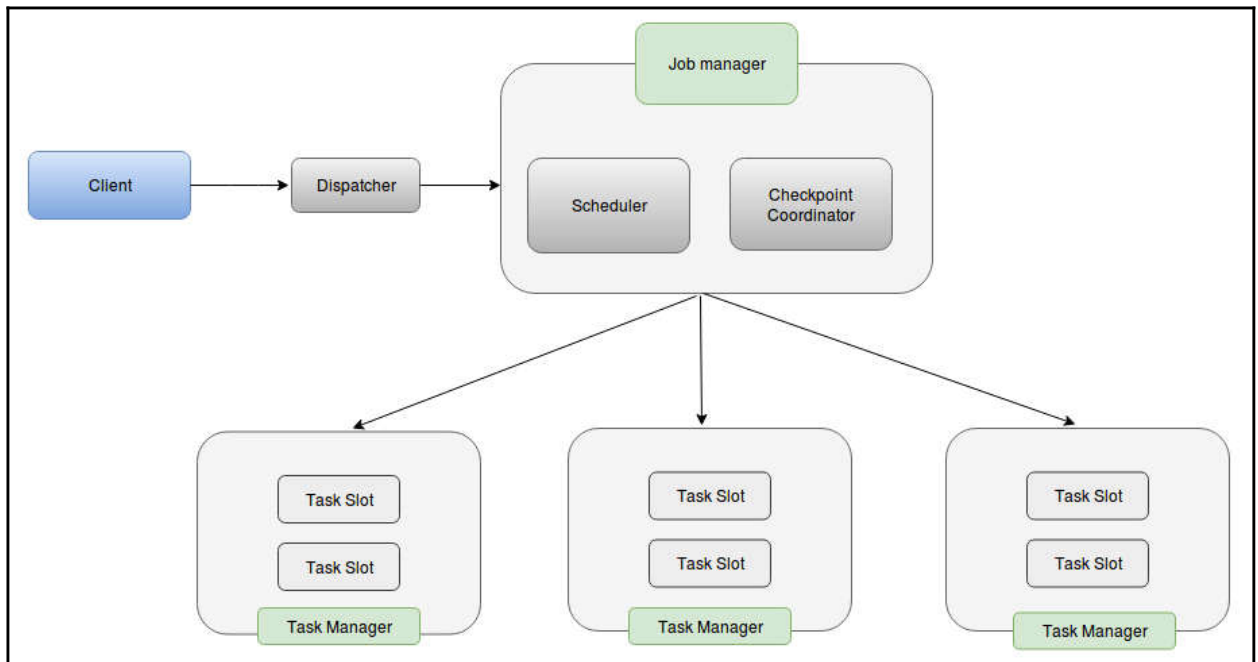
Chapter 6: Real-Time Processing Engines

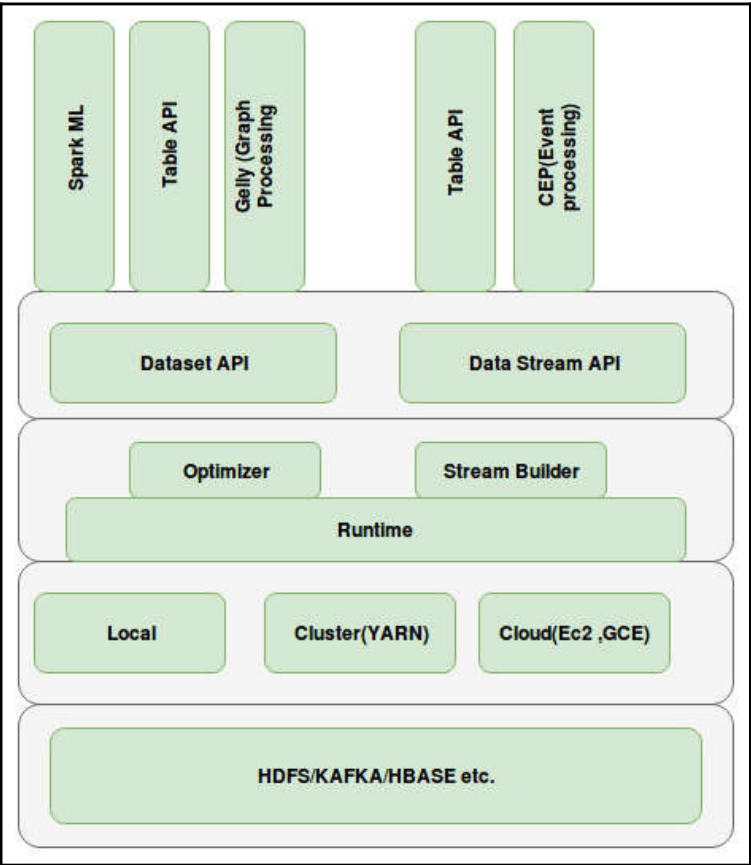


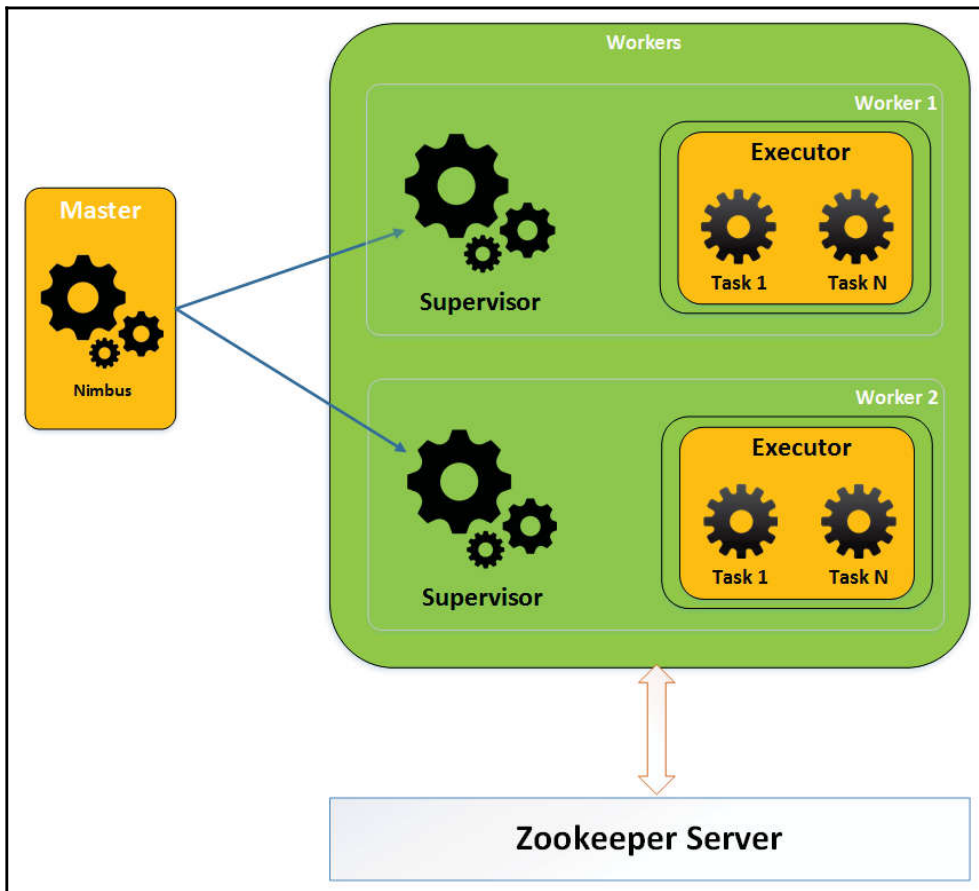
```
exa00077@exa00071:~/Softwares/spark/bin$ ./spark-shell
2018-09-15 20:30:19 WARN Utils:66 - Your hostname, exa00071 resolves to a loopback address: 127.0.1.1; using 192.168.0.103 instead
ce wlp5s0)
2018-09-15 20:30:19 WARN Utils:66 - Set SPARK_LOCAL_IP if you need to bind to another address
2018-09-15 20:30:20 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using builtin-java classes
licable
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
Spark context Web UI available at http://192.168.0.103:4040
Spark context available as 'sc' (master = local[*], app id = local-1537023626092).
Spark session available as 'spark'.
Welcome to

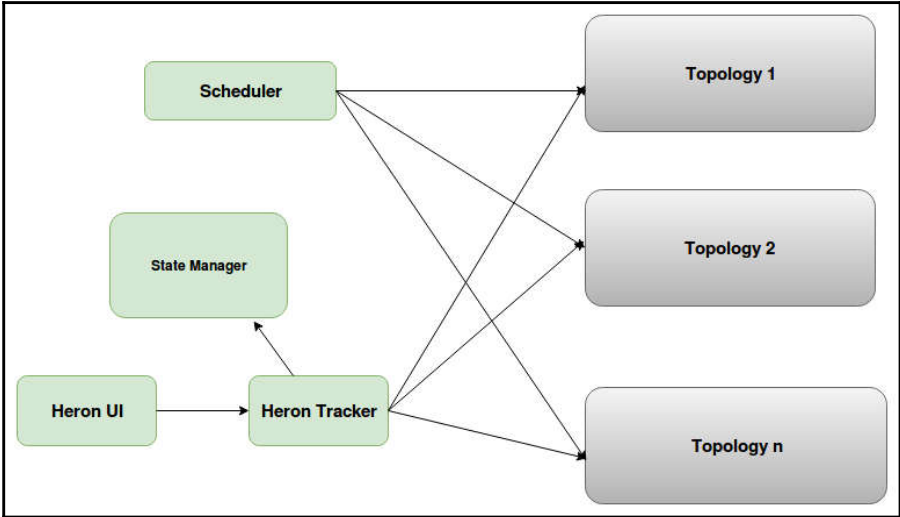
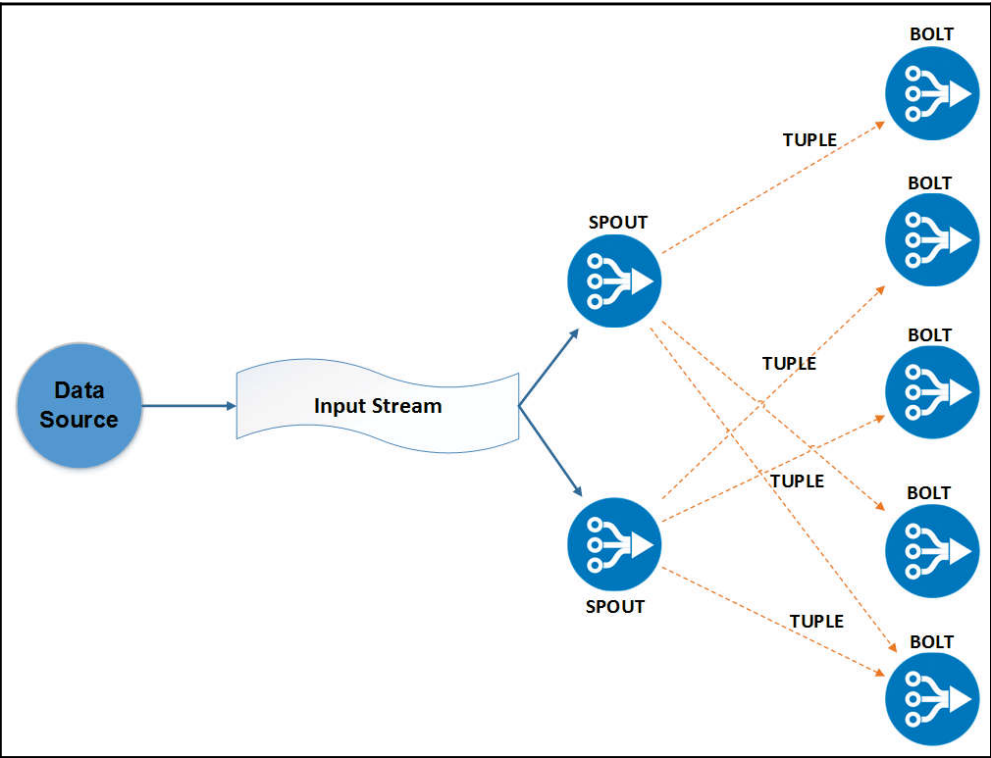
 version 2.3.1

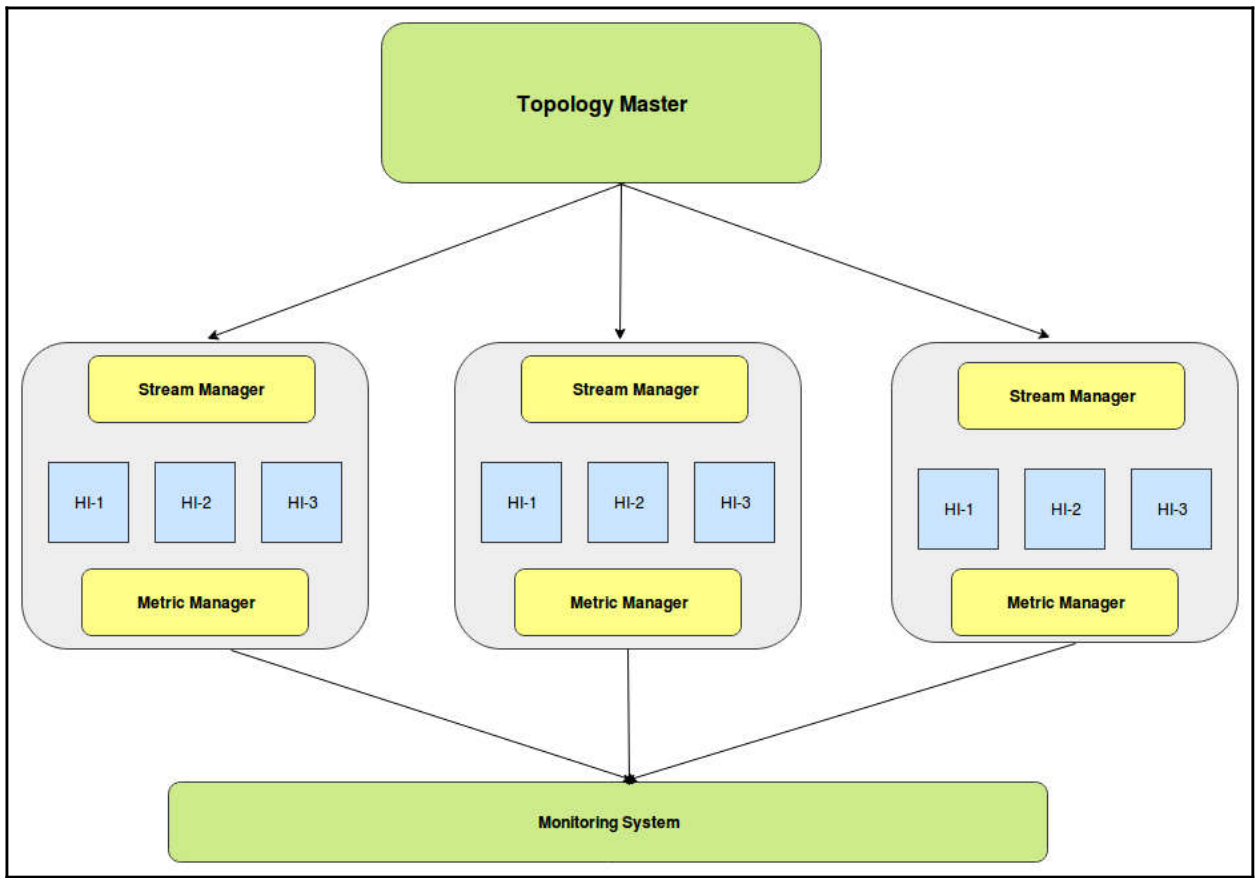
Using Scala version 2.11.8 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_181)
Type in expressions to have them evaluated.
Type :help for more information.
```



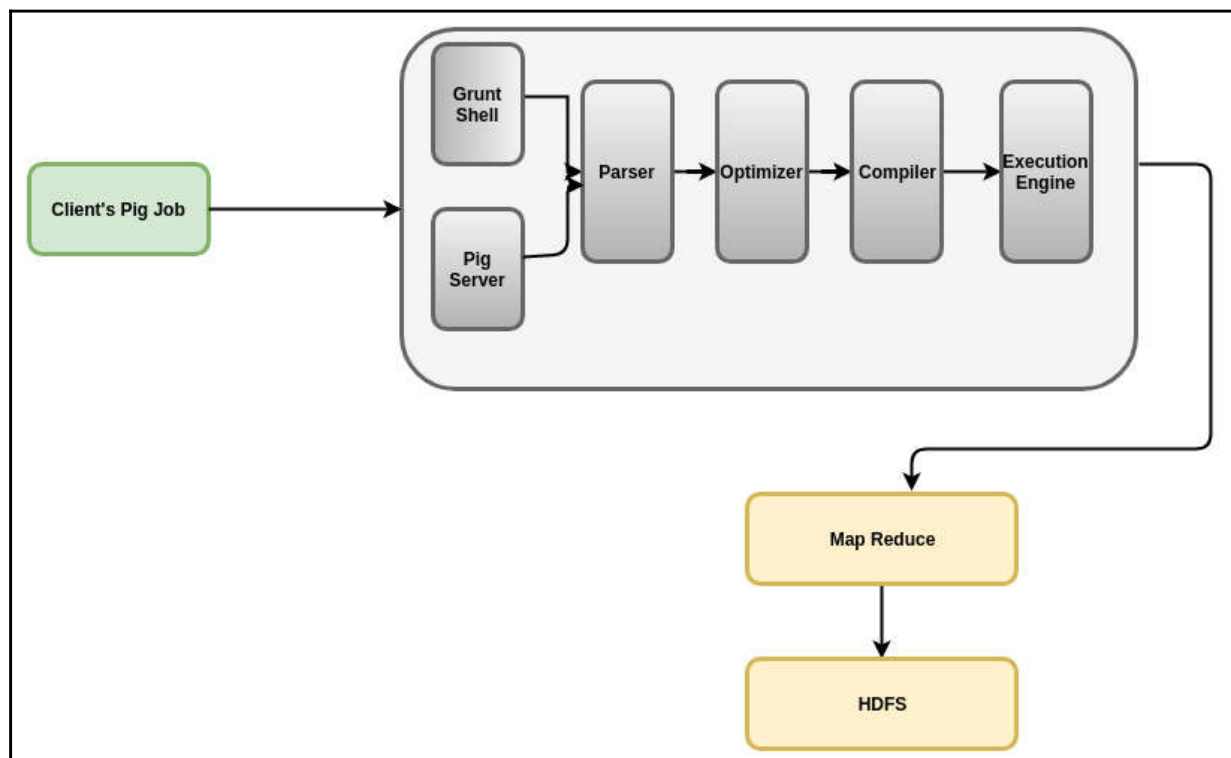


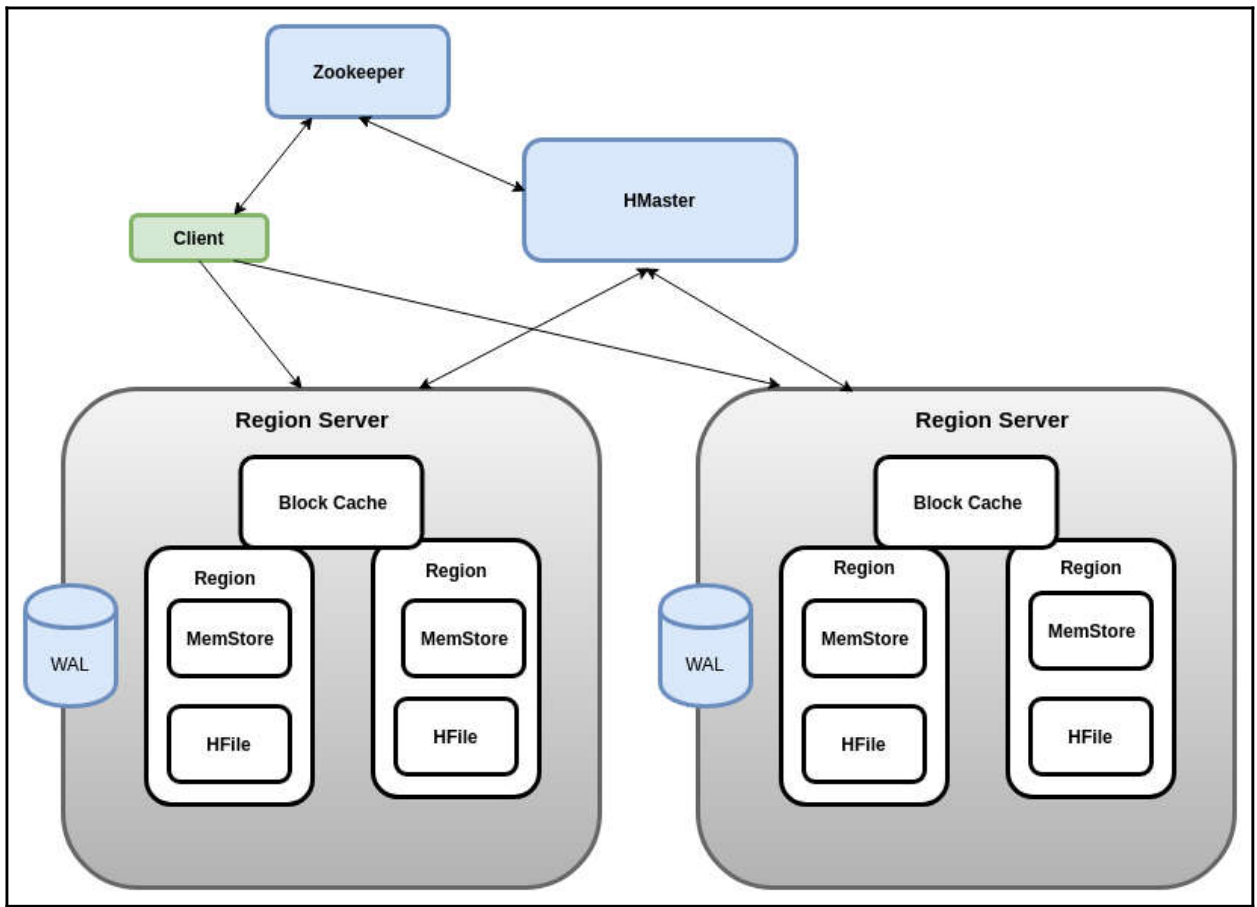


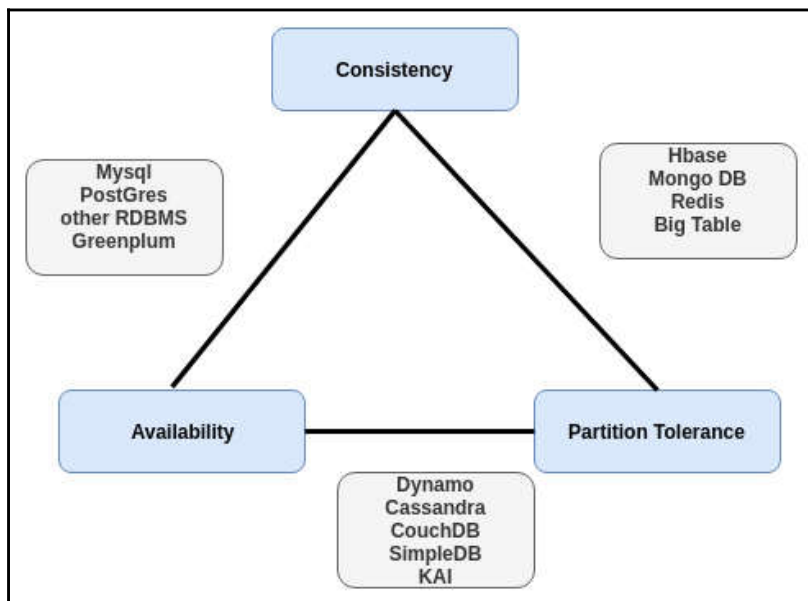




Chapter 7: Widely Used Hadoop Ecosystem Components







```

# Set environment variables here.

# This script sets variables multiple times over the course of starting an hbase process,
# so try to keep things idempotent unless you want to take an even deeper look
# into the startup scripts (bin/hbase, etc.)

# The java implementation to use. Java 1.8+ required.
export JAVA_HOME=/opt/jjava/jdk1.8.0/

# Extra Java CLASSPATH elements. Optional.
# export HBASE_CLASSPATH=

# The maximum amount of heap to use. Default is left to JVM default.
# export HBASE_HEAPSIZE=1G

# Uncomment below if you intend to use off heap cache. For example, to allocate 8G of
# offheap, set the value to "8G".
# export HBASE_OFFHEAPSIZE=1G

# Extra Java runtime options.
# Below are what we set by default. May only work with SUN JVM.
# For more on why as well as other possible settings,
# see http://hbase.apache.org/book.html#performance
export HBASE_OPTS="$HBASE_OPTS -XX:+UseConcMarkSweepGC"

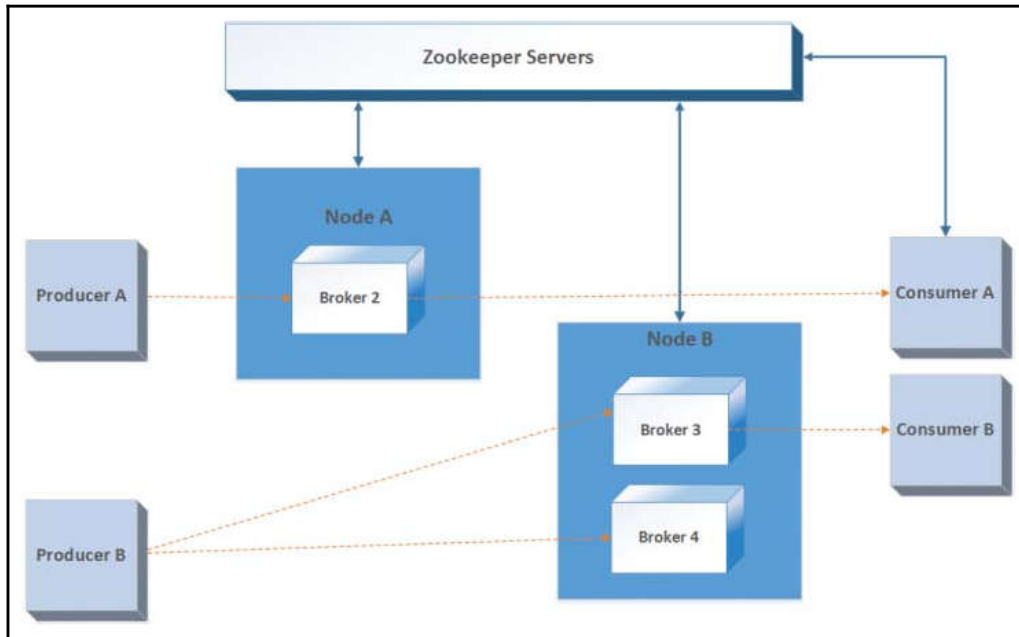
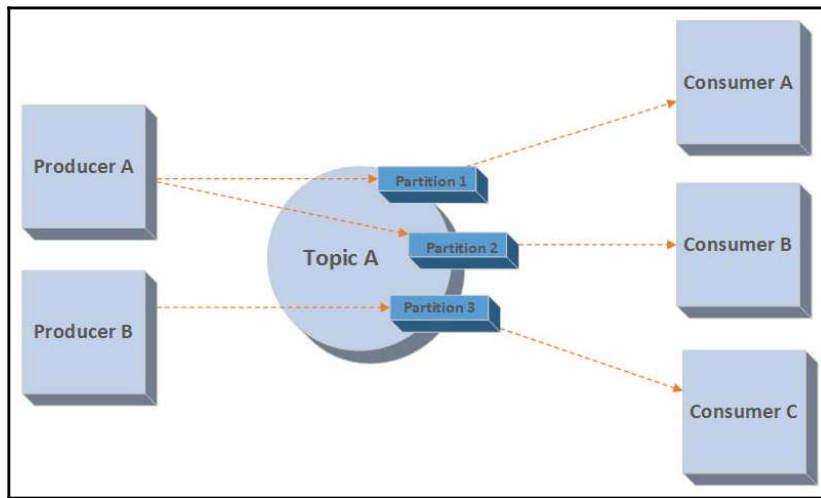
# Uncomment one of the below three options to enable java garbage collection logging for the server-side processes.

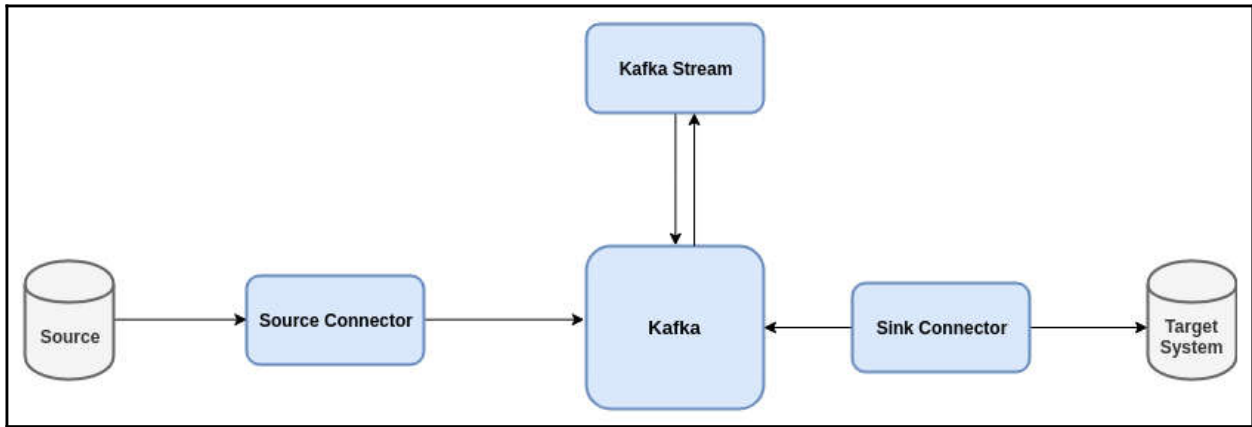
# This enables basic gc logging to the .out file.
# export SERVER_GC_OPTS="-verbose:gc -XX:+PrintGCDetails -XX:+PrintGCDateStamps"

# This enables basic gc logging to its own file.
# If FILE-PATH is not replaced, the log file(.gc) would still be generated in the HBASE_LOG_DIR .
# export SERVER_GC_OPTS="-verbose:gc -XX:+PrintGCDetails -XX:+PrintGCDateStamps -Xloggc:<FILE-PATH>"

# This enables basic GC logging to its own file with automatic log rolling. Only applies to jdk 1.6.0_34+ and 1.7.0_2+.
# If FILE-PATH is not replaced, the log file(.gc) would still be generated in the HBASE_LOG_DIR .
# export SERVER_GC_OPTS="-verbose:gc -XX:+PrintGCDetails -XX:+PrintGCDateStamps -Xloggc:<FILE-PATH> -XX:+UseGCLogFileRotation -XX:NumberOf

```



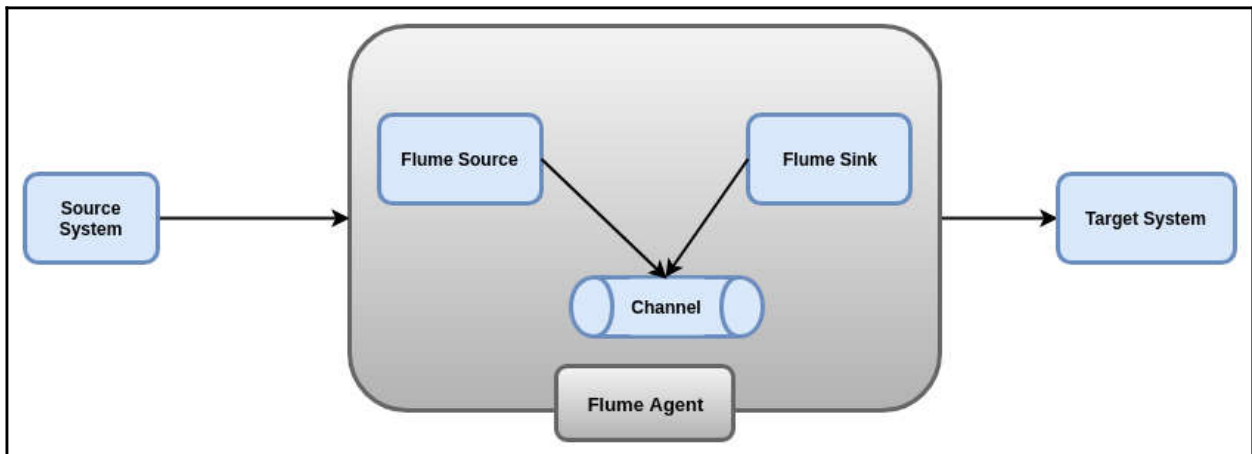


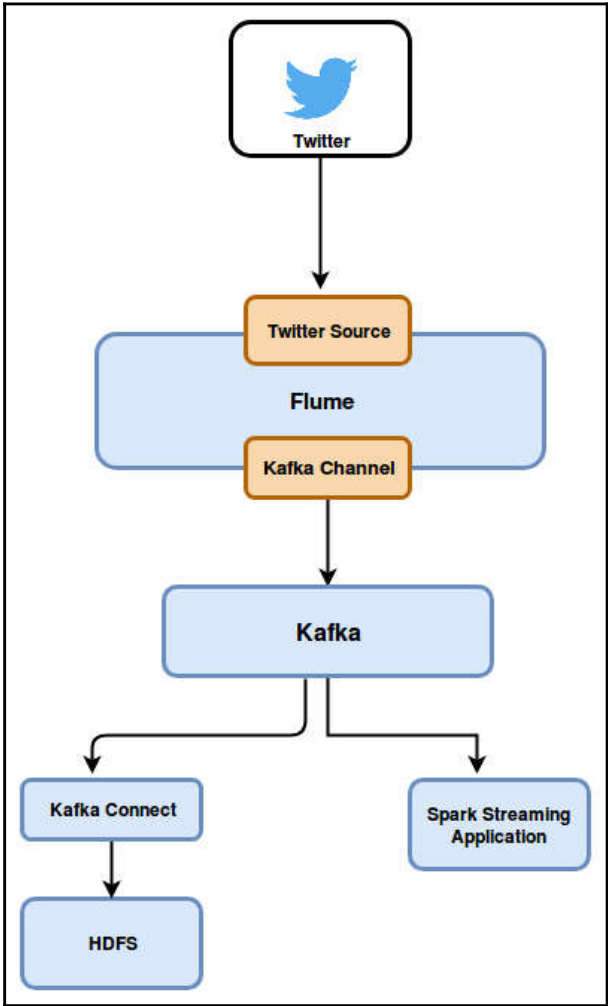
```

SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/chanchal/projects/confluent-3.2.2/share/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/chanchal/projects/confluent-3.2.2/share/aticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
{"id":1,"name":{"string":"Manish"}}
{"id":2,"name":{"string":"Chanchal"}}
  
```

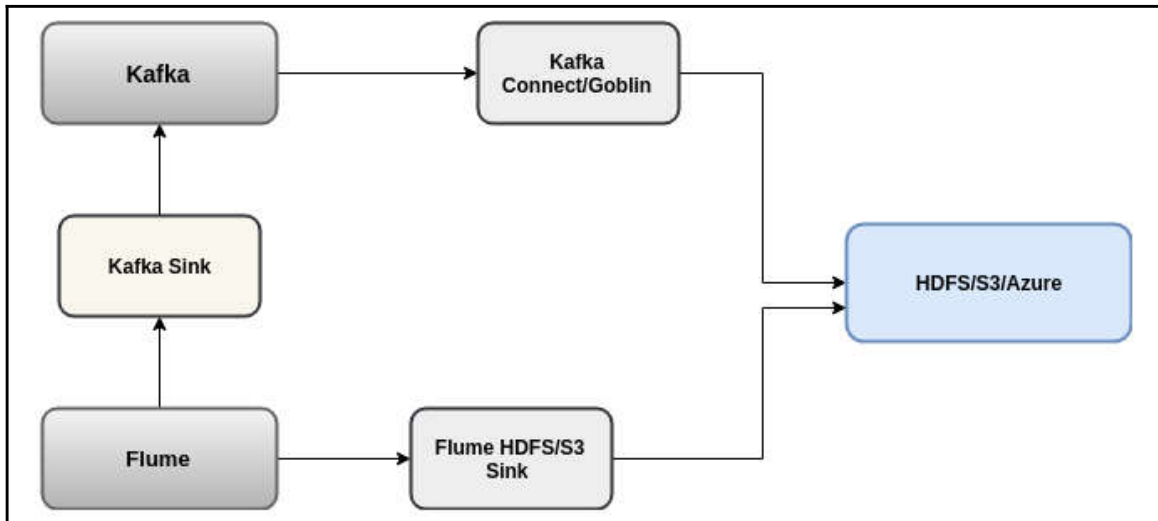
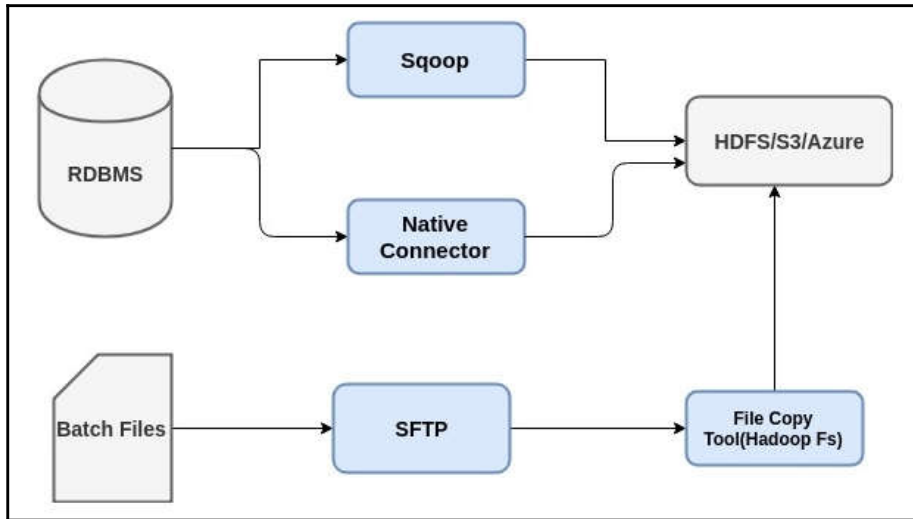
```

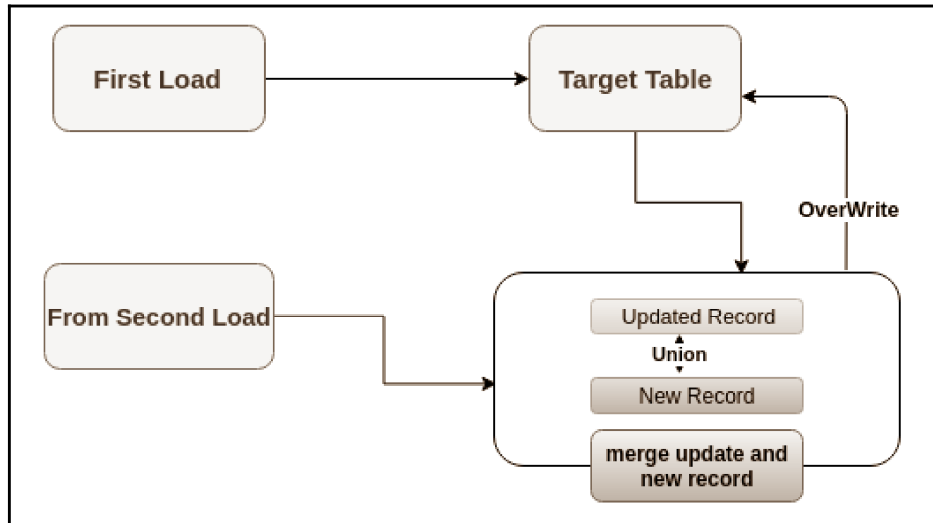
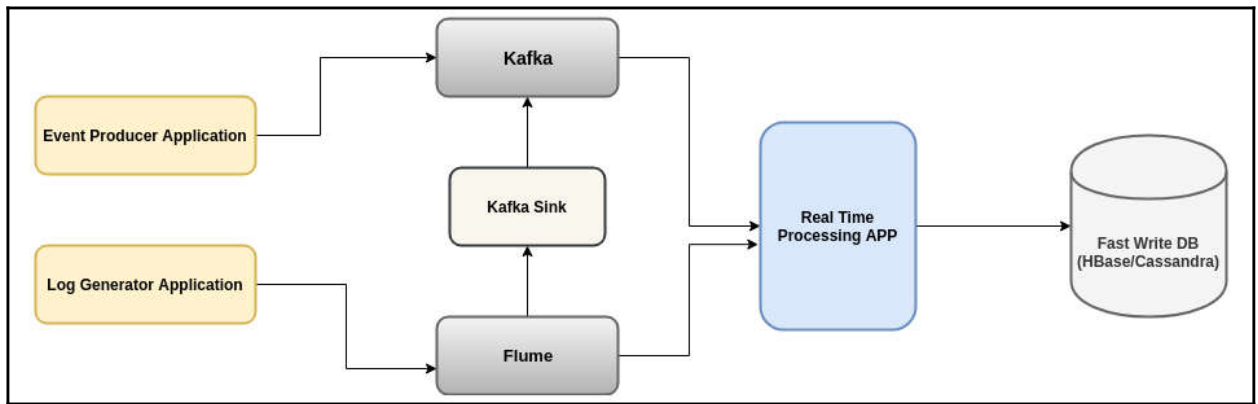
sqlite> select * from authors_sink;
Chanchal|60.0|1|26
Manish|80.0|2|32
sqlite>
  
```

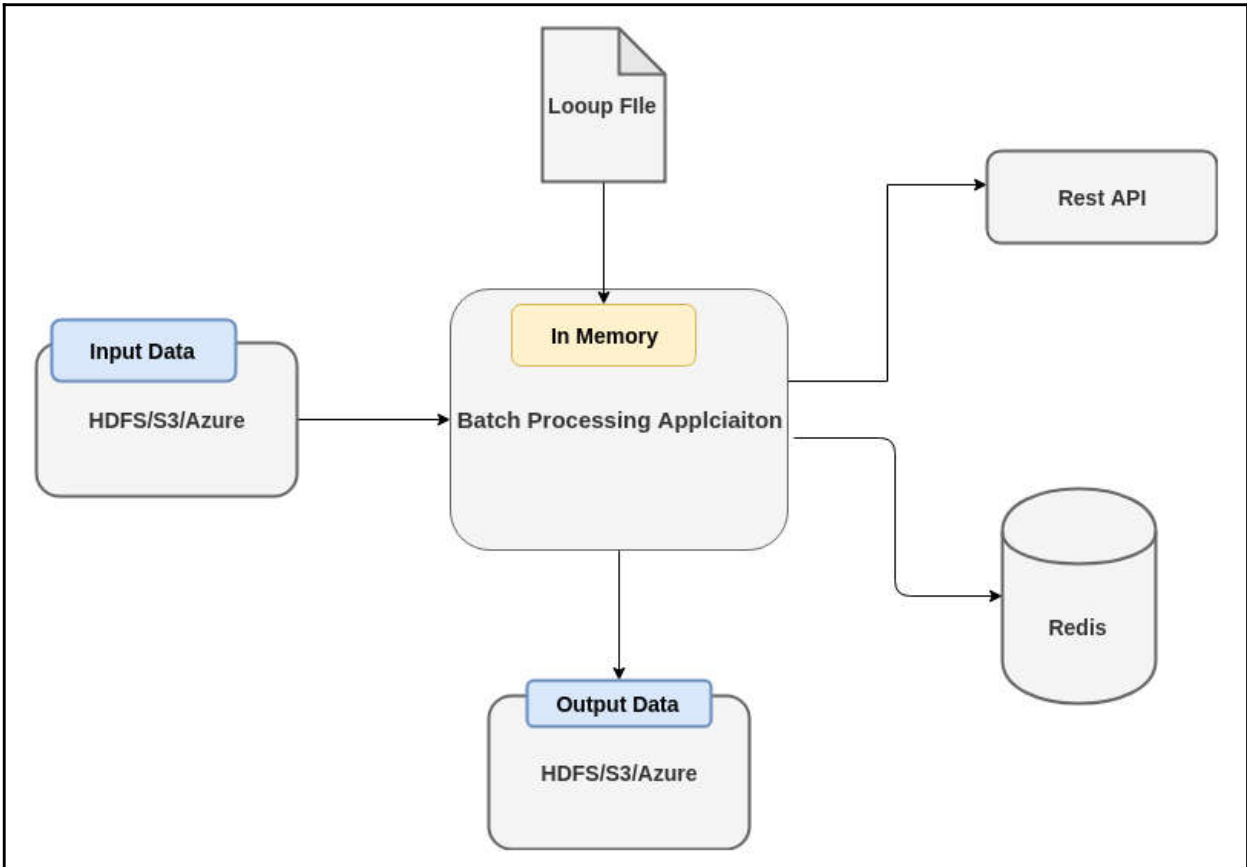
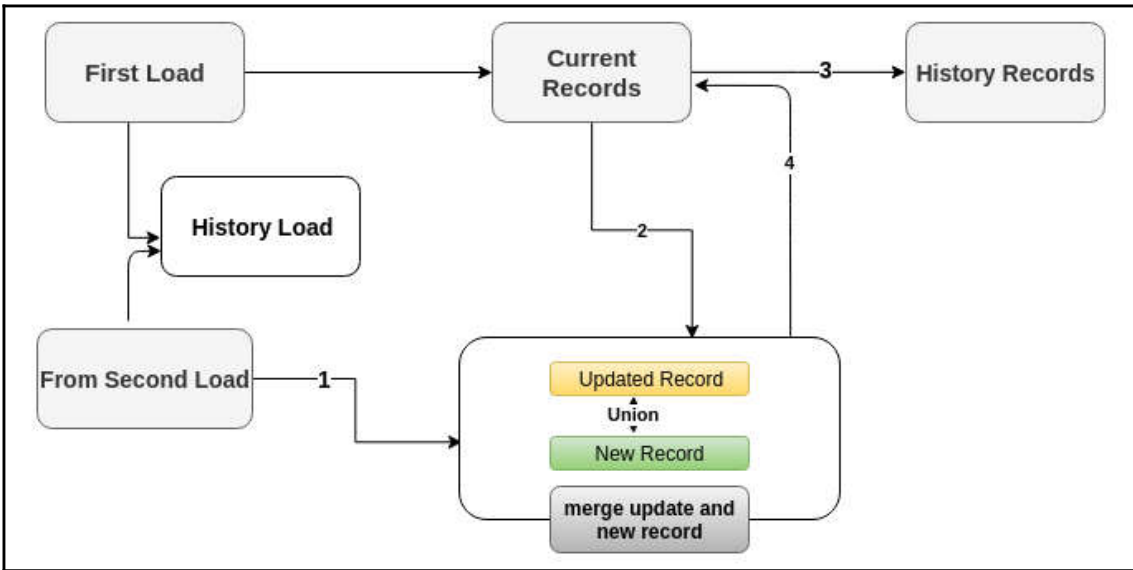




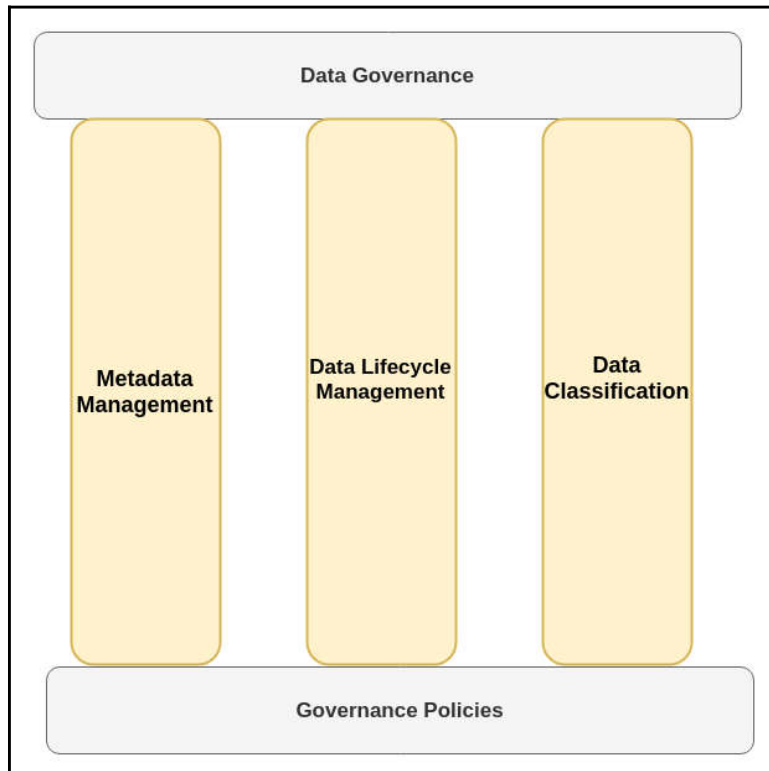
Chapter 8: Designing Applications in Hadoop



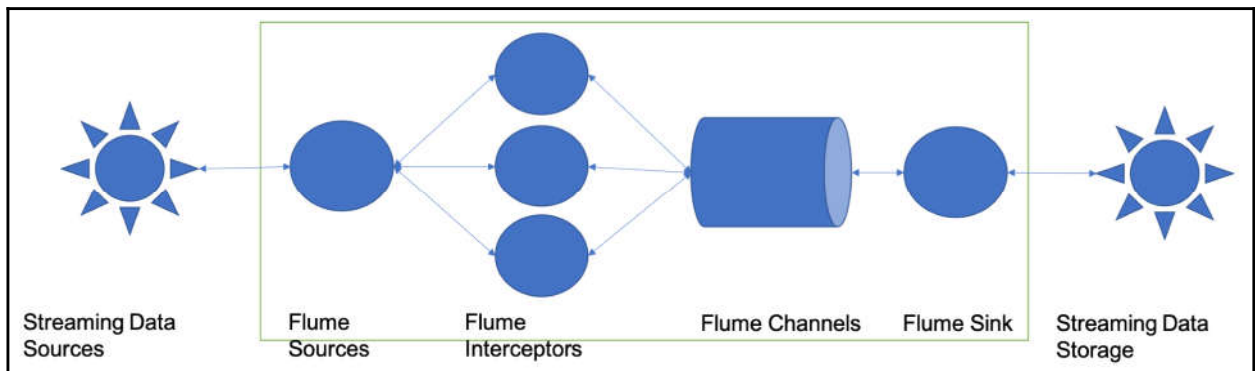
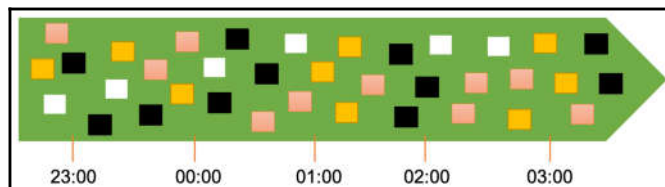
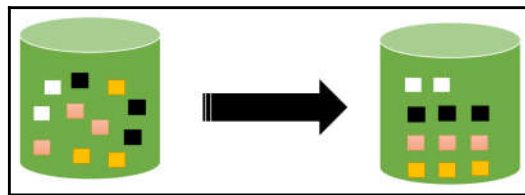


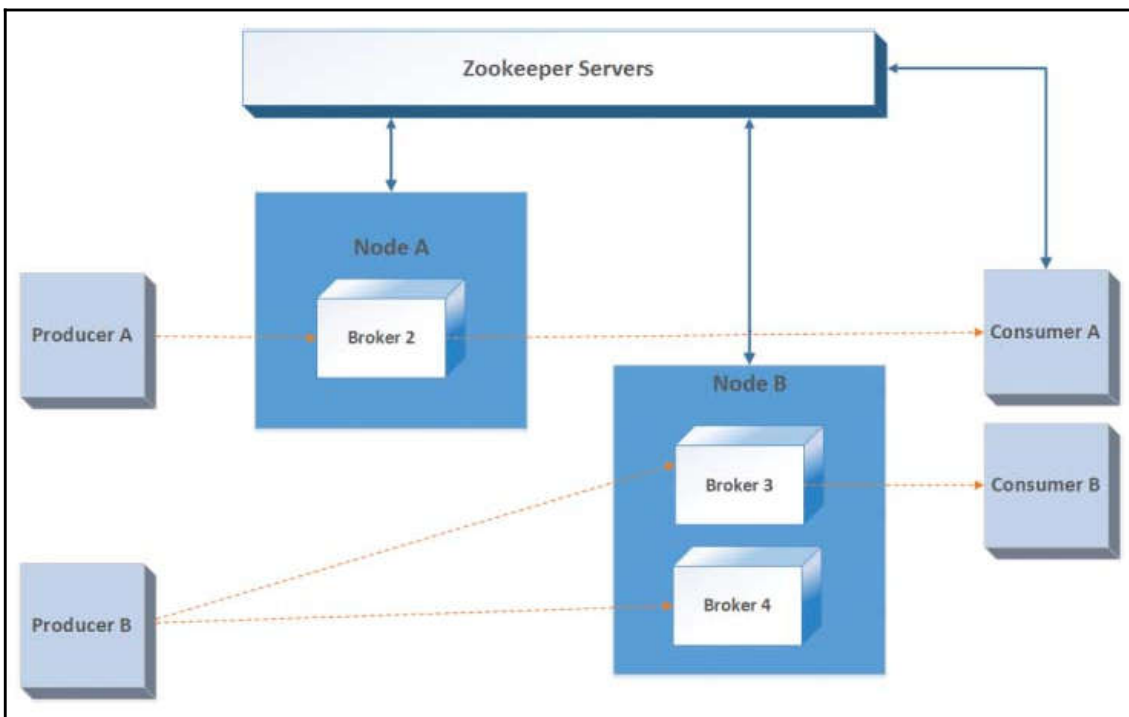
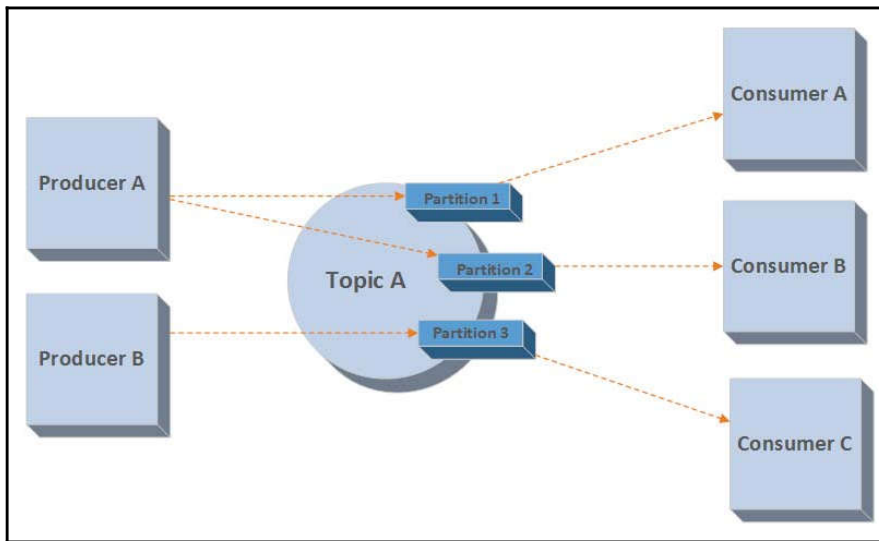


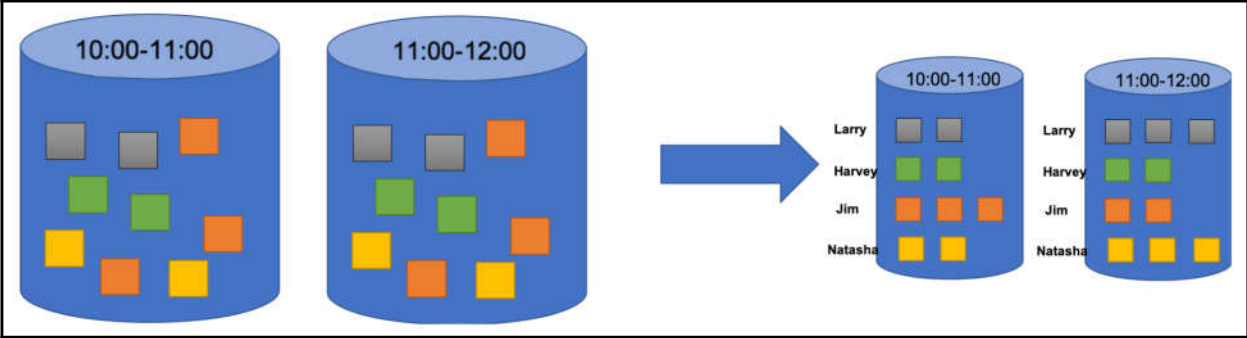
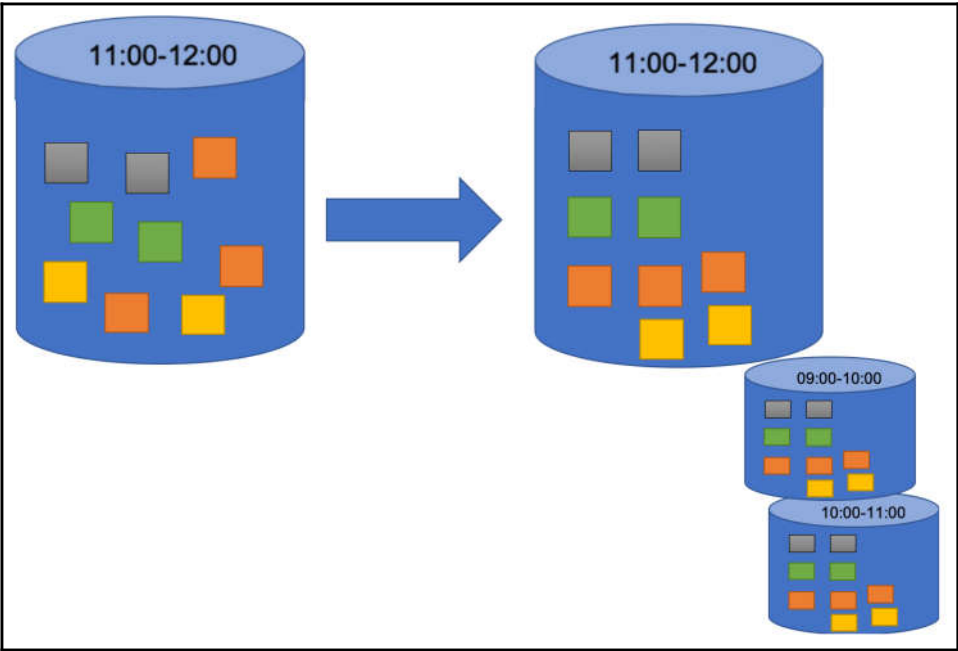
Airflow		DAGs	Data Profiling	Browse	Admin	Docs	About	17:39 UTC
<input type="checkbox"/>	Off	example_bash_operator	0 0 *	airflow				
<input type="checkbox"/>	Off	example_branch_dop_operator_v3	* * * * *	airflow				
<input type="checkbox"/>	Off	example_branch_operator	@daily	airflow				
<input type="checkbox"/>	Off	example_http_operator	1 day, 0:00:00	airflow				
<input type="checkbox"/>	Off	example_passing_params_via_test_command	* * * * *	airflow				
<input type="checkbox"/>	Off	example_python_operator	None	airflow				
<input type="checkbox"/>	Off	example_short_circuit_operator	1 day, 0:00:00	airflow				
<input type="checkbox"/>	Off	example_skip_dag	1 day, 0:00:00	airflow				
<input type="checkbox"/>	Off	example_subdag_operator	@once	airflow				
<input type="checkbox"/>	Off	example_trigger_controller_dag	@once	airflow				
<input type="checkbox"/>	Off	example_trigger_target_dag	None	airflow				
<input type="checkbox"/>	Off	example_xcom	@once	airflow				
<input type="checkbox"/>	Off	latest_only	4:00:00	Airflow				
<input type="checkbox"/>	Off	latest_only_with_trigger	4:00:00	Airflow				
<input type="checkbox"/>	Off	test_utils	None	airflow				



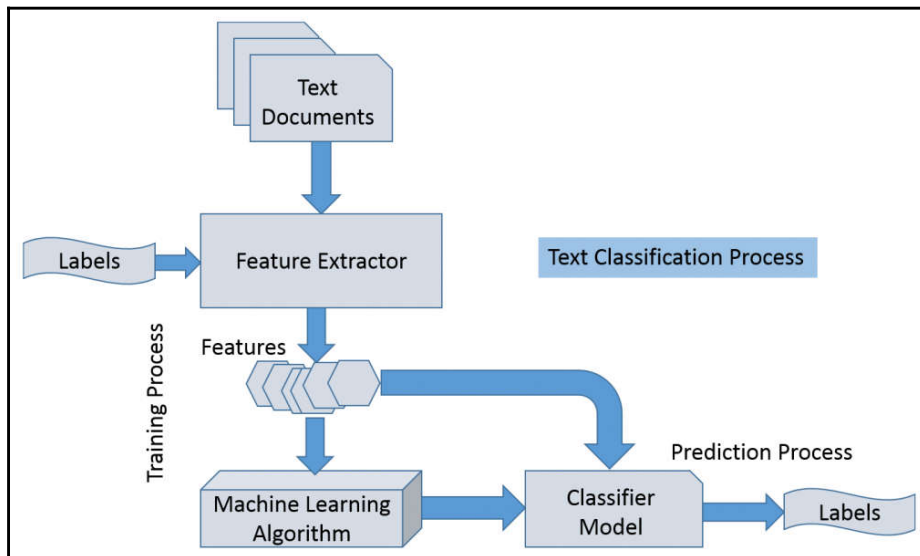
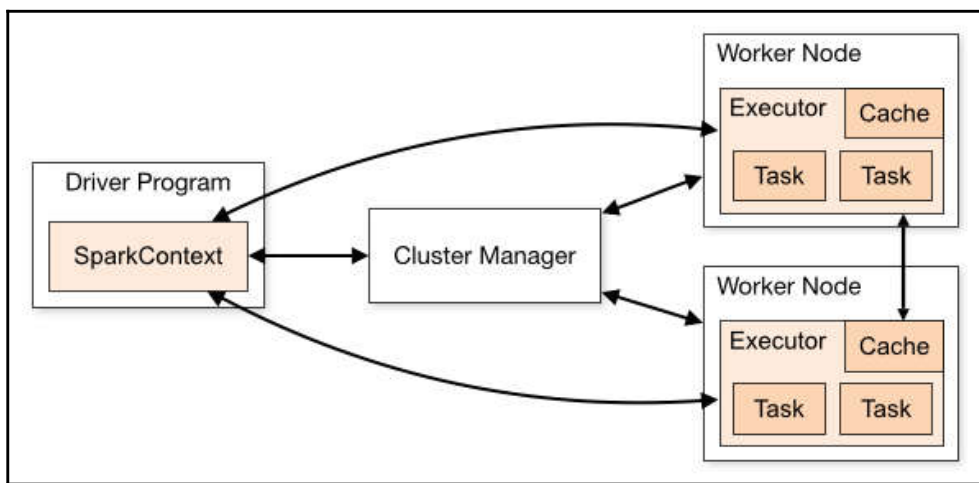
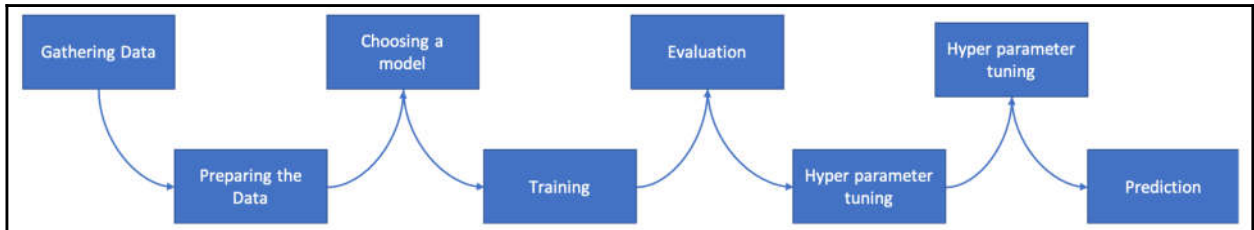
Chapter 9: Real-Time Stream Processing in Hadoop







Chapter 10: Machine Learning in Hadoop



Chapter 11: Hadoop in the Cloud

VPCs > Create VPC

Create VPC

A VPC is an isolated portion of the AWS cloud populated by AWS objects, such as Amazon EC2 instances. You must specify an IPv4 address range for your VPC. Specify the IPv4 address a Classless Inter-Domain Routing (CIDR) block; for example, 10.0.0.0/16. You cannot specify an IPv4 CIDR block larger than /16. You can optionally associate an Amazon-provided IPv6 CIDR with the VPC.

Name tag ⓘ

IPv4 CIDR block* ⓘ

IPv6 CIDR block No IPv6 CIDR Block ⓘ
 Amazon provided IPv6 CIDR block

Tenancy ⓘ

* Required Cancel

VPC Dashboard

Filter by VPC: Add filter

Name	VPC ID	State	IPv4 CIDR	IPv6 CIDR	DHCP options set	Route table
Packt-vpc	vpc-0b620306aa8f825d8	available	10.20.0.0/16	-	dopt-30218455	rtb-0e3edc5

VPC: vpc-0b620306aa8f825d8

VPC ID	vpc-0b620306aa8f825d8	Tenancy	default
State	available	Default VPC	No
IPv4 CIDR	10.20.0.0/16	Classic link	Disabled
IPv6 CIDR	-	DNS resolution	Enabled
Network ACL	acl-064712b338320a349	DNS hostnames	Disabled
DHCP options set	dopt-30218455	ClassicLink DNS Support	Disabled
Route table	rtb-0e3edc5183b231f36	Owner	960226201745

[Subnets](#) > Create subnet

Create subnet

Specify your subnet's IP address block in CIDR format; for example, 10.0.0.0/24. IPv4 block sizes must be between a /16 netmask and /28 netmask, and can be the same size as your VPC. CIDR block must be a /64 CIDR block.

Name tag ⓘ

VPC* ⓘ

VPC CIDRs	CIDR	Status	Status Reason
	10.20.0.0/16	associated	

Availability Zone ⓘ

IPv4 CIDR block* ⓘ

* Required Cancel

VPC Dashboard

[Create VPC](#) [Actions](#)

Filter by VPC: Add filter | 1 to 1 of 1

Name	VPC ID	State	IPv4 CIDR	IPv6 CIDR	DHCP options set	Route table
Packt-vpc	vpc-0b620306aa8f825d8	available	10.20.0.0/16	-	dopt-30218455	rtb-0e3edc5

VPC: vpc-0b620306aa8f825d8

Description

VPC ID: vpc-0b620306aa8f825d8

State: available

IPv4 CIDR: 10.20.0.0/16

IPv6 CIDR: -

Network ACL: acl-064712b338320a349

DHCP options set: dopt-30218455

Route table: rtb-0e3edc5183b231f36

Tenancy: default

Default VPC: No

Classic link: Disabled

DNS resolution: Enabled

DNS hostnames: Disabled

ClassicLink DNS Support: Disabled

Owner: 960226201745

EC2 Dashboard

Events

Tags

Reports

Limits

INSTANCES

Instances

Launch Templates

Spot Requests

Reserved Instances

Dedicated Hosts

Scheduled Instances

Capacity Reservations

IMAGES

AMIs

Bundle Tasks

ELASTIC BLOCK STORE

Volumes

Snapshots

Lifecycle Manager

Create Security Group

Security group name:

Description:

VPC:

Security group rules:

Inbound Outbound

Type	Protocol	Port Range	Source	Description
Custom TCF	TCP	22	Custom 0.0.0.0/0	e.g. SSH for Ad

1. Choose AMI 2. Choose Instance Type 3. Configure Instance 4. Add Storage 5. Add Tags 6. Configure Security Group 7. Review

Step 1: Choose an Amazon Machine Image (AMI)

	Ubuntu Server 16.04 LTS (HVM), SSD Volume Type - ami-076e276d85f524150 (64-bit x86) / ami-05e1b2aec3b47890f (64-bit Arm)	<input type="button" value="Select"/>
Free tier eligible	Ubuntu Server 16.04 LTS (HVM), EBS General Purpose (SSD) Volume Type. Support available from Canonical (http://www.ubuntu.com/cloud/services).	<input checked="" type="radio"/> 64-bit (x86) <input type="radio"/> 64-bit (Arm)
	Microsoft Windows Server 2016 Base - ami-0be9369ce05f0a8ff	<input type="button" value="Select"/>
Windows Free tier eligible	Microsoft Windows 2016 Datacenter edition. [English]	64-bit (x86)
	Deep Learning AMI (Ubuntu) Version 19.0 - ami-05bc59103c52af154	<input type="button" value="Select"/>
	With latest deep learning frameworks pre-installed: MXNet, TensorFlow, PyTorch, Keras, Chainer, Caffe/2, Theano & CNTK, configured with NVIDIA CUDA, cuDNN, NCCL & Intel MKL-DNN. For a fully managed experience, check: https://aws.amazon.com/sagemaker	64-bit (x86)
	Root device type: ebs Virtualization type: hvm ENA Enabled: Yes	

1. Choose AMI 2. Choose Instance Type 3. Configure Instance 4. Add Storage 5. Add Tags 6. Configure Security Group 7. Review

Step 2: Choose an Instance Type

<input type="radio"/>	General purpose	a1.4xlarge	16	32	EBS only	Yes	Up to 10 Gigabit
<input type="radio"/>	Compute optimized	c5n.large	2	5.25	EBS only	Yes	Up to 25 Gigabit
<input type="radio"/>	Compute optimized	c5n.xlarge	4	10.5	EBS only	Yes	Up to 25 Gigabit
<input type="radio"/>	Compute optimized	c5n.2xlarge	8	21	EBS only	Yes	Up to 25 Gigabit
<input checked="" type="radio"/>	Compute optimized	c5n.4xlarge	16	42	EBS only	Yes	Up to 25 Gigabit
<input type="radio"/>	Compute optimized	c5n.9xlarge	36	96	EBS only	Yes	50 Gigabit
<input type="radio"/>	Compute optimized	c5n.18xlarge	72	192	EBS only	Yes	100 Gigabit
<input type="radio"/>	Compute optimized	c5d.large	2	4	1 x 50 (SSD)	Yes	Up to 10 Gigabit
<input type="radio"/>	Compute optimized	c5d.xlarge	4	8	1 x 100 (SSD)	Yes	Up to 10 Gigabit
<input type="radio"/>	Compute optimized	c5d.2xlarge	8	16	1 x 200 (SSD)	Yes	Up to 10 Gigabit
<input type="radio"/>	Compute optimized	c5d.4xlarge	16	32	1 x 400 (SSD)	Yes	Up to 10 Gigabit

Cancel Previous **Review and Launch** Next: Configure Instance

1. Choose AMI 2. Choose Instance Type 3. Configure Instance 4. Add Storage 5. Add Tags 6. Configure Security Group 7. Review

Step 3: Configure Instance Details

Configure the instance to suit your requirements. You can launch multiple instances from the same AMI, request Spot instances to take advantage of the lower pricing, assign an access management role to the instance, and more.

Number of instances [Launch into Auto Scaling Group](#)

Purchasing option Request Spot instances

Network [Create new VPC](#)

Subnet [Create new subnet](#)
251 IP Addresses available

Auto-assign Public IP

Placement group Add instance to placement group.

Capacity Reservation [Create new Capacity Reservation](#)

IAM role [Create new IAM role](#)

Cancel Previous **Review and Launch** Next: Add

1. Choose AMI 2. Choose Instance Type 3. Configure Instance 4. Add Storage 5. Add Tags 6. Configure Security Group 7. Review

Step 4: Add Storage

Your instance will be launched with the following storage device settings. You can attach additional EBS volumes and instance store volumes to your instance, or edit the settings of the root volume. You can also attach additional EBS volumes after launching an instance, but not instance store volumes. [Learn more](#) about storage options in Amazon EC2.

Volume Type	Device	Snapshot	Size (GiB)	Volume Type	IOPS	Throughput (MB/s)	Delete on Termination	Encrypted
Root	/dev/sda1	snap-0252bea5b37202c35	8	General Purpose SSD (gp2)	100 / 3000	N/A	<input checked="" type="checkbox"/>	Not Encrypted

[Add New Volume](#)

Free tier eligible customers can get up to 30 GB of EBS General Purpose (SSD) or Magnetic storage. [Learn more](#) about free usage tier eligibility and usage restrictions.

[Cancel](#) [Previous](#) [Review and Launch](#) [Next: Add](#)

Step 1: Software and Steps

Step 2: Hardware
Step 3: General Cluster Settings
Step 4: Security

Software Configuration

Release:

<input checked="" type="checkbox"/> Hadoop 2.8.5	<input type="checkbox"/> Zeppelin 0.8.0	<input type="checkbox"/> Livy 0.5.0
<input type="checkbox"/> JupyterHub 0.9.4	<input type="checkbox"/> Tez 0.8.4	<input type="checkbox"/> Flink 1.6.1
<input type="checkbox"/> Ganglia 3.7.2	<input type="checkbox"/> HBase 1.4.7	<input checked="" type="checkbox"/> Pig 0.17.0
<input checked="" type="checkbox"/> Hive 2.3.3	<input type="checkbox"/> Presto 0.212	<input type="checkbox"/> ZooKeeper 3.4.13
<input type="checkbox"/> MXNet 1.3.0	<input type="checkbox"/> Sqoop 1.4.7	<input type="checkbox"/> Mahout 0.13.0
<input checked="" type="checkbox"/> Hue 4.2.0	<input type="checkbox"/> Phoenix 4.14.0	<input type="checkbox"/> Oozie 5.0.0
<input type="checkbox"/> Spark 2.3.2	<input type="checkbox"/> HCatalog 2.3.3	<input type="checkbox"/> TensorFlow 1.11.0

AWS Glue Data Catalog settings (optional)

Use for Hive table metadata

Edit software settings

Enter configuration Load JSON from S3

```
classification=config-file-name,properties=[myKey1=myValue1,myKey2=myValue2]
```

Add steps (optional)

Step type: [Configure](#)

Create Cluster - Advanced Options [Go to quick options](#)

Step 1: Software and Steps
Step 2: Hardware
 Step 3: General Cluster Settings
 Step 4: Security

Hardware Configuration ⓘ

If you need more than 20 EC2 instances, [see this topic](#).

Instance group configuration

- Uniform instance groups**
Specify a single instance type and purchasing option for each node type.
- Instance fleets**
Specify target capacity and how Amazon EMR fulfills it for each node type. Mix instance types and purchasing options. [Learn more](#)

Network [Create a VPC](#) ⓘ

EC2 Subnet

No Amazon S3 endpoint or route to AWS public IP Range detected in subnet-023ae6c85c8a2c5d2. S3 endpoints are recommended for EMR access data in S3. For EMR to communicate with Amazon DynamoDB (and for EMRFS consistent view), AWS KMS, and Amazon Kinesis, your cluster must be able to reach this IP range. You can create a NAT instance to create a route for this traffic. [Learn more](#)

[Add S3 endpoint and NAT instance](#)

Root device EBS volume size GiB ⓘ

EC2 Dashboard

- Events
- Tags
- Reports
- Limits
- INSTANCES
 - Instances
 - Launch Templates
 - Spot Requests
 - Reserved Instances
 - Dedicated Hosts
 - Scheduled Instances
 - Capacity Reservations
- IMAGES
 - AMIs
 - Bundle Tasks
- ELASTIC BLOCK STORE
 - Volumes
 - Snapshots
 - Lifecycle Manager
- NETWORK & SECURITY

Resources

You are using the following Amazon EC2 resources in the US East (N. Virginia) region:

155 Running Instances	14 Elastic IPs
0 Dedicated Hosts	3 Snapshots
208 Volumes	10 Load Balancers
9 Key Pairs	107 Security Groups
1 Placement Groups	

[Learn more about the latest in AWS Compute from AWS re:Invent by viewing the EC2 Videos.](#)

Create Instance

To start using Amazon EC2 you will want to launch a virtual server, known as an Amazon EC2 instance.

[Launch Instance](#)

Note: Your instances will launch in the US East (N. Virginia) region

Service Health

Service Status:
 US East (N. Virginia): ✔

Availability Zone Status:

Scheduled Events

US East (N. Virginia):
No events

Account Attributes

Supported Platforms
VPC

Default VPC
vpc-d1b52eb5

Resource ID length management
Console experiments

Additional Information

- Getting Started Guide
- Documentation
- All EC2 Resources
- Forums
- Pricing
- Contact Us

AWS Marketplace

Find free software trial products in AWS Marketplace from the [EC2 Launch Wizard](#). Or try these popular AMIs:

[Barracuda CloudGen Firewall for AWS](#)

EC2 Dashboard

Launch Instance Connect Actions

Filter by tags and attributes or search by keyword 1 to 50 of 214

Name	Instance ID	Instance Type	Availability Zone	Instance State	Status Checks
qubole-bastion	i-0792f317a35ca4721	t2.small	us-east-1d	running	2/2 checks ...
node0060	i-0851fbec385501735	r3.4xlarge	us-east-1d	terminated	
node0059	i-03936e9f04715aabd	r3.4xlarge	us-east-1d	terminated	
node0058	i-079cb3988801ca63d	r3.4xlarge	us-east-1d	running	2/2 checks ...
node0057	i-0bf2d2eb9e2190882	r3.4xlarge	us-east-1d	terminated	
node0056	i-04febdd30a9bd93e0	r3.4xlarge	us-east-1d	terminated	
node0055	i-0970e16d31ec407e8	r3.4xlarge	us-east-1d	terminated	
node0054	i-0f7f470c73e516dee	r3.4xlarge	us-east-1d	terminated	
node0053	i-02a3f8b2804ffaa1c	r3.4xlarge	us-east-1d	terminated	
node0052	i-0ee5429b25e6de0da	r3.4xlarge	us-east-1d	terminated	

Select an instance above

EC2 Dashboard

Launch Instance Connect Actions

Filter by tags and attributes or search by keyword 1 to 50 of 214

Name	Instance ID	Instance Type	Availability Zone	Instance State	Status Checks
node0059	i-03936e9f04715aabd	r3.4xlarge	us-east-1d	terminated	

Instance: i-079cb3988801ca63d (node0058) Public DNS: ec2-34-204-45-169.compute-1.amazonaws.com

Description Status Checks Monitoring Tags

Instance ID	i-079cb3988801ca63d	Public DNS (IPv4)	[REDACTED]
Instance state	running	IPv4 Public IP	[REDACTED]
Instance type	r3.4xlarge	IPv6 IPs	-
Elastic IPs		Private DNS	ip-10-200-1-215.ec2.internal
Availability zone	us-east-1d	Private IPs	10.200.1.215
Security groups	@sc-qbol_acc68_cl16324, analytics-qubole-presto-sg. view inbound rules. view outbound rules	Secondary private IPs	
Scheduled events	No scheduled events	VPC ID	vpc-a483e5c0
AMI ID	qbol-hvm-2018-12-18 121-11-52 (ami-0be70a8fd074a5c6f)	Subnet ID	subnet-9db060eb
Platform	-	Network interfaces	eth0
IAM role	[REDACTED]	Source/dest. check	True

EC2 Dashboard

Launch Instance | Connect | Actions

Filter by tags and attributes or search by keyword

Name	Instance ID	Instance Type	Availability Zone	Instance State	Status Checks
qubole-bastion	i-0792f317a35ca4721	t2.small	us-east-1d	running	2/2 checks ...
node0060	i-0851fbec385501735	r3.4xlarge	us-east-1d	terminated	
node0059	i-03936e9f04715aabd	r3.4xlarge	us-east-1d	terminated	
node0058	i-03988801ca63d	r3.4xlarge	us-east-1d	running	2/2 checks ...
node0057	i-02eb9e2190882	r3.4xlarge	us-east-1d	terminated	
node0056	i-0dd30a9bd93e0	r3.4xlarge	us-east-1d	terminated	
node0055	i-0e16d31ec407e8	r3.4xlarge	us-east-1d	terminated	
node0054	i-070c73e516dee	r3.4xlarge	us-east-1d	terminated	
node0053	i-0f8b2804ffaa1c	r3.4xlarge	us-east-1d	terminated	
node0052	i-070c73e516dee	r3.4xlarge	us-east-1d	terminated	

Instance: i-079cb3988801ca63d

Connect

- Get Windows Password
- Create Template From Instance
- Launch More Like This
- Instance State
- Instance Settings
- Image
- Networking
- CloudWatch Monitoring

- Add/Edit Tags
- Attach to Auto Scaling Group
- Attach/Replace IAM Role
- Change Instance Type
- Change Termination Protection
- View/Change User Data
- Change Shutdown Behavior
- Change T2/T3 Unlimited
- Get System Log
- Get Instance Screenshot
- Modify Instance Placement

Description | Status Checks | Monitoring | Tags

Instance ID: i-079cb3988801ca63d

Instance state: running

Public DNS (IPv4): ec2-34-204-45-169.compute-1.amazonaws.com

4 Public IP: 34.204.45.169

EC2 Dashboard

Launch Instance | Connect | Actions

Filter by tags and attributes or search by keyword

Name	Instance ID	Instance Type	Availability Zone	Instance State	Status Checks
qubole-bastion	i-0792f317a35ca4721	t2.small	us-east-1d	running	2/2 checks ...
node0060	i-0851fbec385501735	r3.4xlarge	us-east-1d	terminated	
node0059	i-03936e9f04715aabd	r3.4xlarge	us-east-1d	terminated	
node0058	i-03988801ca63d	r3.4xlarge	us-east-1d	running	2/2 checks ...
node0057	i-02eb9e2190882	r3.4xlarge	us-east-1d	terminated	
node0056	i-0dd30a9bd93e0	r3.4xlarge	us-east-1d	terminated	
node0055	i-0e16d31ec407e8	r3.4xlarge	us-east-1d	terminated	
node0054	i-070c73e516dee	r3.4xlarge	us-east-1d	terminated	
node0053	i-0f8b2804ffaa1c	r3.4xlarge	us-east-1d	terminated	
node0052	i-070c73e516dee	r3.4xlarge	us-east-1d	terminated	

Instance: i-079cb3988801ca63d

Connect

- Get Windows Password
- Create Template From Instance
- Launch More Like This
- Instance State
- Instance Settings
- Image
- Networking
- CloudWatch Monitoring

- Change Security Groups
- Attach Network Interface
- Detach Network Interface
- Disassociate Elastic IP Address
- Change Source/Dest. Check
- Manage IP Addresses

Description | Status Checks | Monitoring | Tags

Instance ID: i-079cb3988801ca63d

Instance state: running

Public DNS (IPv4): [REDACTED]

IPv4 Public IP: [REDACTED]

EC2 Dashboard

Launch Instance Connect Actions

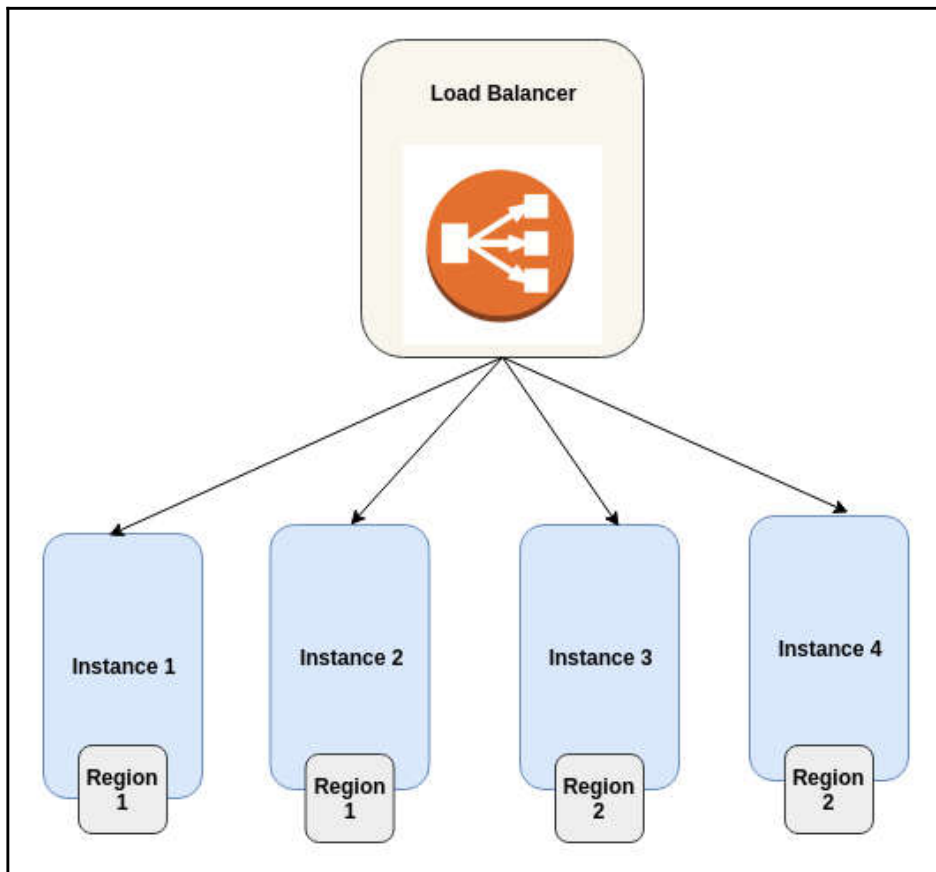
Filter by tags and attributes or search by keyword

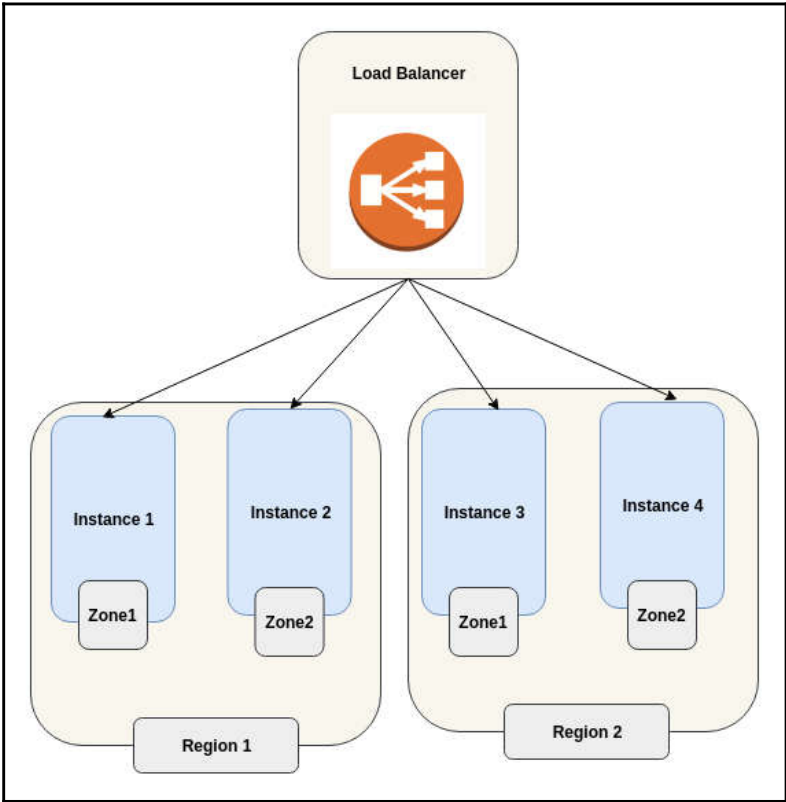
Name	Instance ID	Instance Type	Availability Zone	Instance State	Status Checks
qubole-bastion	i-0792f317a35ca4721	t2.small	us-east-1d	running	2/2 checks ...
node0060	i-0851fbec385501735	r3.4xlarge	us-east-1d	terminated	
node0059	i-03936e9f04715aabd	r3.4xlarge	us-east-1d	terminated	
node0058	i-0398801ca63d	r3.4xlarge	us-east-1d	running	2/2 checks ...
node0057	i-02eb9e2190882	r3.4xlarge	us-east-1d	terminated	
node0056	i-0dd30a9bd93e0	r3.4xlarge	us-east-1d	terminated	
node0055	i-0e16d31ec407e8	r3.4xlarge	us-east-1d	terminated	
node0054	i-0...	r3.4xlarge	us-east-1d	terminated	
node0053	i-0...	r3.4xlarge	us-east-1d	terminated	
node0052	i-0...	r3.4xlarge	us-east-1d	terminated	

Instance: i-079cb398801ca63d Public DNS (IPv4) [Redacted]

Instance ID: i-079cb398801ca63d Public DNS (IPv4) [Redacted]

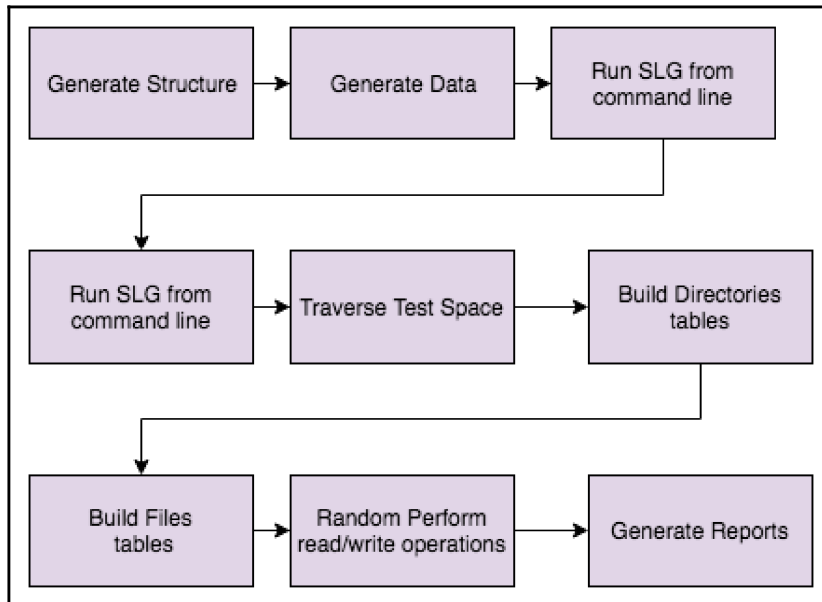
Instance state: running IPv4 Public IP [Redacted]

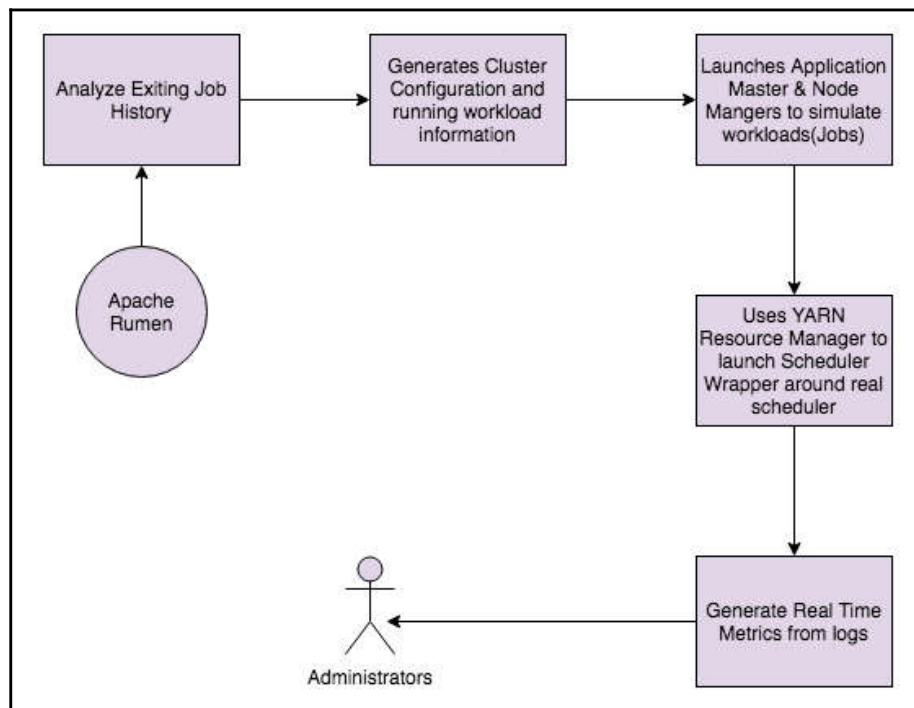




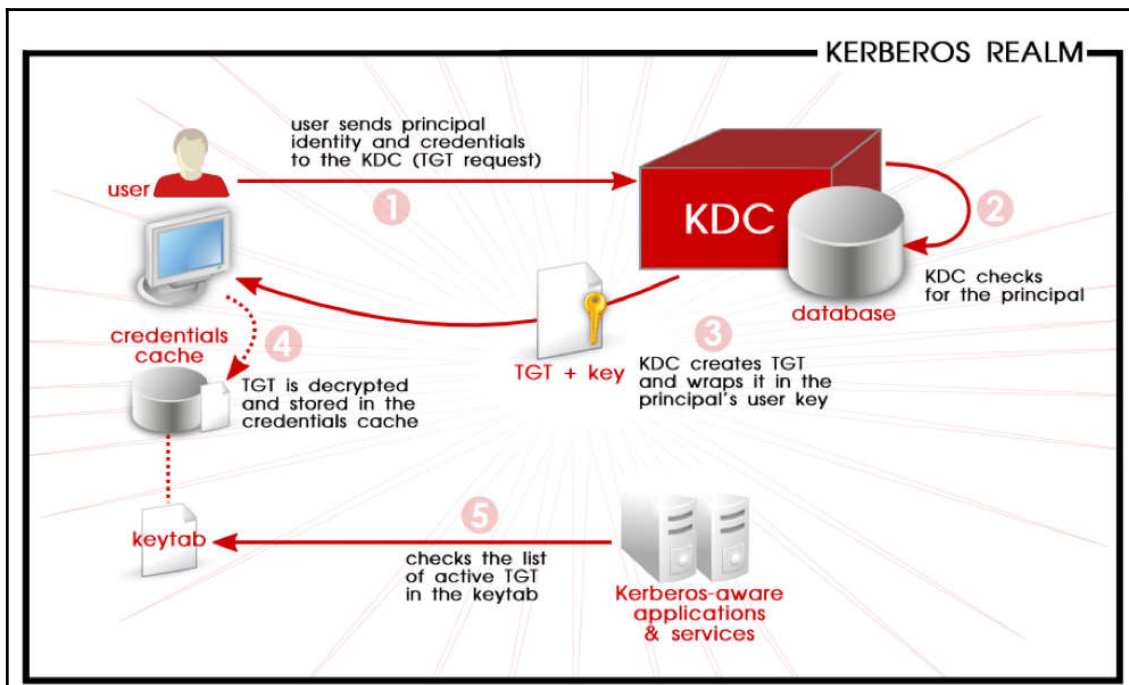
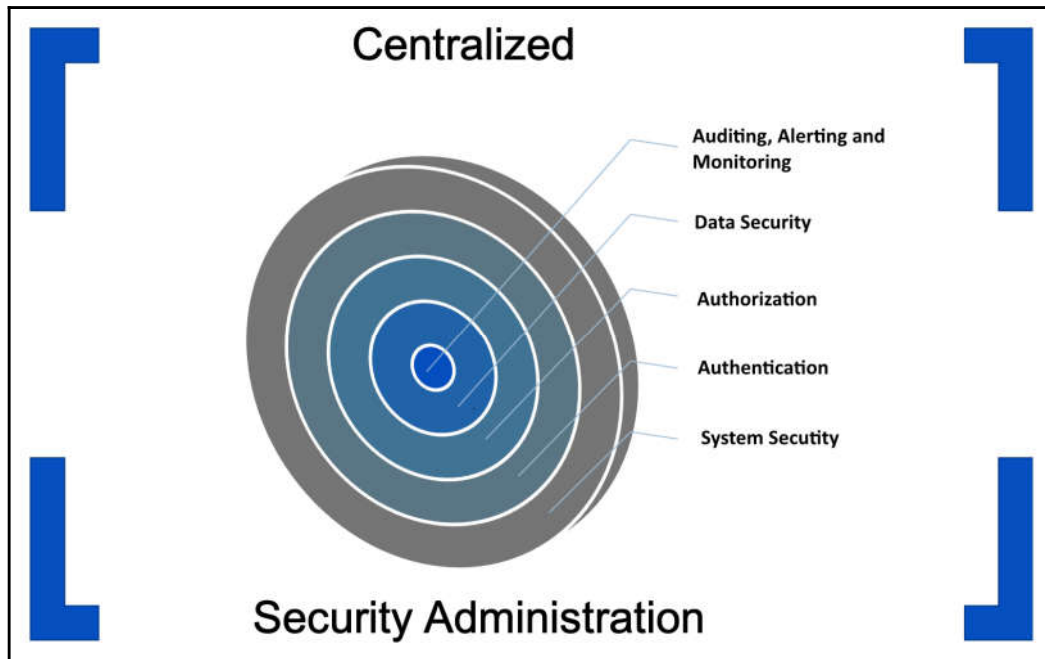
Chapter 12: Hadoop Cluster Profiling

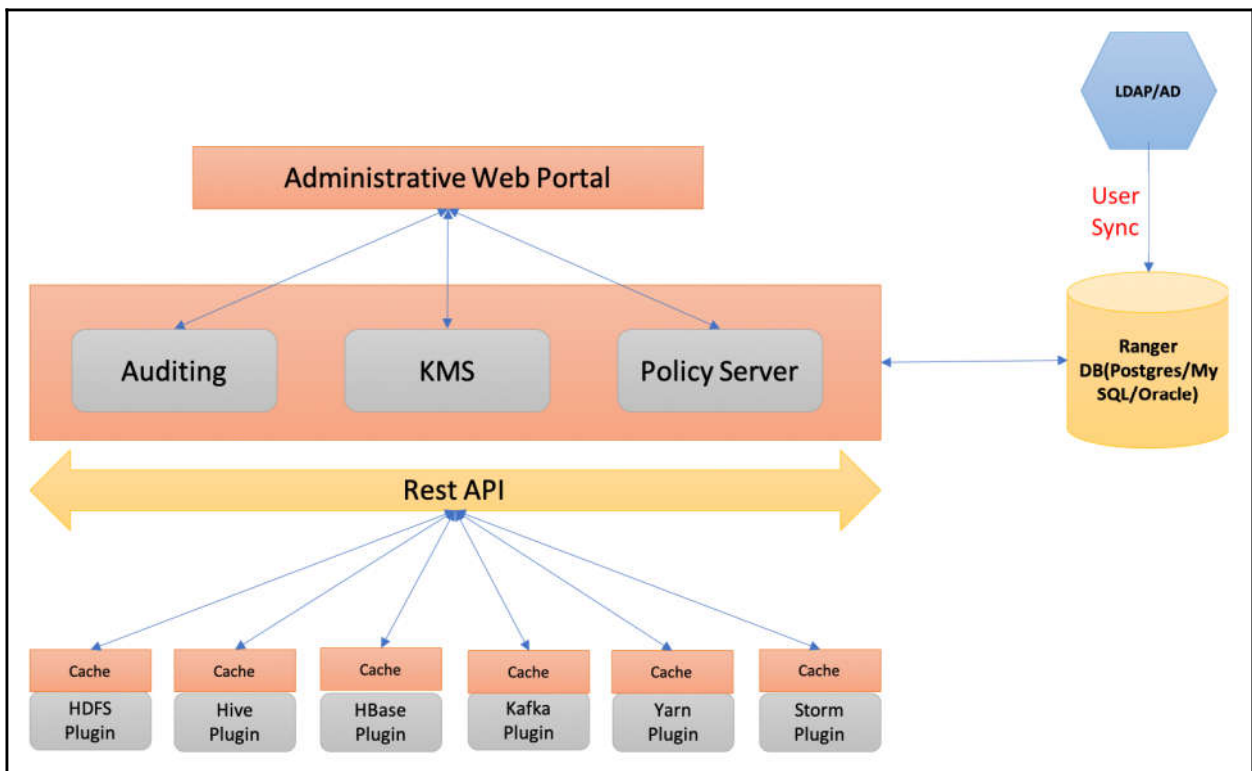
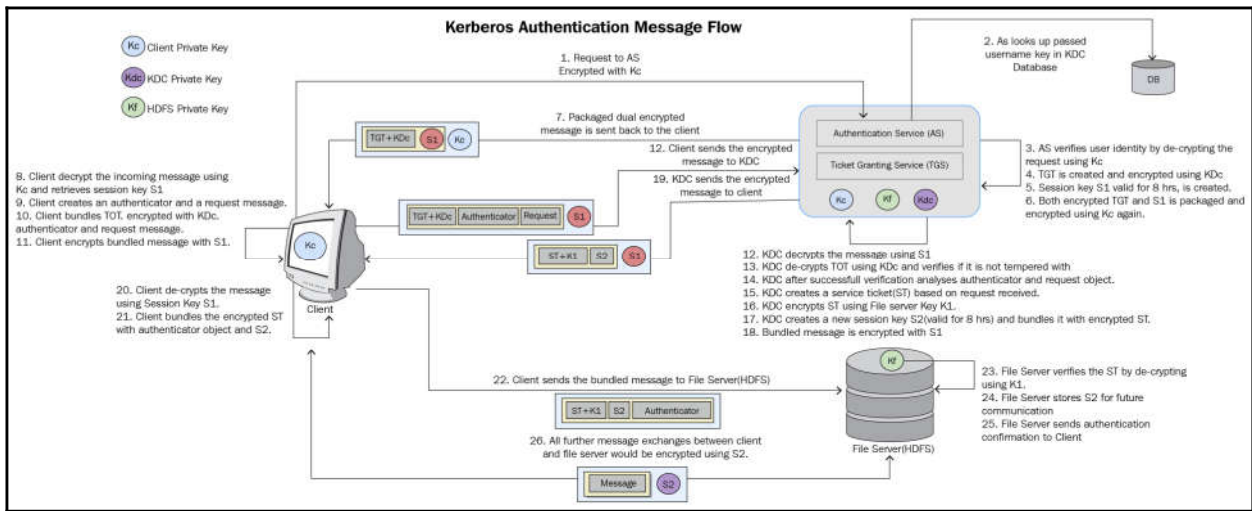
Benchmarking Hadoop Cluster		
HDFS	TestDFSIO	
Namenode	NNBench	NNThroughputBenchmark
	Synthetic Load Generator (SLG)	
YARN	YARN Scheduler Load Simulator (SLS)	
HIVE	TPC-DS	TPC-H
MIX-WORKLOADS	GRIDMIX	RUMEN

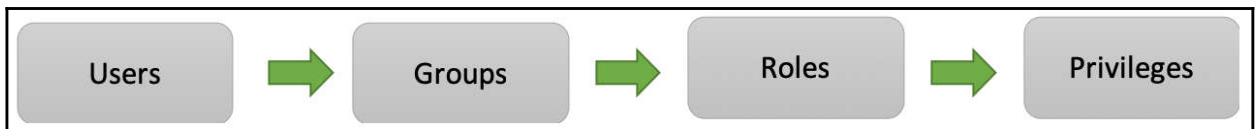
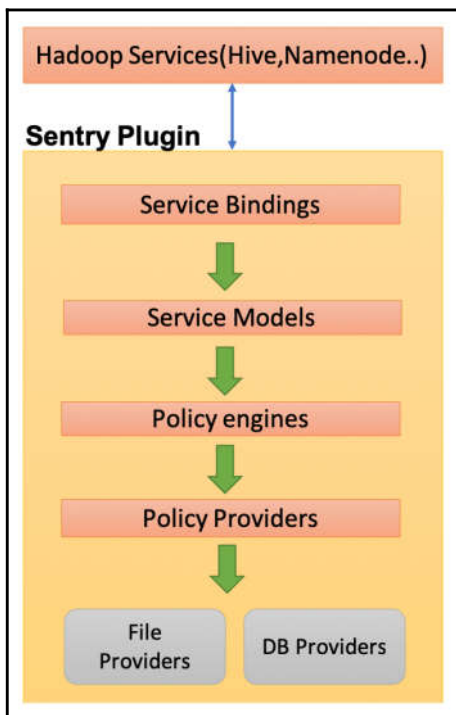




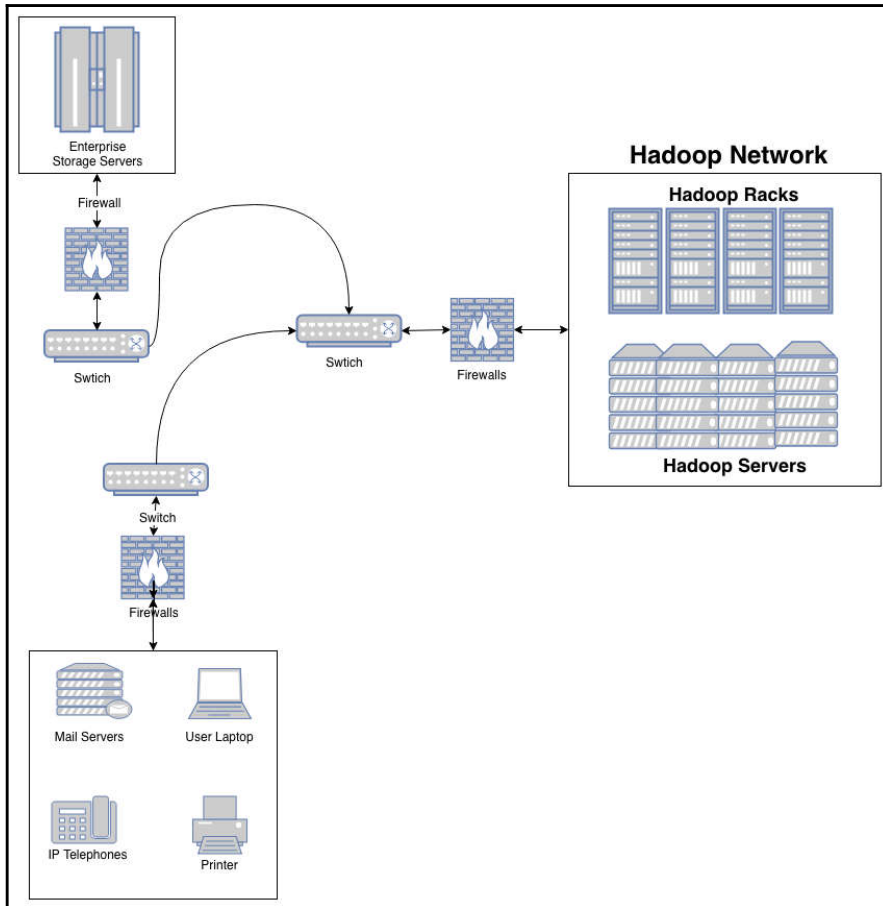
Chapter 13: Who Can Do What in Hadoop

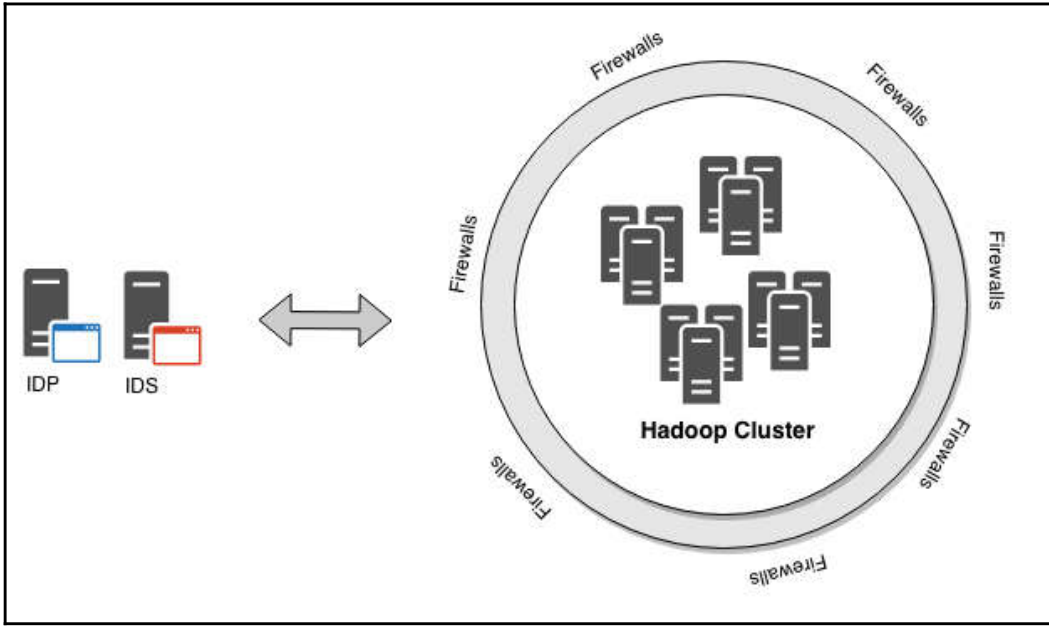


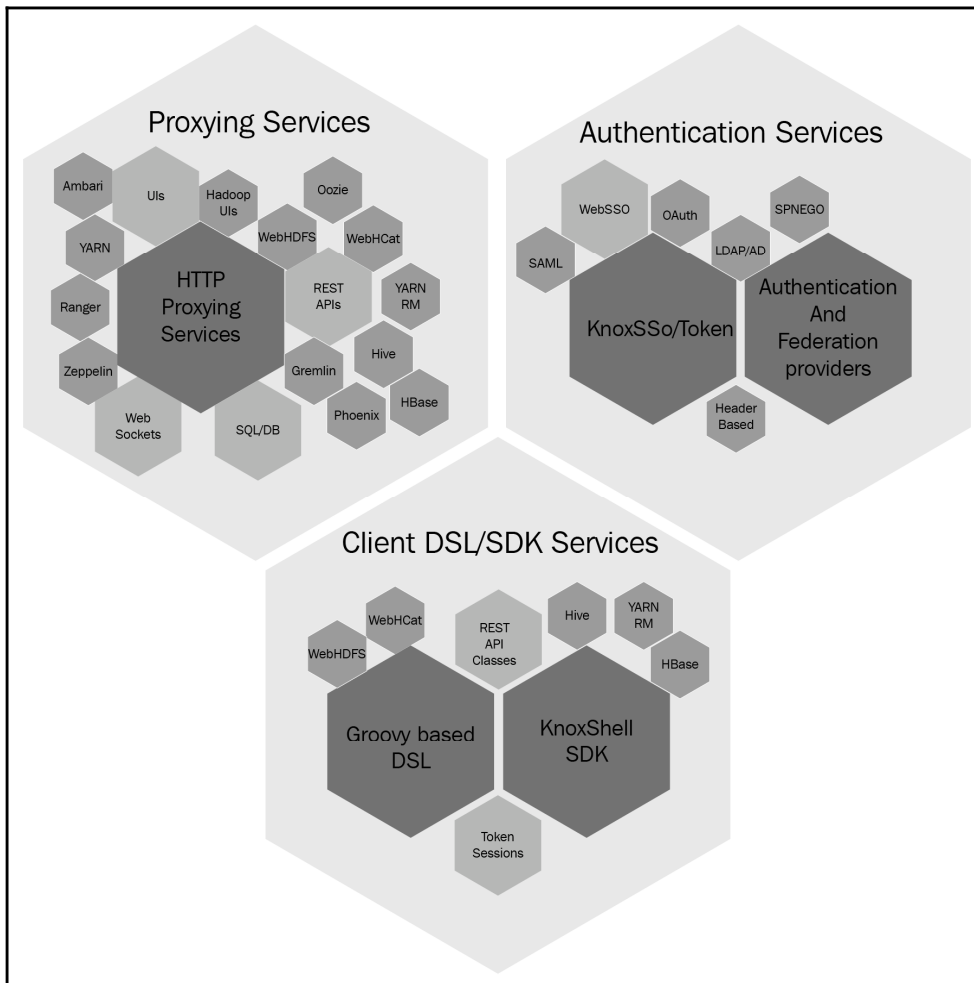




Chapter 14: Network and Data Security

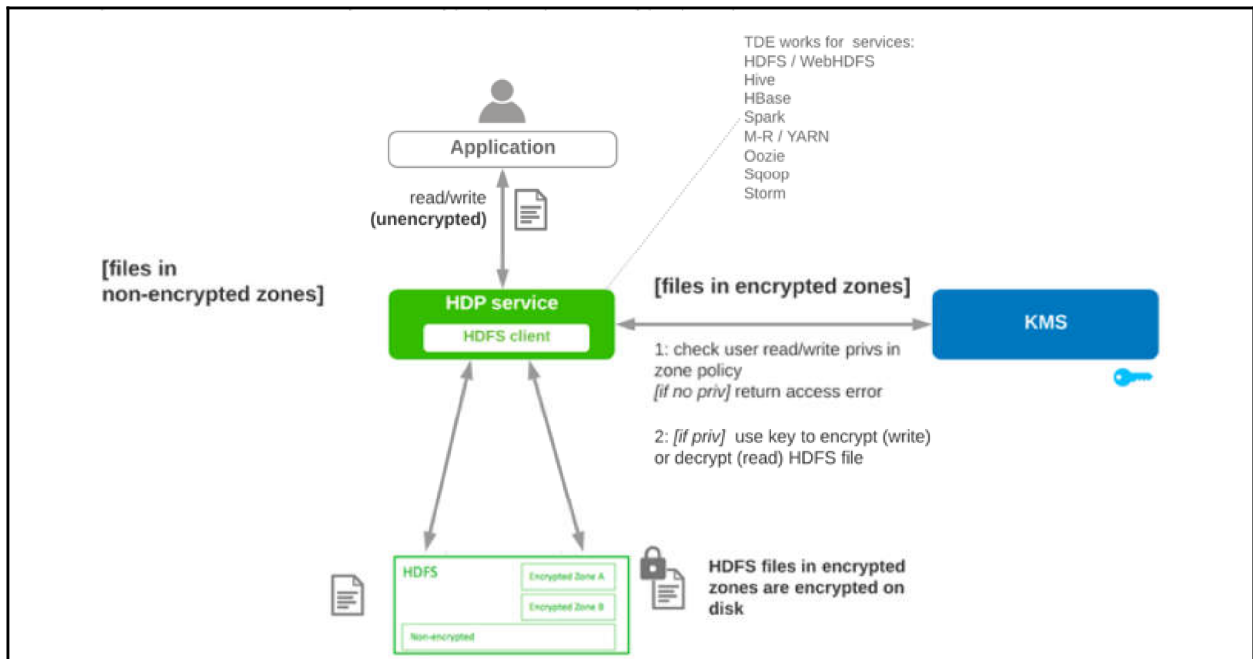






Hadoop Data In Transit Security

Transit Protocols		
RPC	TCP/IP	HTTP
Mapreduce	HDFS CLIENTS	Mapreduce Shuffles
JobTracker		Web Interfaces
TaskTracker		
NameNode		



Data Masking


ID	Employee Name	Employee Salary
1	John	10,000
2	Tim	20,000

↓

ID	Employee Name	Employee Salary
1	John	XXXXX
2	Tim	XXXXX

Data Masking By Random Substitution


ID	Employee Name	Employee Salary
1	John	10,000
2	Tim	20,000



ID	Employee Name	Employee Salary
1	John	5000
2	Tim	1000

Data Masking By Encryption


ID	Employee Name	Employee Salary
1	John	10,000
2	Tim	20,000



ID	Employee Name	Employee Salary
1	John	AB2H345EDNE98TYUO
2	Tim	SDF2096FT32UO7I9OP

Row Level Filtering

ID	Employee Name	Employee Salary
1	John	10,000
2	Tim	20,000



ID	Employee Name	Employee Salary
1	John	10,000

Data Masking By Random Substitution

ID	Employee Name	Employee Salary
1	John	10,000
2	Tim	20,000



ID	Employee Name
1	John
2	Tim

Chapter 15: Monitoring Hadoop



