# Assignment-12-CSP-554

**Exercise 1**

**Read the article "A Big Data Modeling Methodology for Apache Cassandra" available on the blackboard in the 'Articles' section. Provide a ½ page summary including your comments and impressions.**

**Summary:**

The paper covers traditional data modeling, Cassandra data modeling, conceptual and logical data modeling, application workflow, query-driven mapping from a conceptual to a logical data model, and physical data modeling.

**Cassandra Data Model:**

A CQL table can be considered a grouping of divisions containing rows with similar structures. A partition key is distinct from each partition in a table, whereas a clustering key is distinct from each row within a partition. A primary key is a combination of a partition key and a clustering key that uniquely identifies a database row. A table schema is a collection of columns that contains a primary key. Each column's data type is either primitive (int, text, etc.), complex (set, list, or map), or counter. CQL, which has a syntax similar to SQL, is used to express queries over tables. CQL does not support binary operations like joins and instead relies on a set of query predicates rules to ensure efficiency and scalability.

**Conceptual data modeling and application workflow:**

Understanding the data to be maintained and how a data-driven application needs to access it is required when designing a Cassandra database schema. The ER diagram depicts the former Application workflow diagrams, which define data access patterns for application tasks and capture the latter.

**Query driven mapping Data Modeling Principles:**

The four data modeling principles listed below serve as a foundation for translating conceptual data models into logical data models.

DMP1 (Know your data): The first step in successful database design is to understand the data, which is recorded using a conceptual data model.

DMP2 (Know your Questions): Knowing your queries captured by an application process is the second key to a successful database design.

DMP3 (Data Nesting): Data nesting is the third key to a successful database design.

DMP4 (Data Duplication): Data duplication is the fourth key to a successful database design. Mapping Rule: - The following are five mapping rules that facilitate a query-driven transition from a conceptual data model to a logical data model.

MR1 -> (Entities and Relationships): In MR1, entities and relationships map to table rows, whereas entity and relationship types of the map to tables

MR2 -> (Equality Search Attributes): In a query predicate, equality search attributes correspond to the prefix columns of a table's primary key.

MR3 -> (Inequality Search Attributes): A key column in a table clustering corresponds to an inequality search attribute used in a query predicate.

MR4 -> (Ordering Attributes): Ordering attributes, which are supplied in a query, map to clustering key columns in the query's chosen ascending or descending clustering order. MR5 -> (Key Attributes): Primary key columns are mapped to key attribute types.

Mapping Patterns: Mapping Patterns are used to automate Cassandra database schema design. Physical Data Modeling: The final step is to analyze and optimize a logical data model in a physical data model.

**Exercise 2**

**Using command**

wget https://archive.apache.org/dist/cassandra/3.11.2/apache-cassandra-3.11.2-bin.tar.gz

```
aasth@LAPTOP-HJTR6HMR MINGW64 ~
$ ssh -i Downloads/emr-key-pair.pem hadoop@ec2-52-86-39-223.compute-1.amazonaws.
com
The authenticity of host 'ec2-52-86-39-223.compute-1.amazonaws.com (52.86.39.223
)' can't be established.
ED25519 key fingerprint is SHA256:Iur2dzU77Yni+ARqtSn9UyVpO7juOPmnkEH6x/wlhCM.
This key is not known by any other names
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'ec2-52-86-39-223.compute-1.amazonaws.com' (ED25519)
to the list of known hosts.


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

https://aws.amazon.com/amazon-linux-2/
22 package(s) needed for security, out of 32 available
Run "sudo yum update" to apply all updates.


EEEEEEEEEEEEEEEEEEEE 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:::::M   R:::::::::::RR
  E:::::EEEEEEEEEE   M:::::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
EEEEEEEEEEEEEEEEEEEE MMMMMMM              MMMMMMM RRRRRRR      RRRRRR


[hadoop@ip-172-31-59-222 ~]$ wget https://archive.apache.org/dist/cassandra/3.11
.2/apache-cassandra-3.11.2-bin.tar.gz
--2022-12-01 22:44:33--  https://archive.apache.org/dist/cassandra/3.11.2/apache
-cassandra-3.11.2-bin.tar.gz
Resolving archive.apache.org (archive.apache.org)... 138.201.131.134, 2a01:4f8:1
72:2ec5::2
Connecting to archive.apache.org (archive.apache.org)|138.201.131.134|:443... co
nnected.
HTTP request sent, awaiting response... 200 OK
Length: 38436262 (37M) [application/x-gzip]
Saving to: 'apache-cassandra-3.11.2-bin.tar.gz'

100%[====================================>] 38,436,262  14.7MB/s   in 2.5s

2022-12-01 22:44:36 (14.7 MB/s) - 'apache-cassandra-3.11.2-bin.tar.gz' saved [38
436262/38436262]
```

**tar -xzvf apache-cassandra-3.11.2-bin.tar.gz**

```
[hadoop@ip-172-31-59-222 ~]$ tar -xzvf apache-cassandra-3.11.2-bin.tar.gz
apache-cassandra-3.11.2/bin/
apache-cassandra-3.11.2/conf/
apache-cassandra-3.11.2/conf/triggers/
apache-cassandra-3.11.2/doc/
apache-cassandra-3.11.2/doc/cql3/
apache-cassandra-3.11.2/doc/html/
apache-cassandra-3.11.2/doc/html/_images/
apache-cassandra-3.11.2/doc/html/_sources/
apache-cassandra-3.11.2/doc/html/_sources/architecture/
apache-cassandra-3.11.2/doc/html/_sources/configuration/
apache-cassandra-3.11.2/doc/html/_sources/cql/
apache-cassandra-3.11.2/doc/html/_sources/data_modeling/
apache-cassandra-3.11.2/doc/html/_sources/development/
apache-cassandra-3.11.2/doc/html/_sources/faq/
apache-cassandra-3.11.2/doc/html/_sources/getting_started/
apache-cassandra-3.11.2/doc/html/_sources/operating/
apache-cassandra-3.11.2/doc/html/_sources/tools/
apache-cassandra-3.11.2/doc/html/_sources/troubleshooting/
apache-cassandra-3.11.2/doc/html/_static/
apache-cassandra-3.11.2/doc/html/_static/css/
apache-cassandra-3.11.2/doc/html/_static/fonts/
apache-cassandra-3.11.2/doc/html/_static/js/
apache-cassandra-3.11.2/doc/html/architecture/
apache-cassandra-3.11.2/doc/html/configuration/
apache-cassandra-3.11.2/doc/html/cql/
apache-cassandra-3.11.2/doc/html/data_modeling/
apache-cassandra-3.11.2/doc/html/development/
apache-cassandra-3.11.2/doc/html/faq/
apache-cassandra-3.11.2/doc/html/getting_started/
apache-cassandra-3.11.2/doc/html/operating/
apache-cassandra-3.11.2/doc/html/tools/
apache-cassandra-3.11.2/doc/html/troubleshooting/
apache-cassandra-3.11.2/interface/
apache-cassandra-3.11.2/javadoc/
apache-cassandra-3.11.2/javadoc/org/
apache-cassandra-3.11.2/javadoc/org/apache/
apache-cassandra-3.11.2/javadoc/org/apache/cassandra/
apache-cassandra-3.11.2/javadoc/org/apache/cassandra/auth/
apache-cassandra-3.11.2/javadoc/org/apache/cassandra/auth/class-use/
apache-cassandra-3.11.2/javadoc/org/apache/cassandra/auth/jmx/
```

Note, this will create a new directory (apache-cassandra-3.11.2) holding the Cassandra software release. Then enter this command to start Cassandra (lots of diagnostic messages will appear):

**apache-cassandra-3.11.2/bin/cassandra &**



Open a second terminal connection to the EMR master node. Going forward we will call this terminal connection: Cqlsh-Term. Enter the following into this terminal to start the command line interface csqlsh:

**apache-cassandra-3.11.2/bin/cqlsh**

```
aasth@LAPTOP-HJTR6HMR MINGW64 ~/Downloads
$ scp -i emr-key-pair.pem init.cql hadoop@ec2-52-86-39-223.compute-1.amazonaws.com:/home/hadoop
init.cql

aasth@LAPTOP-HJTR6HMR MINGW64 ~/Downloads
$ scp -i emr-key-pair.pem ex2.cql hadoop@ec2-52-86-39-223.compute-1.amazonaws.com:/home/hadoop
ex2.cql

aasth@LAPTOP-HJTR6HMR MINGW64 ~/Downloads
$ scp -i emr-key-pair.pem ex3.cql hadoop@ec2-52-86-39-223.compute-1.amazonaws.com:/home/hadoop
ex3.cql

aasth@LAPTOP-HJTR6HMR MINGW64 ~/Downloads
$ scp -i emr-key-pair.pem ex4.cql hadoop@ec2-52-86-39-223.compute-1.amazonaws.com:/home/hadoop
ex4.cql

aasth@LAPTOP-HJTR6HMR MINGW64 ~/Downloads
$ scp -i emr-key-pair.pem ex5.cql hadoop@ec2-52-86-39-223.compute-1.amazonaws.com:/home/hadoop
ex5.cql

aasth@LAPTOP-HJTR6HMR MINGW64 ~/Downloads
$ |
```

```
[hadoop@ip-172-31-59-222 ~]$ ls
apache-cassandra-3.11.2   apache-cassandra-3.11.2-bin.tar.gz
[hadoop@ip-172-31-59-222 ~]$ ls
apache-cassandra-3.11.2   apache-cassandra-3.11.2-bin.tar.gz   ex2.cql   ex3.cql   ex4.cql   ex5.cql   init.cql
[hadoop@ip-172-31-59-222 ~]$ vi init.cql
[hadoop@ip-172-31-59-222 ~]$ cat init.cql
CREATE KEYSPACE A20468022 WITH REPLICATION = { 'class' : 'SimpleStrategy','replication_factor' : 1 };
[hadoop@ip-172-31-59-222 ~]$ vi ex2.cql
```

Execute the below command:

USE A20468022;

source './ex2.cql'; DESCRIBE TABLE Music;

```
a20468022   system_schema   system_auth   system   system_distributed   system_traces

cqlsh> USE A20468022;
cqlsh:a20468022> source './ex2.cql';
cqlsh:a20468022> DESCRIBE TABLE Music;

CREATE TABLE a20468022.music (
    artistname text,
    albumname text,
    cost int,
    numbersold int,
    PRIMARY KEY (artistname, albumname)
) WITH CLUSTERING ORDER BY (albumname ASC)
    AND bloom_filter_fp_chance = 0.01
    AND caching = {'keys': 'ALL', 'rows_per_partition': 'NONE'}
    AND comment = ''
    AND compaction = {'class': 'org.apache.cassandra.db.compaction.SizeTieredCompactionStrategy', 'max_threshold': '32', 'min_threshold': '4'}
    AND compression = {'chunk_length_in_kb': '64', 'class': 'org.apache.cassandra.io.compress.LZ4Compressor'}
    AND crc_check_chance = 1.0
    AND dclocal_read_repair_chance = 0.1
    AND default_time_to_live = 0
    AND gc_grace_seconds = 864000
    AND max_index_interval = 2048
    AND memtable_flush_period_in_ms = 0
    AND min_index_interval = 128
    AND read_repair_chance = 0.0
    AND speculative_retry = '99PERCENTILE';
```

**Exercise 3) (3 points)**

a) Execute ex3.cql. Provide the content of this file as the result of this exercise.

```
[hadoop@ip-172-31-59-222 ~]$ vi ex3.cql
[hadoop@ip-172-31-59-222 ~]$ cat ex3.cql
INSERT INTO Music (artistName, albumName, numberSold, cost) VALUES ('Mozart', 'Greatest Hits', 100000,
10);
INSERT INTO Music (artistName, albumName, numberSold, cost) VALUES ('Taylor Swift','Fearless', 2300000,
15);
INSERT INTO Music (artistName, albumName, numberSold, cost) VALUES ('Black Sabbath', 'Paranoid',
534000, 12);
INSERT INTO Music (artistName, albumName, numberSold, cost) VALUES ('Katy Perry', 'Prism', 800000, 6);
INSERT INTO Music (artistName, albumName, numberSold, cost) VALUES ('Katy Perry', 'Teenage Dream',
750000, 14);
cqlsh:a20468022> source './ex3.cql';
cqlsh:a20468022> 'SELECT * FROM Music;'
Invalid syntax at line 1, char 1
```

```
cqlsh:a20468022> SELECT * FROM Music;

 artistname     | albumname     | cost | numbersold
----------------+---------------+------+------------
        Mozart | Greatest Hits |   10 |     100000
 Black Sabbath |      Paranoid |   12 |     534000
   Taylor Swift |      Fearless |   15 |    2300000
     Katy Perry |         Prism |    6 |     800000
     Katy Perry | Teenage Dream |   14 |     750000

(5 rows)
cqlsh:a20468022> source './ex4.cql';
```

**Exercise 4) (2 points)**

```
750000, 14);
[hadoop@ip-172-31-59-222 ~]$ vi ex4.cql
[hadoop@ip-172-31-59-222 ~]$ cat ex4.cql
SELECT * FROM Music where artistName = 'Katy Perry';
```

```
cqlsh:a20468022> source './ex4.cql';

 artistname | albumname     | cost | numbersold
------------+---------------+------+------------
 Katy Perry |         Prism |    6 |     800000
 Katy Perry | Teenage Dream |   14 |     750000

(2 rows)
```

**Exercise 5) (2 points)**

```
[hadoop@ip-172-31-59-222 ~]$ vi ex5.cql
[hadoop@ip-172-31-59-222 ~]$ cat ex5.cql
SELECT * FROM Music where numberSold >= 700000 ALLOW FILTERING;
```

```
cqlsh:a20468022> source './ex5.cql';

 artistname     | albumname     | cost | numbersold
----------------+---------------+------+------------
   Taylor Swift |      Fearless |   15 |    2300000
     Katy Perry |         Prism |    6 |     800000
     Katy Perry | Teenage Dream |   14 |     750000

(3 rows)
```

**Submitted By: -**

**Aastha Dhir**

**CWID-A20468022**

**adhir2@hawk.iit.edu**