# Tutorial 7

# Apache Cassandra

To get started with Cassandra NoSQL database, we will step through a single-node, local installation on VM.

**1)** The following points are the requirements to run Cassandra locally: Move to the Apache download site for the Cassandra project (http://cassandra.apache.org/download/), choose 3.11.16, and select a mirror to download the latest version of Cassandra. When complete, copy the .tar or .gzip file to a location that your user has read and write permissions for. This example will assume that this is going to be the ~/Downloads/ directory on ubuntu VM.

**Note:** If you could not understand from the command instructions, then please check the screenshot for better understanding. If you like to explore the details of each command along with examples, then check the documentation of Cassandara on website: https://cassandra.apache.org/doc/latest/

Download Apache Cassandra from the following link as mentioned below

**https://www.apache.org/dyn/closer.lua/cassandra/3.11.16/apache-cassandra-3.11.16-bin.tar.gz**

**2) $cd Downloads**

Follow the instructions to unzip on the below screenshots and change the name of the folder as you did during the Hadoop, HBase and Spark installations (Tutorials, 2, 4, 6).

**3)** Configuration: At this point, you could start your node with no further configuration. However, it is good to get into the habit of checking and adjusting the properties that are indicated as follows using instructions as shown in the screenshot in step no. 2.

```
$cd /usr/local
$cd cassandra
$cd conf
$nano cassandra.yaml
```
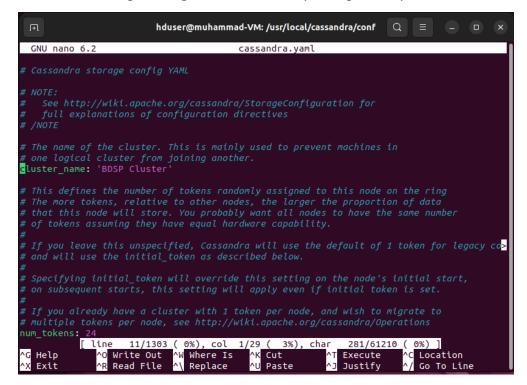
It is usually a good idea to rename your cluster. Inside the conf/cassandra.yaml file, specify a new cluster_name property, overwriting the default Test Cluster as shown below in screenshot:

```
cluster_name: 'BDSP Cluster'
```

The num_tokens property default of 256 has proven to be too high for the newer, 3.x versions of Cassandra. Go ahead and set that to 24:

```
num_tokens: 24
```

save the file using nano/ gedit editor after updating above parameters.



*Press Alt+C to display the line number using nano editor.*

**4)** By default, Cassandra will come up bound to localhost or 127.0.0.1. For your own local development machine
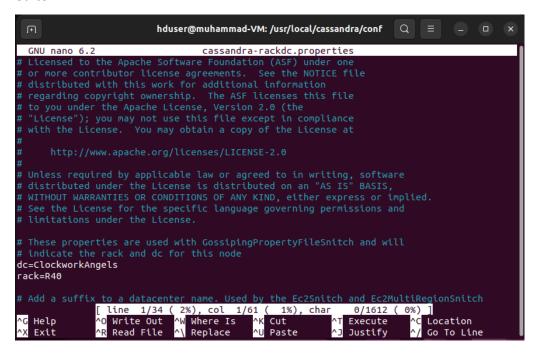


```
cassandra-rackdc.properties
```

In terms of NoSQL databases, Apache Cassandra handles multi-data center awareness better than any other. To configure this, each node must use **GossipingPropertyFileSnitch** (as previously mentioned in the preceding **cassandra.yaml** configuration process) and must

have its local data center (and rack) settings defined. Therefore, we set the dc and rack properties in the **conf/cassandra-rackdc.properties** file:

```
dc=ClockworkAngels
rack=R40
```

If these properties are already set as shown below in the screenshot. You can exit from the editor.



**5)** Starting Cassandra: To start Cassandra locally, execute the Cassandra script. If no arguments are passed, it will run in the foreground. To have it run in the background, send the **-p flag** with a destination file for the **Process ID (PID)**:



```
$bin/cassandra -p cassandra.pid
```

Or

```
$bin/cassandra -f
```

When the Cassandra started, then leave this terminal as shown below

Let this window stay as is for the moment, don't close it. However, if you would like to stop Cassandra, press Ctrl + C to stop the Cassandra process or you can use the command to kill the process as mentioned below on the screen. **You must perform these steps if the Cassandra showed an error in the case of start-up using this command (`bin/cassandra -p cassandra.pid` or `bin/cassandra -f`).**



This will store the PID of the Cassandra process in a file named `cassandra.pid` in the `local/cassandra` directory. Several messages will be dumped to the screen.

Open a new terminal by pressing **ctrl+ Alt + t**, Check the status of Cassandra by using the following commands as mentioned below in the screenshot

**6) Cassandra Cluster Manager:** If you want an even faster way to install Cassandra, you can use an open-source tool called **CCM**. **CCM** installs Cassandra for you, with very minimal configuration. In addition to ease of installation, CCM also allows you to run multiple Cassandra nodes locally.

Open a new terminal by pressing (ctrl + Alt + t). Install python 2.7 before execution of next commands.

```
$sudo apt install python2
```

Press Y for the installation.



First, let's clone the CCM repository from GitHub, and cd into the directory:

```
$cd        <- Press enter Key
```

```
$sudo apt install git
```

```
$git clone https://github.com/riptano/ccm.git
```

```
$cd ccm
```

Next, we will run the setup program to install CCM:

```
$nano ./setup.py
```

And change the first line word "python" to "python2"as mentioned below



```
$sudo ./setup.py install
```



## A quick introduction to the data model

Now that we have a Cassandra cluster running on our local machine, we will demonstrate its use with some quick examples. We will start with **cqlsh**, and use that as our primary means of working with the Cassandra data model.

7) Using Cassandra with **cqlsh**: To start working with Cassandra, let's start the **Cassandra Query Language** (**CQL**) shell. The shell interface will allow us to execute CQL commands to define, query, and modify our data. As this is a new cluster and we have turned on authentication and authorization, we will use the default cassandra and cassandra username and password, as follows:

```
$cd /usr/local/cassandra/
```

**`$bin/cqlsh`**



One terminal showed that the Cassandra is running, and you can execute the Cassandra database commands on the other terminal.

**`cassandra@cqlsh> describe cluster;`**



**`cqlsh>DESCRIBE KEYSPACES;`**

Check all the tables that are defined in the keyspace.

**`cqlsh>DESCRIBE KEYSPACE system;`**

**`cqlsh>CREATE KEYSPACE vehicle_tracker WITH REPLICATION = { 'class' : 'SimpleStrategy', 'replication_factor' : 1 };`**

**`cqlsh>DESCRIBE KEYSPACES;`**

Check the screenshot of this command on the next page of the tutorial. If you like to drop the keyspace

**`cqlsh>DROP KEYSPACE vehicle_tracker;`**

If you would like to know the details of the commands, please check the website: https://cassandra.apache.org/doc/latest/cassandra/developing/cql/ddl.html

```
cqlsh>USE home_Security;
```

Follow the screenshot to create the Table in the collection **'home_security'**



Create another table named as 'activity' inside the collection 'home_security' and the screenshots are mentioned below

```
hduser@muhammad-vm: /usr/local/cassandra

cqlsh:home_security> CREATE TABLE activity (home_id text, datetime timestamp, code_used text, event text, PRIMARY KEY (datetime));
cqlsh:home_security> INSERT INTO activity (home_id, datetime, code_used, event) VALUES ('H01474777', '2014-05-21 07:32:16', '5999'
, 'alarm set');
cqlsh:home_security> SELECT * FROM activity;

 datetime                        | code_used | event     | home_id
---------------------------------+-----------+-----------+-----------
 2014-05-21 06:32:16.000000+0000 |      5999 | alarm set | H01474777

(1 rows)
cqlsh:home_security>
```

8) Copy the data from **csv** file. Download the file **'events.csv'** and **'homes.csv'** from Moodle in the **'Downloads'** folder on VM (This is not Hadoop and it is your local Ubuntu machine) and write the command as mentioned in the screenshot.

```
hduser@muhammad-vm: /usr/local/cassandra

cqlsh:home_security> copy activity (home_id, datetime, code_used, event) FROM '/home/hduser/Downloads/events.csv' WITH header
 = true AND delimiter = '|';
Using 1 child processes

Starting copy of home_security.activity with columns [home_id, datetime, code_used, event].
Processed: 32 rows; Rate:      56 rows/s; Avg. rate:      82 rows/s
32 rows imported from 1 files in 0.389 seconds (0 skipped).
cqlsh:home_security> SELECT * FROM activity;

 datetime                        | code_used              | event     | home_id
---------------------------------+------------------------+-----------+-----------
 2014-05-22 11:44:07.000000+0000 | alarm reset by office  |      null | H01474777
 2014-05-23 18:06:58.000000+0000 |       alarm turned off |      1566 | H02257222
 2014-05-23 08:28:16.000000+0000 |              alarm set |      8889 | H01545551
 2014-05-21 07:32:16.000000+0000 |              alarm set |      5599 | H01474777
 2014-05-22 19:10:56.000000+0000 |       alarm turned off |      1245 | H00999943
 2014-05-22 11:23:59.000000+0000 |         alarm breached |      null | H01474777
 2014-05-22 07:45:28.000000+0000 |              alarm set |      2121 | H01033638
 2014-05-22 17:22:15.000000+0000 |       alarm turned off |      5599 | H01474777
 2014-05-21 13:02:11.000000+0000 |       alarm turned off |      1919 | H01033638
 2014-05-23 08:52:19.000000+0000 |              alarm set |      1245 | H00999943
 2014-05-22 21:59:44.000000+0000 |       alarm turned off |      1566 | H02257222
 2014-05-22 11:25:00.000000+0000 |          police called |      null | H01474777
 2014-05-21 09:05:54.000000+0000 |              alarm set |      1245 | H00999943
 2014-05-23 07:44:23.000000+0000 |              alarm set |      5599 | H01474777
 2014-05-21 19:03:33.000000+0000 |       alarm turned off |      1245 | H00999943
 2014-05-21 18:41:02.000000+0000 |       alarm turned off |      8889 | H01545551
 2014-05-23 07:49:36.000000+0000 |              alarm set |      1566 | H02257222
 2014-05-21 18:30:33.000000+0000 |       alarm turned off |      5599 | H01474777
 2014-05-21 16:58:39.000000+0000 |              alarm set |      1919 | H01033638
 2014-05-21 07:50:43.000000+0000 |       alarm turned off |      2121 | H01033638
 2014-05-23 18:56:23.000000+0000 |       alarm turned off |      1245 | H00999943
 2014-05-22 07:44:13.000000+0000 |              alarm set |      5599 | H01474777
 2014-05-23 18:14:53.000000+0000 |              alarm set |      8889 | H01545551
 2014-05-21 07:55:58.000000+0000 |              alarm set |      2121 | H01033638
 2014-05-21 06:32:16.000000+0000 |                   5999 | alarm set | H01474777
 2014-05-22 08:55:10.000000+0000 |              alarm set |      1245 | H00999943
 2014-05-21 08:30:14.000000+0000 |              alarm set |      8889 | H01545551
 2014-05-23 18:28:41.000000+0000 |       alarm turned off |      5599 | H01474777
 2014-05-21 19:01:46.000000+0000 |       alarm turned off |      2121 | H01033638
 2014-05-21 05:29:47.000000+0000 |              alarm set |      1566 | H02257222
 2014-05-21 07:50:22.000000+0000 |              alarm set |      2121 | H01033638
 2014-05-22 08:32:22.000000+0000 |              alarm set |      8889 | H01545551
 2014-05-22 18:35:29.000000+0000 |       alarm turned off |      8889 | H01545551

(33 rows)
cqlsh:home_security>
```
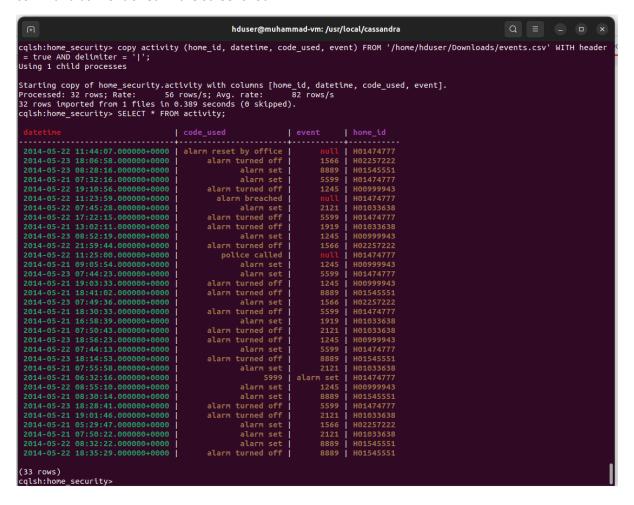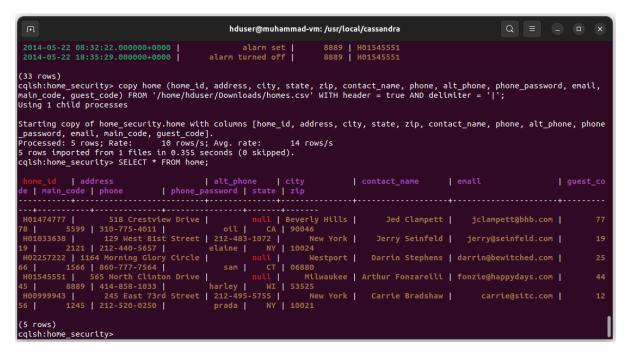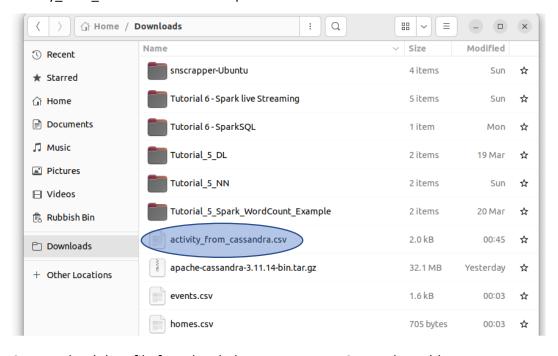
9) Export the data from the Cassandra table to 'csv' file on your local Ubuntu machine.



The output file will be stored in 'Downloads' folder as shown below on Ubuntu VM. You might see some other files in the Download folder than this screenshot. Make sure that activity_from_cassandra.csv must be present.



Steps to load data file from local ubuntu system to Cassandra Table

1) Create a file named as "employees_data.csv" and insert the records as mentioned below

**$nano employees_data.csv**

employee_id,firstname,lastname,department,city

1,Peter,Mark,Engineering,Dublin

2,Sean,Kelly,Physics,Dublin

3,Derek,Monahan,IT,Galway

4,Miles,Turner,Medical,Cork

5,Sarah,Hayes,Nursing,Cork

Or download the file "employees_data.csv" from Moodle.

2) Follow the sequence of commands to load data into Table "employees_data" and Keyspace named as "employees".

**$cd /usr/local/cassandra**

**hduser@muhammad-VM:/usr/local/cassandra$ bin/cqlsh**

**Connected to BDSP Cluster at 127.0.0.1:9042.**

**cqlsh> CREATE KEYSPACE employees WITH replication = {'class':' SimpleStrategy', 'replication_factor' : 1};**

**cqlsh> CREATE TABLE employees.employees_data (employee_id int PRIMARY KEY, firstname text, lastname text, department text, city text);**

**cqlsh> USE employees;**

**cqlsh:employees> COPY employees.employees_data (employee_id, firstname, lastname, department, city) FROM '/home/hduser/Downloads/ employees_data.csv' WITH HEADER = true;**

**cqlsh:employees> select * from employees.employees_data;**

3) The screenshot showed the sequence of commands as shown below.

4) Learn from the book reference provided in references for further understanding of Cassandra query language and perform queries on the datasets of your choice.


## References:

- https://cassandra.apache.org/doc/latest/cassandra/cql/ddl.html

- Cassandra: The Definitive Guide, (Revised) Third Edition, 3rd Edition, Jeff Carpenter, Eben Hewitt, O'Reilly Media, Inc., January 2022.

- Installation instructions: https://www.cloudduggu.com/cassandra/installation/