

BIG DATA TECHNOLOGIES

ASSIGNMENT-12

Exercise 1

Answer:

Overview:

The article discusses classical data modeling, logical data and conceptual modeling, Cassandra data modeling, application workflow, query-driven mapping from a conceptual to a logical data model, and physical data modeling.

Conceptual data modeling and application workflow:

- When developing a Cassandra database schema, it is necessary to understand the data to be maintained as well as when a data-driven application would access it.
- The ER diagram is a representation of the previous Application workflow diagrams, that specify data entry patterns of application processes and represent both of these.

Cassandra Data Model:

- A CQL database is a collection of partitions having rows with comparable structures. The partition key is unique to each partition in a table, whereas a clustering key is unique to each row inside each partition.
- A table schema is a collection of columns that includes a primary key. A primary key is a combination of a partition key and a clustering key that uniquely identifies a database entry. Every column's data type is usually primitive (int, text, etc.), complex (set, list, or map), or counter.
- CQL, which has a syntax comparable to SQL, is used for expressing queries across tables. CQL does not handle binary operations such as joins and instead focuses on a set of query predicate statements rules to assure performance.

Mapping based on queries:

The four data modeling principles outlined below serve as a foundation for turning conceptual data models to logical data models.

DMP1 - Understanding the data, which is recorded using a conceptual data model, is the first stage in successful database design.

DMP2 - The second essential to a successful database design is to understand your queries recorded by an application process.

DMP3 - The third key to a successful database design is data nesting.

DMP4 - The fourth key to a successful database design is data duplication.

Mapping Patterns are utilized for automating Cassandra database schema designing.

Physical Data Modelling: The final phase is to evaluate and improve the data model's logic in a physical data model.

Mapping Rule:

The five mapping rules listed below assist a query-driven move from a conceptual one to a logical data model.

MR1 - Entities and relationships are mapped to table rows in MR1, while entity & relationship types are mapped to tables.

MR5 - The primary key fields are associated with key attributes.

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

```
Last login: Sun Nov 26 17:34:45 2023 from 287.237.235.14

#
_ _ _ _ _      Amazon Linux 2
 \_ _ _ _ _    \
 ~~~~\#####\ AL2 End of Life is 2025-06-30.
       \|###|
        V--'----->
               A newer version of Amazon Linux is available!

~~~~~\_____/_____
     _m/_

Amazon Linux 2023, GA and supported until 2028-03-15.
https://aws.amazon.com/linux/amazon-linux-2023/
```

```

EEEEEEEEEEEEEEEEEEEE MMMMMM MMMMMM RRRRRRRRRRRR
E:::EEEEEEEEEEEEEEEE M:::M:::M M:::M:::M R:::R:::R:::R:::R
E:::EEEEEEEEEEEEEEEE EEEEE M:::M:::M M:::M:::M R:::R:::R:::R:::R
E:::E M:::M:::M M:::M:::M M:::M:::M R:::R:::R:::R:::R
E:::E M:::M:::M M:::M:::M M:::M:::M R:::R:::R:::R:::R
E:::EEEEEEEEEEEE M:::M:::M M:::M:::M M:::M:::M R:::R:::R:::R:::R
E:::EEEEEEEEEEEE M:::M:::M M:::M:::M M:::M:::M R:::R:::R:::R:::R
E:::E M:::M:::M M:::M:::M M:::M:::M R:::R:::R:::R:::R
E:::E EEEEE M:::M:::M M M M M:::M R:::R:::R:::R:::R
E:::EEEEEEEEEEEEEEEE M:::M:::M M:::M:::M R:::R:::R:::R:::R
EEEEEEEEEEEEEEEEEEEE MMMMMM MMMMMM RRRRRR RRRRRR

```

```
[mc2-user@ipw-172-31-12-156 ~]$ wget https://archive.apache.org/dist/cassandra/3.11.2/apache-cassandra-3.11.2-bin.tar.gz
--2023-11-26 17:39:38-- https://archive.apache.org/dist/cassandra/3.11.2/apache-cassandra-3.11.2-bin.tar.gz
Resolving archive.apache.org (archive.apache.org)... 65.108.294.189, 2a01:493::1:a8b84::2
Connecting to archive.apache.org (archive.apache.org)... 65.108.294.189:443... connected.
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 16.9MB/s in 2.2s
2023-11-26 17:39:40 [16.9 MB/s] - 'apache-cassandra-3.11.2-bin.tar.gz' saved [38436262/38436262]
```

```
lxc-userip-172-31-12-156 -i$ tar -xvzf 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/javadocs/
  apache-cassandra-3.11.2/javadoc/org/
  apache-cassandra-3.11.2/javadoc/org/apache/
```

```
[ec2-user@ip-172-31-12-156 ~]$ java apache-cassandra-3.11.2/bin/cassandra &
[1] 6214
[ec2-user@ip-172-31-12-156 ~]$ OpenJDK 64-bit Server VM warning: Cannot open file apache-cassandra-3.11.2/bin/../logs/gc.log due to No such file or directory

CompilerOracle: dontinline org/apache/cassandra/db/ColumnsSerializer.deserializeLargeSubset (Lorg/apache/cassandra/io/util/DataInputPlus;Lorg/apache/cassandra/db/Columns;)I;Lorg/apache/cassandra/db/Columns;
CompilerOracle: dontinline org/apache/cassandra/db/ColumnsSerializer.serializeLargeSubset (Ljava/lang/Object;Lorg/apache/collection/Lorg/apache/cassandra/db/Columns;Lorg/apache/cassandra/io/util/DataOutputPlus;I)V
CompilerOracle: dontinline org/apache/cassandra/db/CMSerializer.serializeLargeSubset (Ljava/lang/Object;Lorg/apache/collection/Lorg/apache/cassandra/db/Columns;Lorg/apache/cassandra/io/util/DataOutputPlus;I)V
CompilerOracle: dontinline org/apache/cassandra/db/commitlog/AbstractCommitLogSegmentManager.advanceAllocatingFrom (Lorg/apache/cassandra/db/commitlog/CommitLogSegment;)V
CompilerOracle: dontinline org/apache/cassandra/db/transform/BaseIterator.tryGetMoreContents ()Z
CompilerOracle: dontinline org/apache/cassandra/db/transform/StoppingTransformation.stop ()V
CompilerOracle: dontinline org/apache/cassandra/db/transform/StoppingTransformation.stopPartition ()V
CompilerOracle: dontinline org/apache/cassandra/io/util/BufferedDataOutputStream.flush ()V
CompilerOracle: dontinline org/apache/cassandra/io/util/BufferedDataOutputStream.writeSlow ()V
CompilerOracle: dontinline org/apache/cassandra/io/util/BufferedDataOutputStream.writeSlowSlow ()V
CompilerOracle: dontinline org/apache/cassandra/io/util/BufferingInputStream.readPrimitiveSlowly ()I
CompilerOracle: inline org/apache/cassandra/db/rows/UnfilteredSerializer.serializeRowBody (Lorg/apache/cassandra/db/rows/Row;Lorg/apache/cassandra/db/SerializationHeader;Lorg/apache/cassandra/io/util/DataOutputPlus;)V
CompilerOracle: inline org/apache/cassandra/io/util/Memory.checkBounds ()JJV
CompilerOracle: inline org/apache/cassandra/io/util/SafeMemory.checkBounds ()JJIV
CompilerOracle: inline org/apache/cassandra/utils/AsymmetricOrdering.selectBoundary (Lorg/apache/cassandra/utils/AsymmetricOrdering;O;II)I
CompilerOracle: inline org/apache/cassandra/utils/AsymmetricOrdering.strictnessOfLessThan (Lorg/apache/cassandra/utils/AsymmetricOrdering;O;I)I
CompilerOracle: inline org/apache/cassandra/utils/BloomFilter.indexes (Lorg/apache/cassandra/utils/FILTER/FilterKey;)[I
CompilerOracle: inline org/apache/cassandra/utils/BloomFilter.setIndexes (JJIJ)[I
CompilerOracle: inline org/apache/cassandra/utils/ByteBufferUtil.compare (Ljava/nio/ByteBuffer;Ljava/nio/ByteBuffer;)I
CompilerOracle: inline org/apache/cassandra/utils/ByteBufferUtil.compare ((Ljava/nio/ByteBuffer;)I
CompilerOracle: inline org/apache/cassandra/utils/ByteBufferUtil.compareUnsigned (Ljava/nio/ByteBuffer;Ljava/nio/ByteBuffer;)I
CompilerOracle: inline org/apache/cassandra/utils/FastByteOperations$UnsafeOperations.compareTo (Ljava/lang/Object;JLJava/lang/Object;)II
CompilerOracle: inline org/apache/cassandra/utils/FastByteOperations$UnsafeOperations.compareTo (Ljava/lang/Object;JLJava/lang/Object;)J
CompilerOracle: inline org/apache/cassandra/utils/FastByteOperations$UnsafeOperations.compareTo (Ljava/lang/Object;JLJava/lang/Object;)I
CompilerOracle: inline org/apache/cassandra/vint/VIntDecoder.encodeVint ()II][8
INFO [main] 2023-11-26 17:41:05.080 VMConfigurationLoader.java:89 - Configuration location: file:/home/ec2-user/apache-cassandra-3.11.2/conf/cassandra.yaml
INFO [main] 2023-11-26 17:41:05.465 Config.java:495 - Node configuration:[allocate_tokens_for_keyspace=null; authenticator=AllowAllAuthenticator; authorizer=AllowAllAuthorizer; auto_bootstrap=true; auto_snapshot=true; backpressure_enabled=false; batch_pressure_strategy=org.apache.cassandra.net.RateBasedBackPressure{threshold_i=0.9, factor=5, flowFAST}; batch_size_min_flushhold_in_Kb=50; batch_size_warn_threshold_in_Kb=5; batchlog_replay_throttle_in_Kb=1024; broadcast_address=null; broadcast_rpc_address=null; buffer_pool_use_heap_if_exhausted=true; cache_contention_timeout_in_ms=1000; cdc_enabled=false; cdc_free_space_in_MB=16; cdc_log_group_name=cdc; client_encryption_options={dynamic_snitch:null; snitch=null; encryption_options=null; dynamic_snitch_update_interval_in_ms=1000; enable_materialized_views=true; failure_policy=stop; commitlog_compression=null; commitlog_directory=null; commitlog_max_compression_buffers_in_pool=3; commitlog_periodic_queue_size=1; commitlog_segment_size_in_Mb=32; commitlog_sync_period_in_ms=NAN; commitlog_sync_batch_window_in_ms=NAN; commitlog_sync_period_in_ms=10000; commitlog_total_size_in_Mb=null; compaction_large_partition_warning_threshold_mb=100; compaction_throughput_mb_per_sec=6; concurrent_compactors=null; concurrent_counter_writes=32; concurrent_materialized_view_writes=32; current_reads=32; current_replicates=null; current_writes=32; counter_cache_keys_to_save=2147483647; counter_cache_size=7200; counter_cache_size_in_Mb=null; counter_write_request_timeout_in_ms=5000; credentials_cache_max_entries=1000; credentials_update_interval_in_ms=-1; credentials_validity_in_ms=2000; cross_node_timeout=5; data_file_directories=java.lang.String@f00b7; disk_access_mode=auto; disk_failure_policy=stop; disk_optimization_estimate_percentile=95; disk_optimization_page_cross_chance=0.1; disable_dynamic_snitch=true; dynamic_snitch_load_ratio=1.0; dynamic_snitch_weighted_vnode_load=1.0; dynamic_snitch_warning_threshold=1.0; dynamic_snitch_update_interval_in_ms=1000; enable_materialized_views=true; enable_scripted_user_defined_functions=false; enable_user_defined_functions=true; enable_user_defined_functions_threads=true; encryption_options=null; endpoint_snitch=SimpleSnitch; file_cache_round_up=null; file_cache_size_in_Mb=null; gc_log_threshold_in_ms=200; gc_warn_threshold_in_ms=1000; hinted_handoff_disabled_datacenters=[]; hinted_handoff_enabled=true; hinted_handoff_throttle_in_Kb=1024; hints_compression_algorithm=com.lmax.disruptor.atomic.boundedblockingqueue; hints_index_in_ms=10000; incremental_backups=false; index_interval=null; index_summary_capacity_in_Mb=null; index_summary_resize_interval_in_minutes=60; initial_token=null; inter_dc_stream_throughput_outbound_megabits_per_sec=200; inter_dc_tcp_nodelay=false; internode_authenticator=null; internode_compression=snappy; internode_recv_buff_size_in_bytes=65536; internode_send_buff_size_in_bytes=65536; ip_tos=0x00; jmx_remote_port=7199; jmx_server_host=null; listen_interface=null; listen_interface_prefer_ipv6=false; listen_on_broadcast_address=false; max_hints_in_window_in_ms=10000000; max_hints_delivery_threads=2; max_hints_file_size_in_Mb=128; maxmutation_size_in_Kb=null; max_streaming_retries=3; max_value_size_in_Mb=256; memtable_allocation_type=heap_buffers; memtable_cleanup_threshold=null; memtable_flush_writers=0; memtable_heap_space_in_Mb=null; memtable_offheap_space_in_Mb=null; min_free_space_per_drive_in_Mb=50; native_transport_max_concurrent_connections=1; native_transport_max_concurrent_connections_per_ip=1; native_transport_max_frame_size_in_Mb=256; native_transport_max_threads=128; native_transport_port=9042; native_transport_port_ssl=null; num_tokens=256; otc_background_expiration_interval_ms=200; otc_coalescing_enough_coalesced_messages=8; otc_coalescing_strategy=DISABLED; otc_coalescing_window_us=200; partitioner=org.apache.cassandra.dht.Murmur3Partitioner; permissions_cache_max_entries=1000; permissions_update_interval_in_ms=-1; permissions_validation_interval_in_ms=2000; phi_convict_threshold=10; prepared_statements_cache_size_in_Mb=null; range_request_timeout_in_ms=10000; read_repair_delay_in_ms=1000; request_scheduler=null; request_scheduler_id=null; request_scheduler_options=null; request_timeout_in_ms=10000; role_manager=org.apache.cassandra.role.Manager; roles_cache_max_entries=1000; roles_update_interval_in_ms=-1; roles_validation_interval_in_ms=2000; row_cache_class_name=org.apache.cassandra.cache.OHCProvider; row_cache_keys_to_save=2147483647;
```

```
Last login: Sun Nov 26 11:37:19 on tty008
(base) sailavanyanarthu@Sailavanyas-Air ~ % ssh -i keyPair.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com

Warning: Identity file keyPair.pem not accessible: No such file or directory.
ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com: Permission denied (publickey,gssapi-keyex,gssapi-with-mic).
(base) sailavanyanarthu@Sailavanyas-Air ~ % cd downloads
(base) sailavanyanarthu@Sailavanyas-Air downloads % chmod 400 keyPair.pem
(base) sailavanyanarthu@Sailavanyas-Air downloads % ssh -i keyPair.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com

Last login: Sun Nov 26 17:37:51 2023 from 207.237.235.14

      #_
     _#_
    ~\  #####
   ~ ~ \_#####\
        \###/
       ~#\#
      ~ ~ V--'  -->
           /
          /
         /
        /_/_/
       /_/_/
      /_/_/

Amazon Linux 2

AL2 End of Life is 2025-06-30.

A newer version of Amazon Linux is available!

Amazon Linux 2023, GA and supported until 2028-03-15.
https://aws.amazon.com/linux/amazon-linux-2023/
```

Step C

apache-cassandra-3.11.2/bin/cqlsh

```
[[ec2-user@ip-172-31-12-156 ~]$ apache-cassandra-3.11.2/bin/cqlsh
Connected to Test Cluster at 127.0.0.1:9042.
[cqlsh 5.0.1 | Cassandra 3.11.2 | CQL spec 3.4.4 | Native protocol v4]
Use HELP for help.
```

Step D

a)

```
[[ec2-user@ip-172-31-12-156 ~]$ vi init.cql
[[ec2-user@ip-172-31-12-156 ~]$ cat init.cql
CREATE KEYSPACE A20516764 WITH REPLICATION = { 'class' : 'SimpleStrategy', 'replication_factor' : 1 };
[[ec2-user@ip-172-31-12-156 ~]$
```

b) and c)

```
cqlsh> source './init.cql';
cqlsh> describe keyspaces;
```

```
system_schema  system_auth  system  system_distributed  system_traces  a20516764
```

d)

```
[cqlsh> USE A20516764;
```

```
CREATE TABLE A20516764.Music (
  artistName text,
  albumName text,
  numberSold int,
  Cost int,
  PRIMARY KEY (artistName, albumName))
WITH CLUSTERING ORDER BY (albumName DESC);

-- INSERT --
```

7,43

A11

```
[[ec2-user@ip-172-31-12-156 ~]$ vi ex2.cql
```

This is the EMR Connection for Third Terminal

[illegible]

```

EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRRRRRRRRRR
E:EEEEEEEEEEEEEEEE: M:MMMM:M M:MMMM:M R:RRRRRRRRRRRR:
EE:EEEEEEEEEEEEEEEE:E M:MMMM:M M:MMMM:M R:RRRRRRRRRRRR:
E:EE:E EEEEE M:MMMM:M M:MMMM:M RR:RRR R:RR:R
E:EE:E M:MMMM:M:M M:MMMM:M R:RR R:RR:R
E:EEEEEEEEEEEEEE M:MM:M M:M:M M:M:M R:RRRRRRRRRRRR:
E:EEEEEEEEEEEEEE M:MM:M M:M:M:M M:M:M R:RRRRRRRRRRRR:
E:EEEEEEEEEEEEEE M:MM:M M:M:M M:M:M R:RRRRRRRRRRRR:
E:EEEEEEEEEEEEEE M:M:M M:M:M M:M:M R:RR R:RR:R
E:EE:E EEEEE M:MM:M M:M:M R:RR:R
EE:EEEEEEEEEEEEEE:E M:MM:M M:M:M R:RR:R R:RR:R
EEEEEEEEEEEEEEEEEE:E M:MM:M M:M:M RRR:RRR
EEEEEEEEEEEEEEEEEEEE MMMMMMMM MMMMMMMM RRRRRRRR RRRRRRR

```

Exercise 3

a)

```

--yPair.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com ... --air.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com ... --r.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com +
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, 16);
INSERT INTO Music (artistName, albumName, numberSold, Cost) values ('Katy Perry', 'Teenage Dream', 750000, 14);
-- INSERT --

```

```
[[ec2-user@ip-172-31-12-156 ~]$ vi ex3.cql
```

b)

```
[cqlsh:a20516764> source './ex3.cql';  
[cqlsh:a20516764> SELECT * FROM Music;
```

artistname	albumname	cost	numbersold
Taylor Swift	Fearless	15	2300000
MOZART	Greatest Hits	10	100000
Black Sabbath	Paranoid	12	534000
Katy Perry	Teenage Dream	14	750000
Katy Perry	Prism	16	800000

(5 rows)

Exercise 4

```
...yPair.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com ...  
...air.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com  
...r.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com +  
SELECT * FROM Music WHERE artistName in ('Katy Perry');  
~  
~  
~  
~  
~  
~  
~
```

```
[[ec2-user@ip-172-31-12-156 ~]$ vi ex4.cql
```

```
[cqlsh:a20516764> source './ex4.cql';
```

artistname	albumname	cost	numbersold
Katy Perry	Teenage Dream	14	750000
Katy Perry	Prism	16	800000

(2 rows)

Exercise 5

```
...yPair.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com ...  
...air.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com  
...r.pem ec2-user@ec2-3-144-12-121.us-east-2.compute.amazonaws.com +  
SELECT * FROM Music WHERE numberSold >= 700000 allow filtering;  
~  
~  
~  
~  
~  
~  
~
```

```
[[ec2-user@ip-172-31-12-156 ~]$ vi ex5.cql
```

```
[cqlsh:a20516764> source './ex5.cql';
```

artistname	albumname	cost	numbersold
Taylor Swift	Fearless	15	2300000
Katy Perry	Teenage Dream	14	750000
Katy Perry	Prism	16	800000

```
(3 rows)
```

```
cqlsh:a20516764> █
```

Submitted by:
Sailavanya Narthu
A20516764
snarthu@hawk.iit.edu