# Measuring Cassandra Cluster

Rohit Mehra

*University of Colorado, Boulder*

**Abstract**

Apache Cassandra is a highly scalable, distributed, high performance database that falls under the NoSQL category. It is designed to handle large amounts of data across many commodity servers, providing high availability with no single point of failure.Due to replication and clustering feature, it is highly scalable and its performance increases linearly as we add the nodes.

The aim of the project is to get familiar with Docker and on using NTP and setting Cassandra cluster. The project focuses on getting measurements in terms of read-write latencies and ping latencies across Cassandra Clusters setup in multiple topologies and networks. We have implemented 4 topologies which is elaborated in the sections below.

*Keywords:*

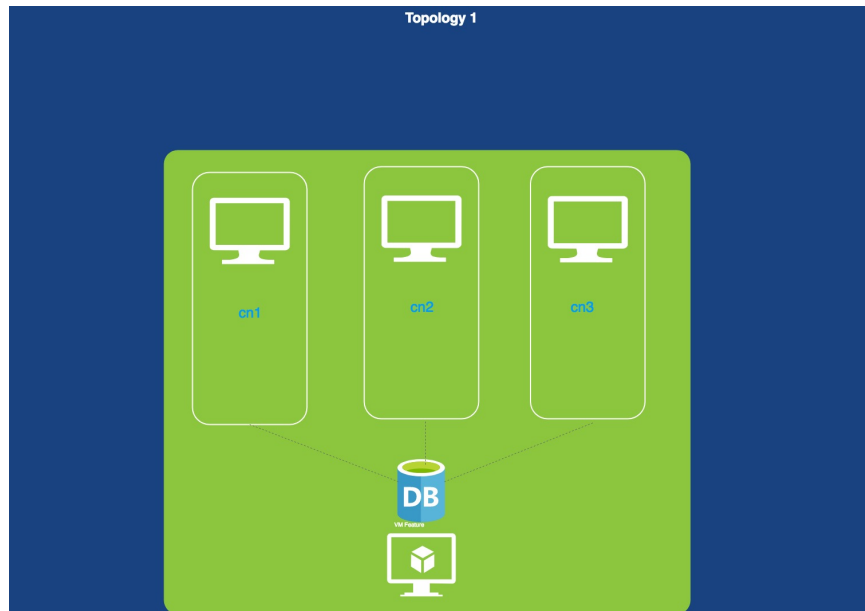Cassandra, NoSQL, Docker, VirtualMachine

## 1. The First Topology



Figure 1: Topology 1

In the topology, we had a single host system running three different docker containers. Each container had a Cassandra node in it and was called cn1, cn2 and cn3.Also, we booted up 1 machine which was used as server to measure clock skew. The steps we followed were as follows.

1. Dell XPS was used as host system with 16GB RAM

2. The next step was to spring up the Docker instances. The instructions for which could be found in the docker documentation.

3. The third step was to install NTP and Cassandra.

4. The last step was to establish a cluster using all the client nodes.

5. Cassandra stress command was used to measure read, write and mix latencies among a cluster node.

6. One client and server combination was used to test clock skew.

The figure describing the topology is $Figure 1$.

### 1.1. Testing Topology 1

1. All the docker containers are assigned IP from the docker subnet.

2. We used public NTP server IP in the server node so that server can be in sync with a public server.

3. NTP.conf file of all the clients points to the server IP so that it can be in sync.

4. ntpstat can be used to check the synchronization state of the client and the server.

5. The reason we needed to use NTP to configure the clock skew is as follows: Computer's use Cesium 133 for clock synchronization.

As we go down the tree as shown in $Figure 2$, we start to get skew between the root node and the nodes along the tree. To ensure that all the VMs/Containers follow the same time, we use NTP for Synchronization.

The NTP conf can be configured at \etc\ntp.conf. Once the conf files are set, the servers are rebooted.The sync status can be seen using the $ntpstat$ command.
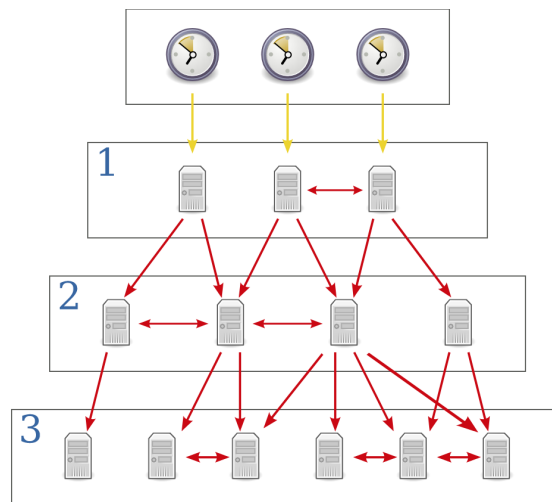


Figure 2: NTP Hierarchy

6. Now we wanted to check if the key spaces are shared between the Cassandra ring.

7. We used the $cassandra - stress$ tool to get the read, write and mixed latencies of the same.

8. As the number of iterations increased, the read and write latencies started to gradually increase as well.
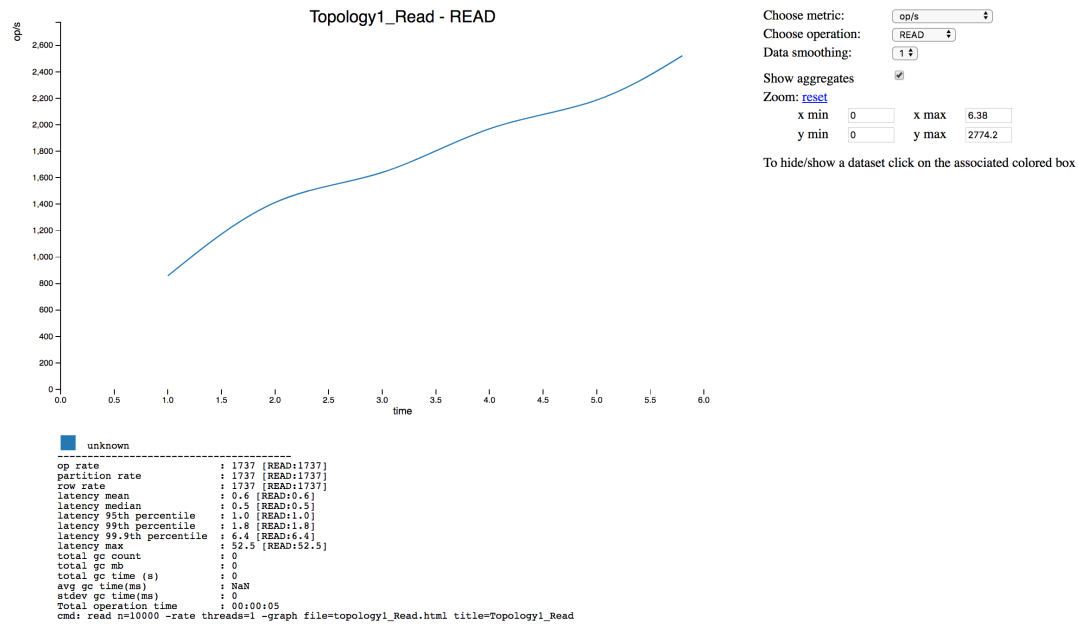
9. This can be seen in $Figure3$, $Figure4$, $Figure5$



**Topology1_Read - READ**

Choose metric: op/s
Choose operation: READ
Data smoothing: 1

Show aggregates ✓
Zoom: reset

| | | | |
|---|---|---|---|
| x min | 0 | x max | 6.38 |
| y min | 0 | y max | 2774.2 |

To hide/show a dataset click on the associated colored box

■ unknown
----------------------------------------
```
op rate                   : 1737 [READ:1737]
partition rate            : 1737 [READ:1737]
row rate                  : 1737 [READ:1737]
latency mean              : 0.6 [READ:0.6]
latency median            : 0.5 [READ:0.5]
latency 95th percentile   : 1.0 [READ:1.0]
latency 99th percentile   : 1.8 [READ:1.8]
latency 99.9th percentile : 6.4 [READ:6.4]
latency max               : 52.5 [READ:52.5]
total gc count            : 0
total gc mb               : 0
total gc time (s)         : 0
avg gc time(ms)           : NaN
stdev gc time(ms)         : 0
Total operation time      : 00:00:05
cmd: read n=10000 -rate threads=1 -graph file=topology1_Read.html title=Topology1_Read
```

Figure 3: Cassandra Stress Topology 1 Read



**Topology1_Write - WRITE**

Choose metric: op/s
Choose operation: WRITE
Data smoothing: 1

Show aggregates ✓
Zoom: reset

| | | | |
|---|---|---|---|
| x min | 0 | x max | 6.38 |
| y min | 0 | y max | 3081.1 |

To hide/show a dataset click on the associated colored box

■ unknown
----------------------------------------
```
op rate                   : 1710 [WRITE:1710]
partition rate            : 1710 [WRITE:1710]
row rate                  : 1710 [WRITE:1710]
latency mean              : 0.6 [WRITE:0.6]
latency median            : 0.5 [WRITE:0.5]
latency 95th percentile   : 1.1 [WRITE:1.1]
latency 99th percentile   : 1.9 [WRITE:1.9]
latency 99.9th percentile : 6.2 [WRITE:6.2]
latency max               : 34.0 [WRITE:34.0]
total gc count            : 0
total gc mb               : 0
total gc time (s)         : 0
avg gc time(ms)           : NaN
stdev gc time(ms)         : 0
Total operation time      : 00:00:05
cmd: write n=10000 -rate threads=1 -graph file=topology1_Write.html title=Topology1_Write
```

Figure 4: Cassandra Stress Topology 1 Write

Figure 5: Cassandra Stress Topology 1 Mixed

10. The increase in all the measurements is linear as all the containers are on the host system so there is not much latency.

11. The last step was to measure the clock skew. This was done using the client-server code described in $Figure 4$



Figure 6: Client-Server code

12. The results of the same can be found in $Figure 5$
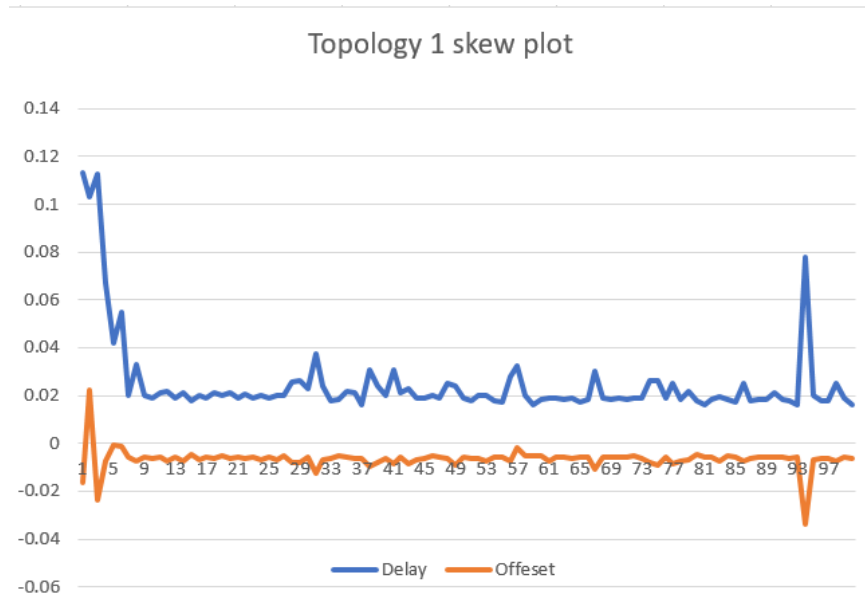
Figure 7: Skew output

13. Offset and Delay is almost constant apart from the start and the end.
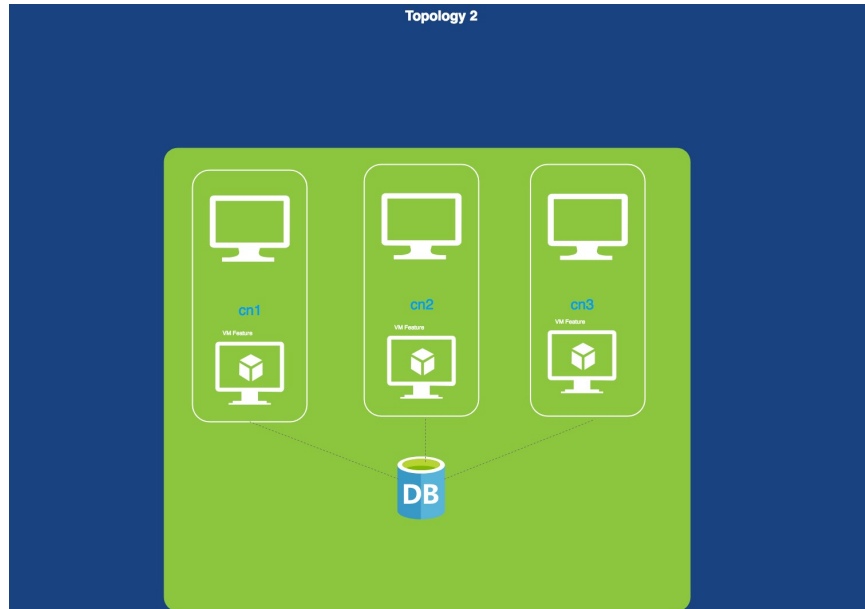
## 2. The Second Topology



Figure 8: Second Topology

In the topology, we had a three VMs running three different docker containers. Each container had a Cassandra node in it and was called cn1, cn2 and cn3. The VM was sprung up using Oracle's VirtualBox technology.

The steps we followed were as follows.

1. Setup three Virtual machine using Oracle's Virtual Box. The Operating System used was Ubuntu 16.04, with 20GB disk space and 4GB RAM.

2. NAT Network was used as network setting with different MAC address for all the VM. This was done to assign different IP to each VM.

3. The next step was to install Cassandra, NTP on all the VM's.

4. The last step was to establish connections between the containers containing Cassandra to form a ring

The figure describing the topology is $Figure 8$.

*2.1. Testing Topology 2*

1. The first step was to ensure that the docker container can ping one another. This was done using the $ping$ command, using the Docker's IP that can be obtained using the $ifconfig$ command. Once the ping is received, we can go ahead with the next steps.

2. We first installed NTP on each of the docker containers in the each VM. Next we treated one of the containers as the NTP server and the rest as the NTP clients.

3. To ensure that all the VMs/Containers follow the same time, we use NTP for Synchronization.

4. Like in the previous topology ntpstat was used to check if clocks are in sync.

   The NTP conf can be configured at \etc\ntp.conf. Once the conf files are set, the servers are rebooted.The sync status can be seen using the $ntpstat$ command.

5. Now we wanted to check if the key spaces are shared between the Cassandra ring.

6. We used the $cassandra - stress$ tool to get the read, write and mixed latencies of the same.

7. As the number of iterations increased, the read and write latencies started to gradually increase as well.

8. This can be seen in $Figure 10$, $Figure 11$ and $Figure 12$

9. The last step was to measure the clock skew. This was done using the client-server code described in $Figure 4$

10. The results of the same can be found in $Figure 9$

11. As different VM's are used, our operations per second for read, write and mix latencies increase.

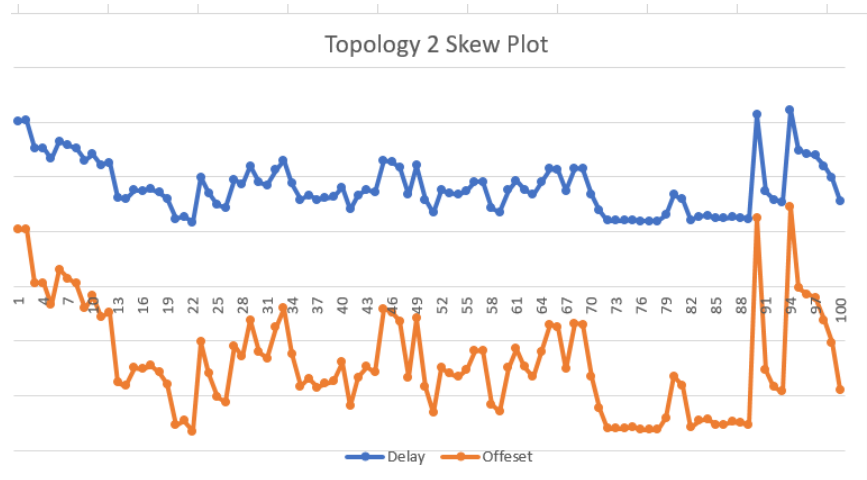12. Offset and Delay increase variably as the network delay increases.
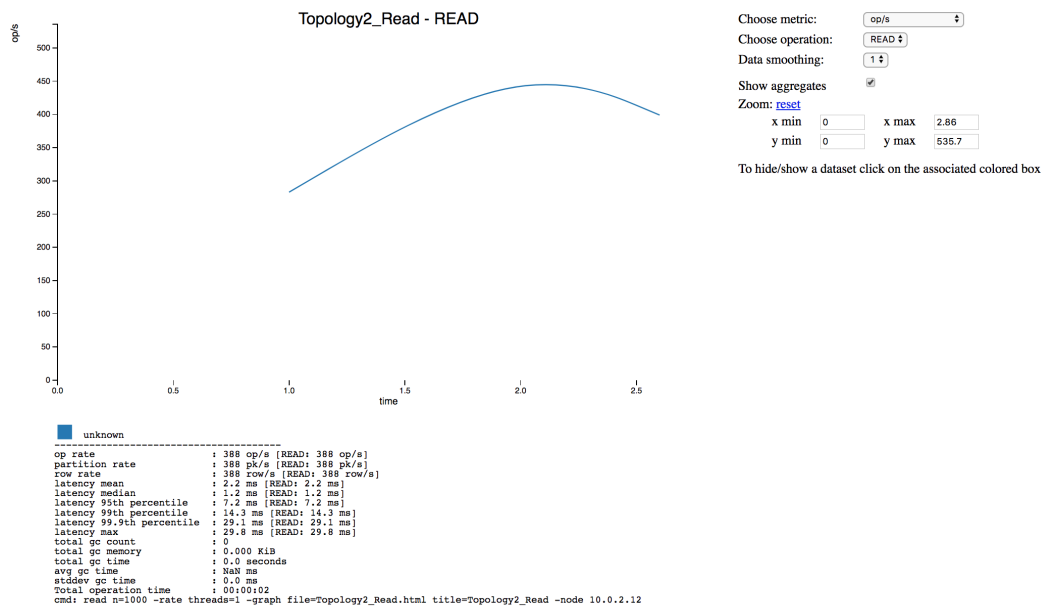


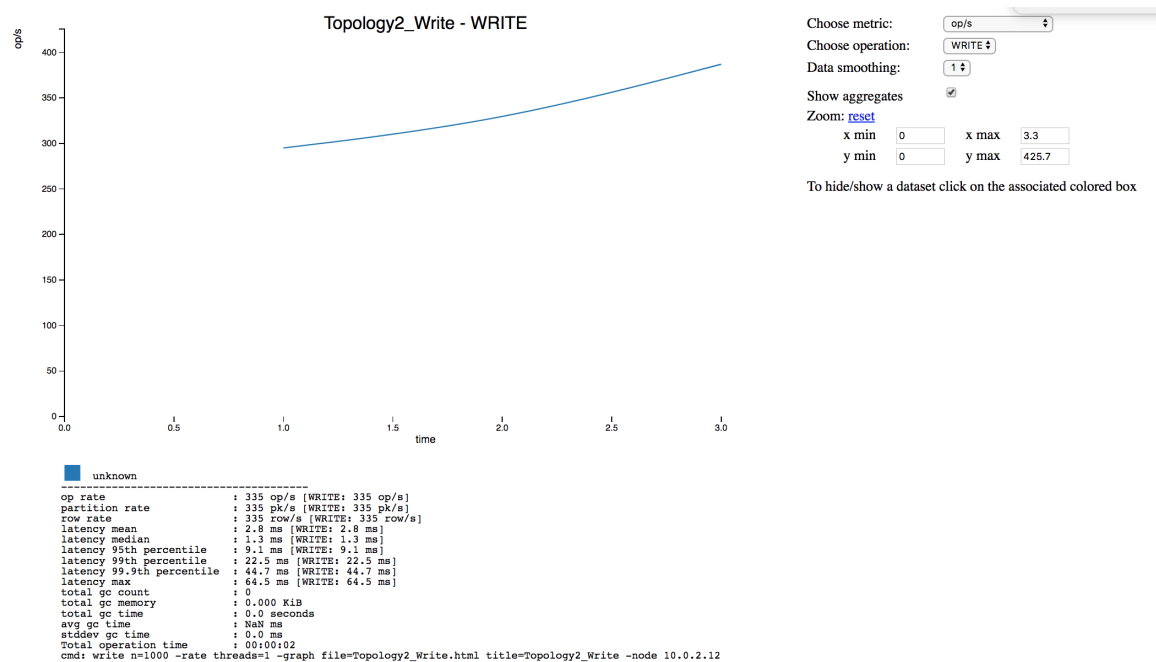Figure 9: Skew output

Figure 10: Cassandra Stress Topology 2 Read


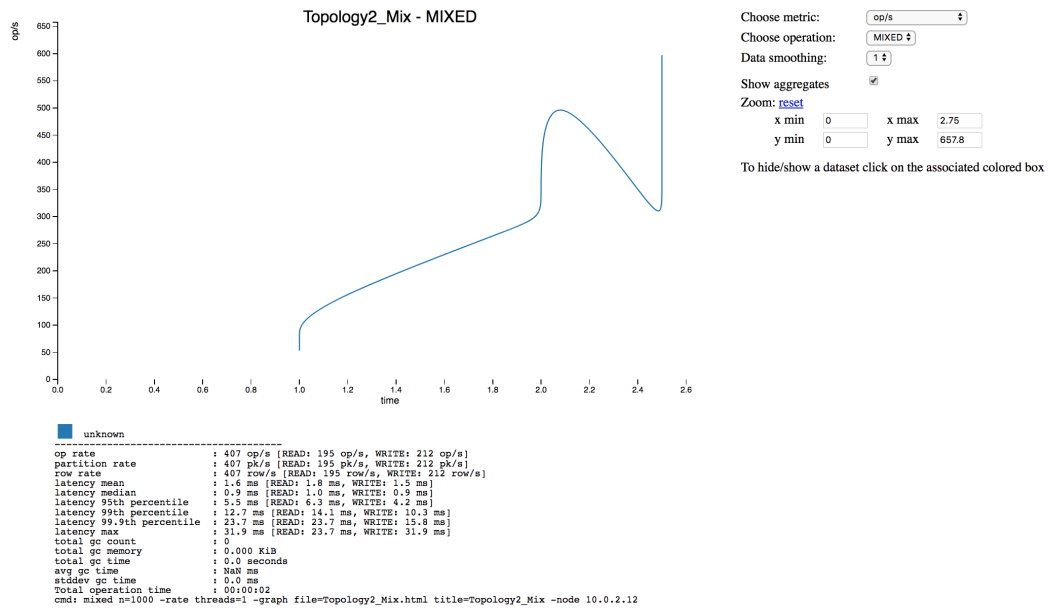
Figure 11: Cassandra Stress Topology 2 Write

Figure 12: Cassandra Stress Topology 2 Mixed

## 3. The Third Topology


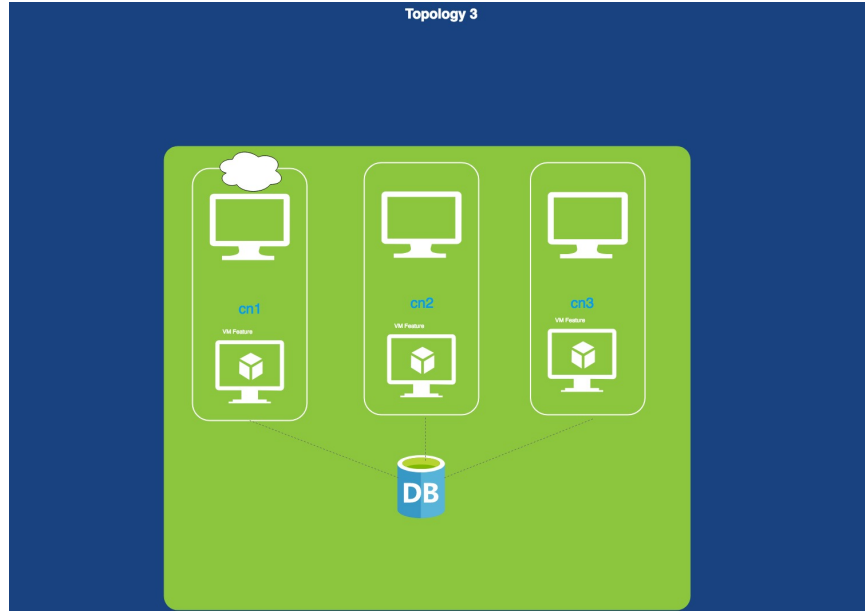
Figure 13: Toplogy 3

In the topology, we had a three VMs running three different docker containers in three different datacenters in Clemson, Utah and Wisconsin. Each container had a Cassandra node in it and was called cn1, cn2 and cn3. The VM was sprung up using CloudLab.

The steps we followed were as follows.

1. Setup a Virtual machines using Cloudlab. The Operating System used was Ubuntu 14.04.

2. The second step was to install Cassandra in each of the machines.

3. The last step was to establish connections between the cloud instances containing Cassandra to form a ring.

4. We had to open ports on the system for communication.

The figure describing the topology is $Figure 13$.

### 3.1. Testing Topology 3

1. The first step was to ensure that all machines can ping one another. This was done using the $ping$ command, using the machine's Public IP that can be obtained using the $ifconfig$ command. Once the ping is received, we can go ahead with the next steps.

2. To calculate time skew, we made one of the machine as server and modified NTP.conf file to point to NTP public IP.

3. We had to remove NTP from DHCP so that changes in NTP file can persist.

4. NTP.conf file of all the clients was edited to point to IP of the server so that they can sync their time with server.

5. We used the $cassandra - stress$ tool to get the read, write and mixed latencies of the same.

6. As the number of iterations increased, the read and write latencies started to gradually increase as well.The main reason was the machines were in different physical locations.

7. This can be seen in $Figure 15$, $Figure 16$ and $Figure 17$

8. The last step was to measure the clock skew. This was done using the client-server code described in $Fig4$

9. The results of the same can be found in $Figure14$

10. Delay and offset also increase due to change in the machines geography.
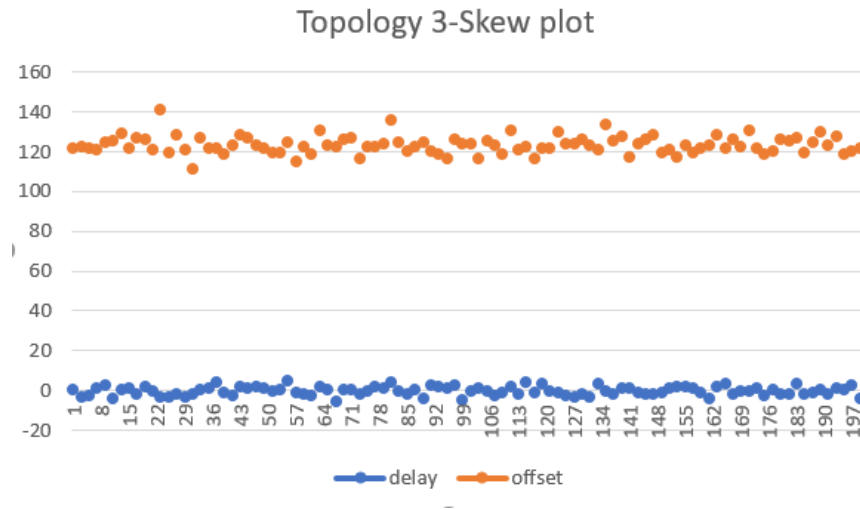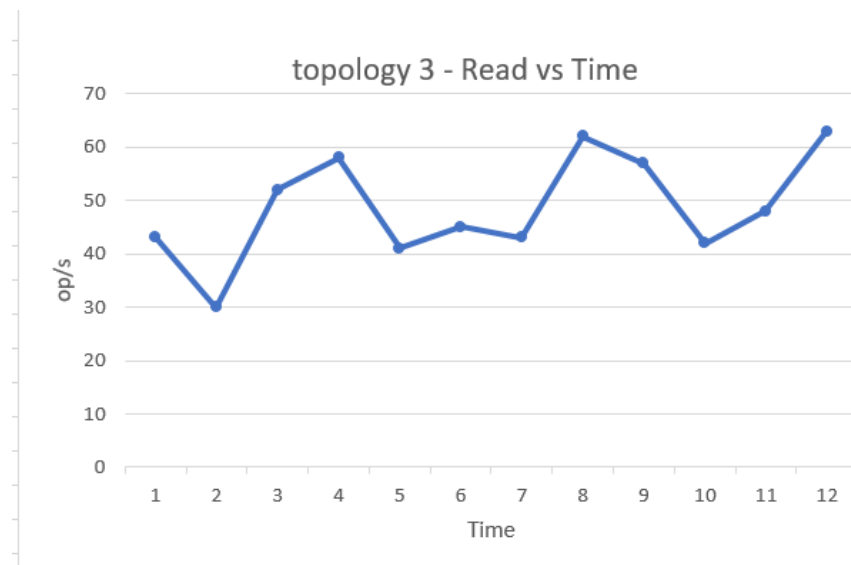
Figure 14: Skew output
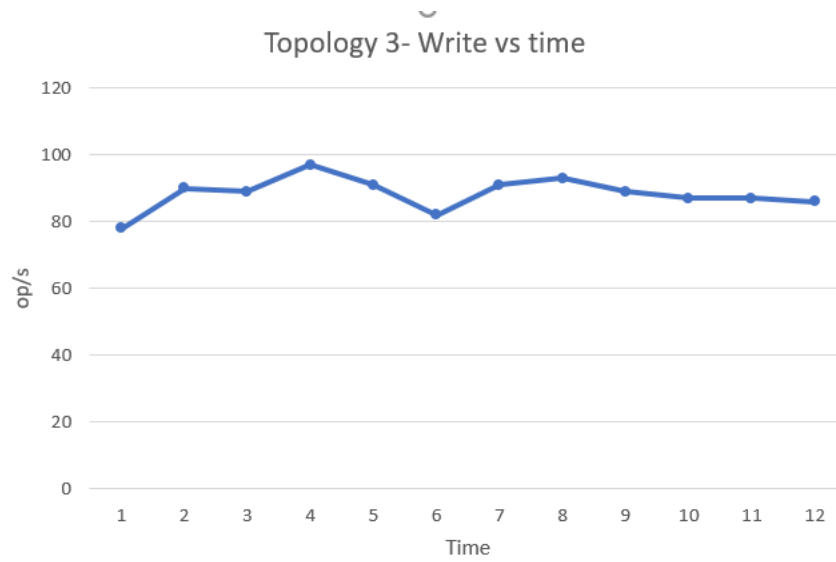
Figure 15: Cassandra Stress Topology 3 Read

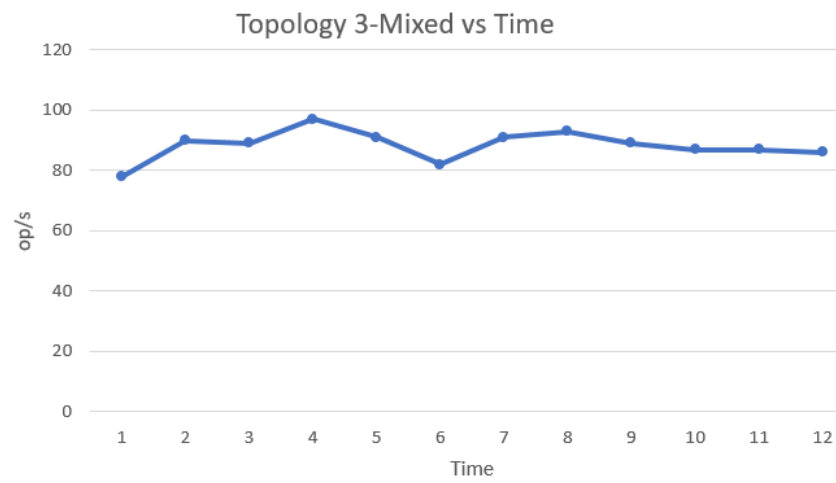Figure 16: Cassandra Stress Topology 3 Write



Figure 17: Cassandra Stress Topology 3 Mixed
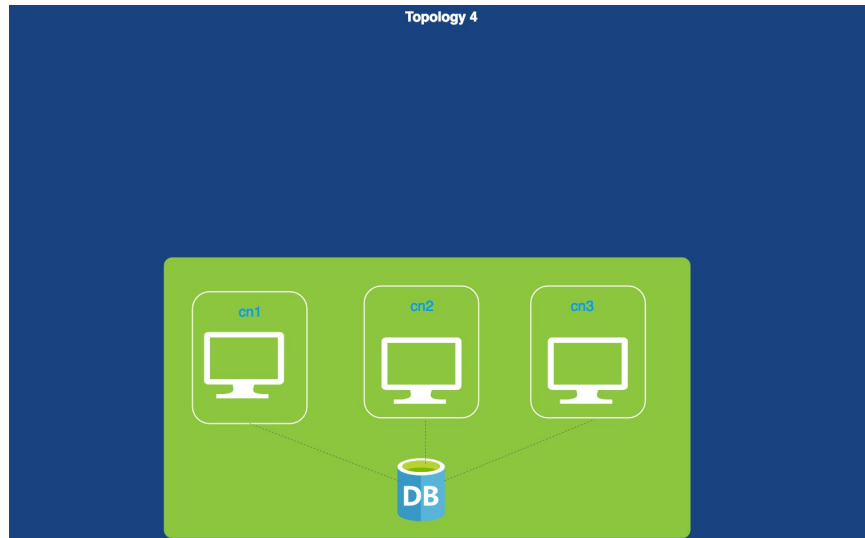
## 4. The Fourth Topology



Figure 18: Topology 4

In the topology, we had a three physical machines running in the same network. Each machine had a Cassandra node in it and was called cn1, cn2 and cn3. The VM was sprung up using Oracle's Virtualbox.

The steps we followed were as follows.

1. Install cassandra and ntp in all the machines.

2. Open ports on all the machines for communication.

The figure describing the topology is $Figure 18$.

### 4.1. Testing Topology 4

1. The first step was to ensure that the machines can ping one another. This was done using the $ping$ command, using the machine's IP that can be obtained using the $ifconfig$ command. Once the ping is received, we can go ahead with the next steps.

2. We used 1 machine as the server and synced its time with the public NTP server.

3. We modified the NTP.conf file of all the client machines to point to server IP.The sync status can be seen using the $ntpstat$ command.

4. Now we wanted to check if the key spaces are shared between the Cassandra ring.

5. We used the $cassandra - stress$ tool to get the read, write and mixed latencies of the same.

6. As the number of iterations increased, the read and write latencies linearly increased.

7. It was similar to the first topology as machines were in same LAN.

8. This can be seen in $Figure 19$, $Figure 20$ and $Figure 21$

9. The last step was to measure the clock skew. This was done using the client-server code described in $Figure 4$

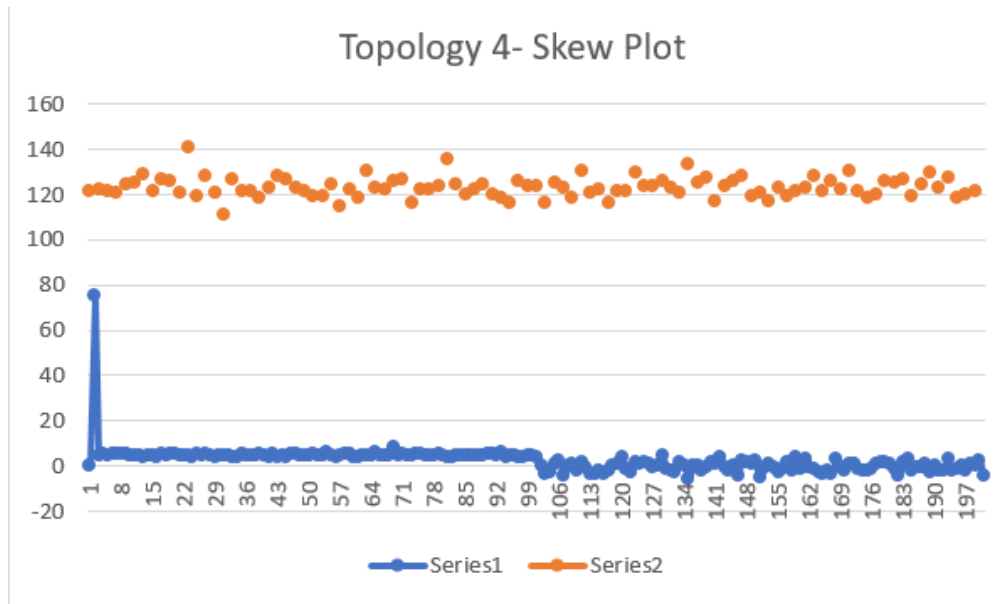10. The results of the same can be found in $Figure 22$
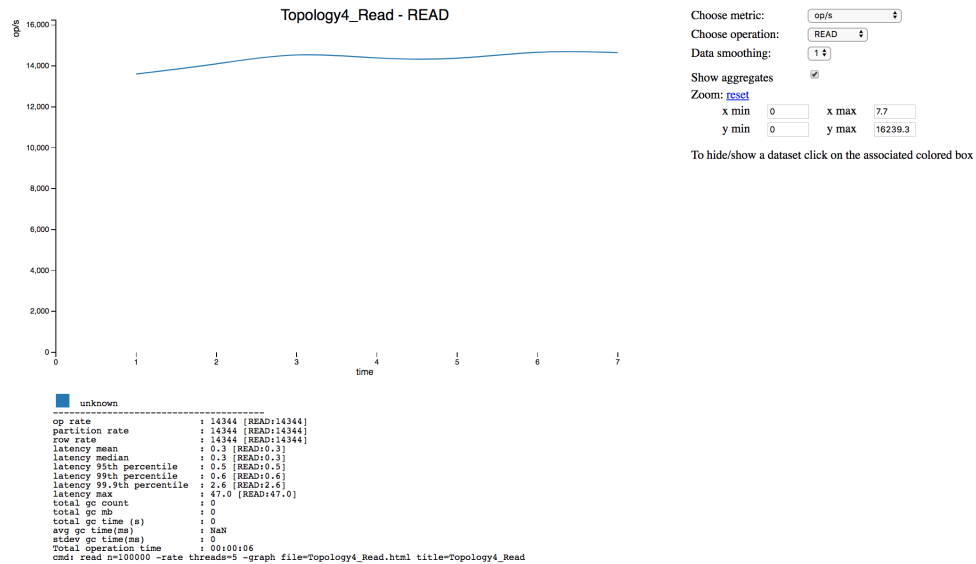
Figure 19: Skew output
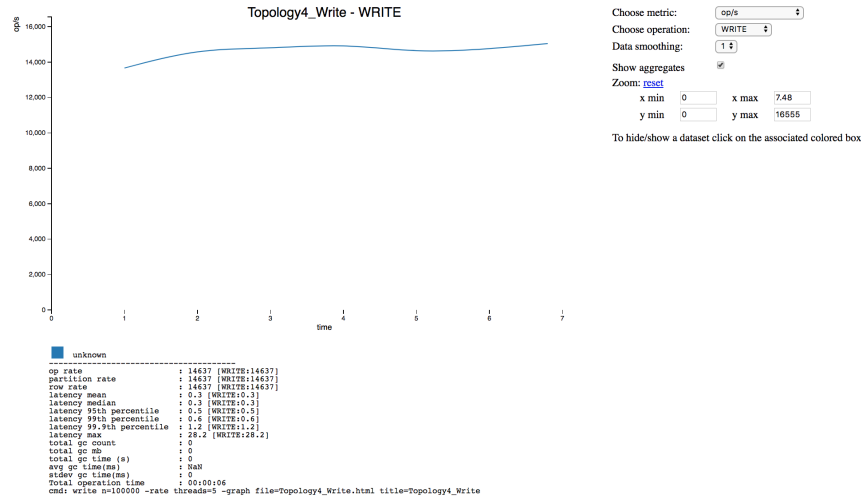


Figure 20: Cassandra Stress Topology 4 Read
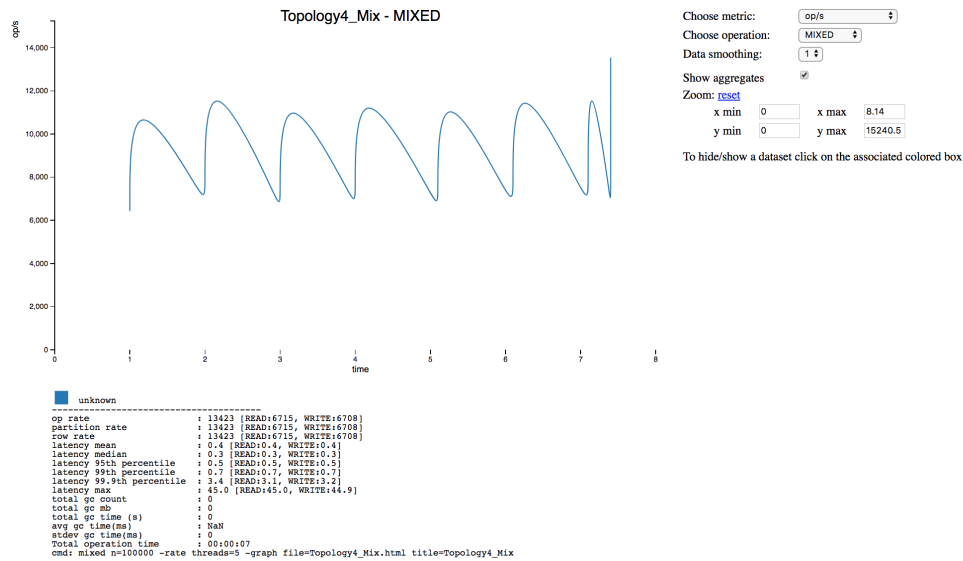
Figure 21: Cassandra Stress Topology 4 Write



Figure 22: Cassandra Stress Topology 4 Mixed

## 5. Conclusion

For the simple topology, read, write and mix latency will be less as compared to the complex topologies. Also delay and offset will be very less for the simple topology and will increase as the topology becomes more complex.We can see in topology 1 and topology 4, either all the containers are in the same system or are in the same LAN on different machines.So the latency, delay and offset is very less for these topologies. For topology 3, the latencies increase and the number of operation decrease as the machines are in Virtual Machines. For topology 4, the latencies increase and number of operations decrease as the machines are in different physical locations. We tested delay and offset after running the stress test so that the machines have some data.

14