# CSE 512 – Distributed Database System – Phase 2

Team Name: DNS

Deepak Soundararajan – 1211181124 - dsounda2@asu.edu
Narendrakumar Sampath Kumar -121106860 - nsampat1@asu.edu
Swathi Kanchan – 1211213195 - skanchan@asu.edu

**Requirement:**

To create a custom Java Function to perform Naive Spatial Join Query using Cartesian Product Algorithm with the given GeoSpark Jar and compare the Execution time/ Average memory / Average CPU utilization of the entire cluster (one master and two worker nodes) for the set of programs.

**Project Description:**

The given data set contains the points and the rectangular boundary in the geographical area. The project finds all the points from each of the rectangle of the query window. For every rectangular point, the method combines all the points. (i.e. Cartesian Product)

1. The function should be integrated into GeoSpark source code. The source code should be forked or downloaded from GeoSpark Github master branch LATEST.

2. The result format should be same as GeoSpark spatial join query: JavaPairRDD<Envelope, HashSet<Point>>.

This phase is continuation of the following tasks which was implemented during Phase 1

- Create GeoSpark SpatialRDD (PointRDD).
- Spatial Range Query: Query the PointRDD using this query window [x1(35.08),y1(-113.79),x2(32.99),y2(-109.73)].
    - Query the PointRDD
    - Build R-Tree index on PointRDD then query this PointRDD.
- Spatial KNN query: Query the PointRDD using this query point [x1(35.08),y1(-113.79)].
    - Query the PointRDD and find 5 Nearest Neighbors.
    - Build R-Tree index on PointRDD then query this PointRDD again.
- Spatial Join query: Create a GeoSpark RectangleRDD and use it to join PointRDD
    - Join the PointRDD using Equal grid without R-Tree index.
    - Join the PointRDD using Equal grid with R-Tree index.
    - Join the PointRDD using R-Tree grid without R-Tree index

**Implementation:**

The data is preloaded to Hadoop File System and the spark is clustered with one master node and two worker node. The method 'SpatialJoinQueryUsingCartesianProduct' is implemented in 'JoinQuery' class in GeoSpark library with the following algorithm (Source : https://gist.github.com/aniquetahir/acb2a781a55cf76a5d2a32d4f0a4d5d6).

1.Create a PointRDD objectRDD;

2.Create a RectangleRDD queryWindowRDD;

3.Collect rectangles from queryWindowRDD to one Java List L;

4. For each rectangle R in L

      do RangeQuery.SpatialRangeQuery(objectRDD, queryEnvelope, 0, true);

End;

5.Collect all results; //"Collect" is a standard function under SparkContext.

6.Parallelize the results to generate a RDD in this format: JavaPairRDD<Envelope, HashSet<Point>>.;//"Parallelize" is a standard function under SparkContext.

7.Return the result RDD;

We used Ganglia as cluster monitoring software to identify the statistics. We get the Execution time, Memory utilization, CPU utilization from Ganglia.

**Task 1: Comparing Task 2a & Task 2b of Phase 1**

In Phase 1, Task2a, we perform Spatial Range Query using the given query window. In Task 2b, We built R-Tree index on PointRDD then query the given PointRDD again. We ran it for 10 iteration and obtained the execution time, Average CPU Utilization, Average Memory Utilization. The statistics of task 2a and task 2b is compared as shown in below table.

The **Execution time** for 2b is lesser than 2a because of R-Tree Indexing

The **Average CPU utilization** for 2b is lesser than 2a because of R-Tree Indexing

The **Average Memory Utilization** is higher for 2b is higher than 2a as Indexing consumes slightly more memory.

|  | Spatial Range Query (2a) | Spatial Range Query (Index) (2b) |
|---|---|---|
| Execution Time | 5.1 Seconds | 4.7 Seconds |
| Average CPU Usage | 43.7% | 30.9% |
| Average Memory Usage | 5.7G | 5.9G |

## Task 2: Comparing Task 3a & 3b of Phase 1

In Task 3a and Task 3b of Phase 1, we perform Spatial KNN query. With the given query point, we find 5 Nearest Neighbors in Task 3a and we built R-Tree index on PointRDD then query this PointRDD again. We
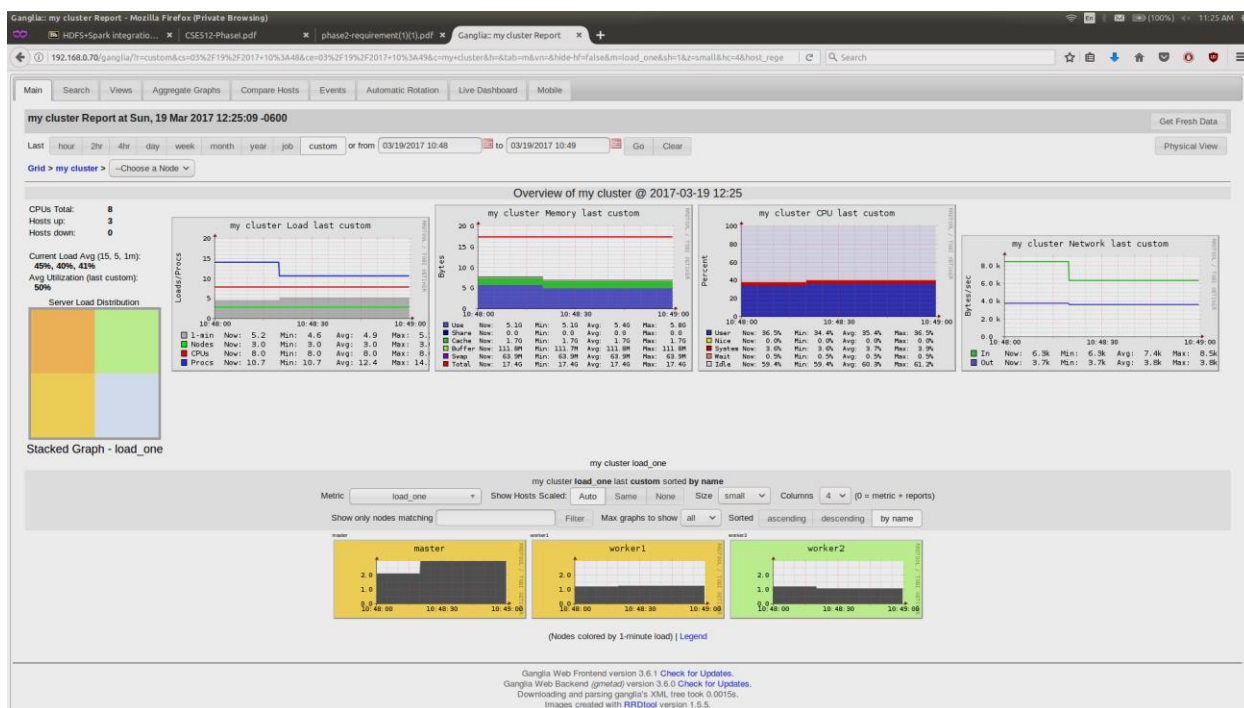
ran it for 10 iteration and obtained the execution time, Average CPU Utilization, Average Memory Utilization. The statistics of task 3a and task 3b is compared as shown below. The following are the observations made.

The **Execution time** to find 5 neighbors is slightly higher than finding neighbors built with R-Tree Index.

The **CPU utilization** is slightly higher for 3a when compared to 3b because 3b uses R-Tree Index

**Average Memory Utilization** of 3b is slightly higher than 3a because of usage of Indexing

|  | Spatial KNN Query | Spatial KNN Query (Index) |
|---|---|---|
| Execution Time | 5.9 Seconds | 5.1 Seconds |
| Average CPU Usage | 35.4% | 33.6% |
| Average Memory Usage | 5.1G | 6.5G |

## Task 3: Comparison of Task 4(a) & (b) of Phase 1

In the first part of this task 4, we compared the Execution times, Average memory, Average CPU utilization for spatial join query implementation of PointRDD using Equal grid without R-Tree index and PointRDD using Equal grid with R-Tree index. We ran the program for 10 iteration and suitable data is obtained.

We found that **Execution time** of Spatial join, pointRDD having equal grid without R-Tree index was higher than Spatial join, pointRDD having equal grid with R-Tree index.

The **Average memory utilization** increases slightly when used with Index. This is because of R-Tree indexing.

The **Average CPU utilization** of 4a is little higher than 4b because of the usage of R-Tree Index.

The difference can be seen in the table below,

|  | Spatial Join Query (Equal Grid) (4a) | Spatial Join Query (Equal Grid with Index) (4b) |
|---|---|---|
| Average Runtime | 313 Seconds | 227 Seconds |
| Average CPU Usage | 36.8% | 35.8% |
| Average Memory Usage | 6.1G | 7.3G |

**Task 3: Comparing Task 4 (a) & (c) of Phase 1**

In the second part of this task 3, we compared the Execution times, Average memory, Average CPU utilization for spatial join query implementation of PointRDD using Equal grid without R-Tree index. and PointRDD using R-Tree grid without R-Tree index. We ran for 10 iteration and suitable statistics is obtained.
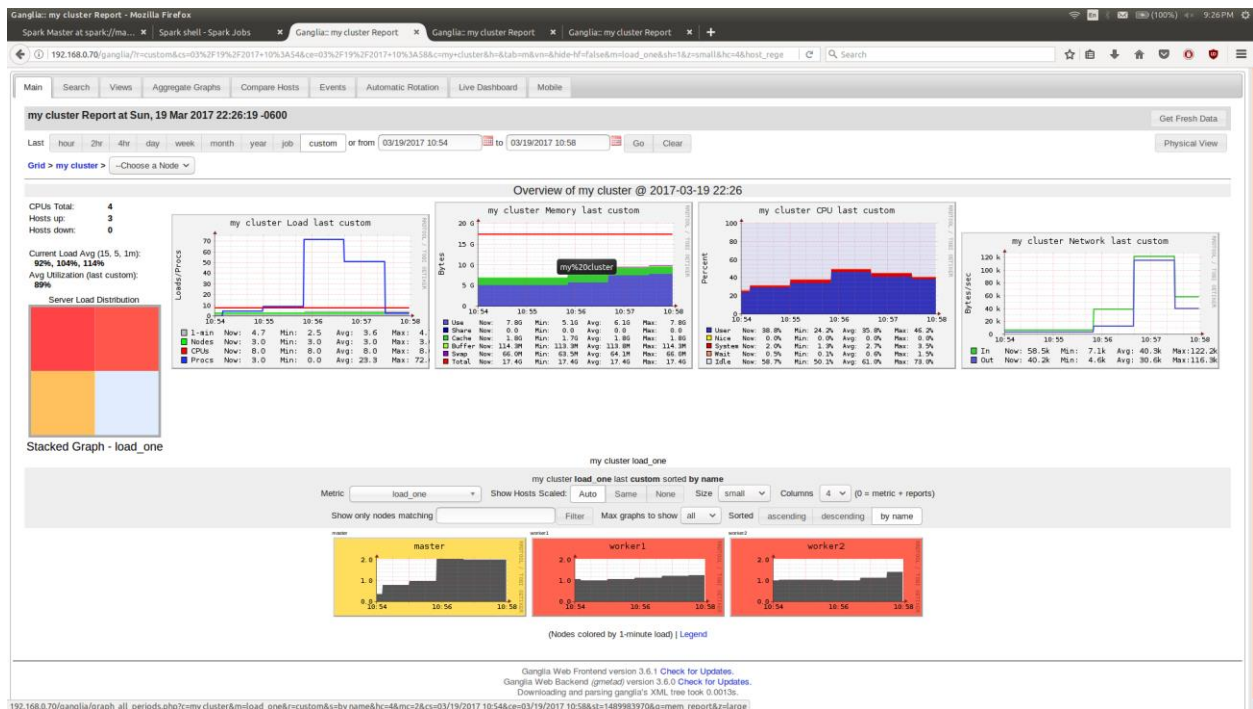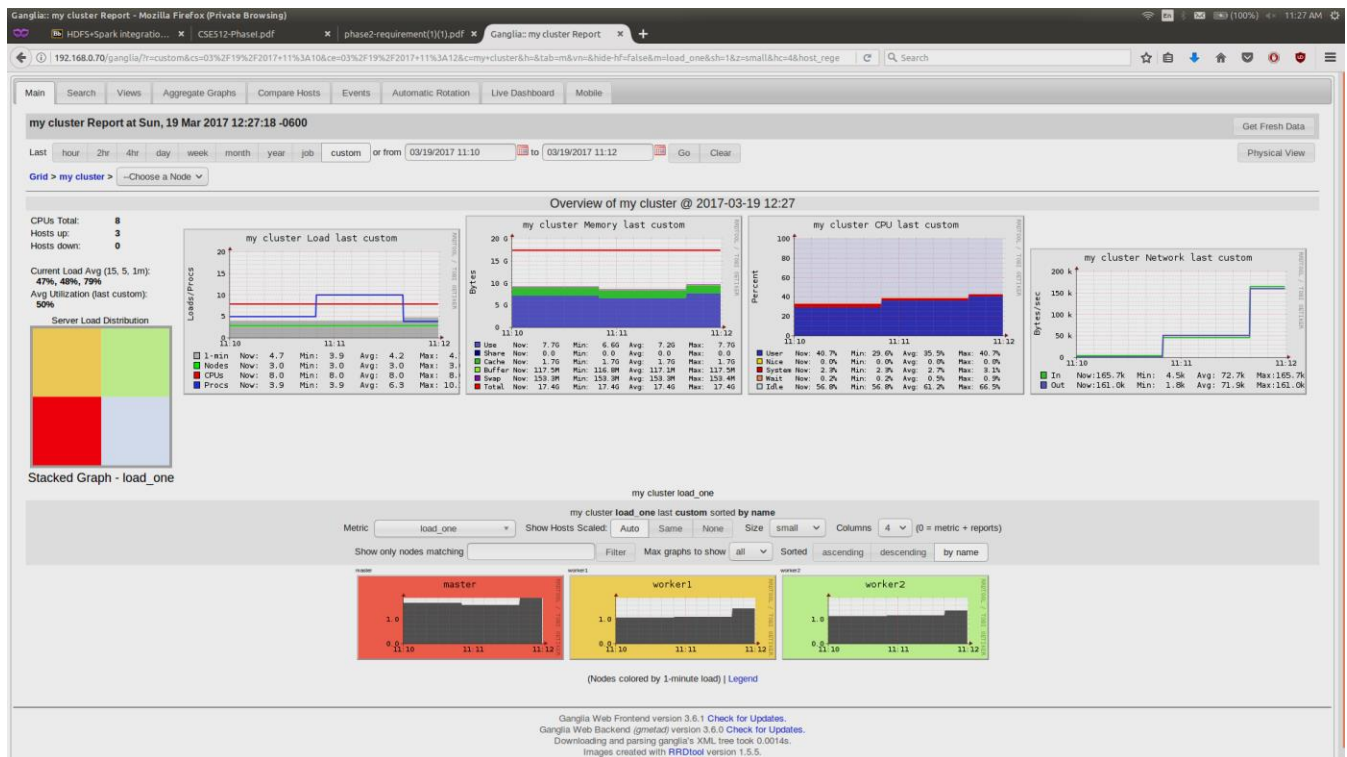
We found that **Execution time** of 4a is much higher than the execution time of 4c, since 4a uses Equal grid and 4c uses R-Tree grid

The **Average memory utilization** increases slightly from 4a to 4c indicating that the R-Tree grid utilizes slightly higher memory compared to Equal grid.

The **Average CPU utilization** of 4c is lesser when compared to 4a and 4b. The difference is because of the R-Tree grid usage whereas in 4a Equal grid is used.

The difference is seen in the table

|  | Spatial Join Query (Equal Grid) (4a) | Spatial Join Query (RTree Grid with Index) (4c) |
|---|---|---|
| Average Runtime | 313 Seconds | 243 Seconds |
| Average CPU Usage | 36.8% | 35.5% |
| Average Memory Usage | 6.1G | 7.2G |

## 2. Comparison of Phase1 Task4 (a) (b) (c) with Phase 2 Task A:

The task 4 of the Phase 1 which uses the normal join from existing GeoSpark Library and task 1 of the Phase 2, where we implemented a Naïve Spatial Join Query using Cartesian Product algorithm. This means that it is running using all the queryEnvelope (RectangleRDD) which takes huge amount of time. The Execution time is much higher than the Spatial Join Query executed in 4 (a) (b) (c)

The **Average memory utilization** of task1 of Phase 2 is much higher than any of the 4 (a) (b) (c) because of the Cartesian product algorithm.

The **Average CPU Utilization** is similar any of the 4 (a) (b) (c) but in phase2Task1 the *Maximum CPU usage goes up to 59.2%* whereas in 4 (a) (b) (c) it was equal to 40%. This is because of the Cartesian product algorithm.

|  | Spatial Join Query (Equal Grid without Index)4 (a) | Spatial Join Query (Index) 4 (b) | Spatial Join Query (R-Tree Grid without Index) 4 (c) | Spatial Join Query (Cartesian) (Phase 2) |
|---|---|---|---|---|
| Average Runtime | 313 Seconds | 227 Seconds | 243 Seconds | 6450 Seconds |
| Average CPU Usage | 36.8% | 35.8% | 35.5% | 36% |
| Average Memory Usage | 6.1G | 7.3G | 7.2G | 8.2G |

Phase 2 Task 1, Ganglia Statistics:

Spark Master at spark://ma... ✕ | Spark shell - Spark Jobs ✕ | Ganglia:: my cluster Report ✕ | Ganglia:: my cluster Report ✕ | Ganglia:: my cluster Report ✕ | +

192.168.0.70/ganglia/?r=custom&cs=03%2F19%2F2017+11%3A43&ce=03%2F19%2F2017+13%3A34&c=my+cluster&h=&tab=m&vn=&hide-hf=false&m=load_one&sh=1&z=small&hc=4&host_rege

Main | Search | Views | Aggregate Graphs | Compare Hosts | Events | Automatic Rotation | Live Dashboard | Mobile

**my cluster Report at Sun, 19 Mar 2017 22:23:26 -0600**

Get Fresh Data

Last  hour  2hr  4hr  day  week  month  year  job  custom  or from  03/19/2017 11:43  to  03/19/2017 13:34  Go  Clear

Physical View

Grid > my cluster >  --Choose a Node ⌄

Overview of my cluster @ 2017-03-19 22:23

CPUs Total: 4
Hosts up: 3
Hosts down: 0

Current Load Avg (15, 5, 1m):
88%, 100%, 126%
Avg Utilization (last custom):
89%

Server Load Distribution

Stacked Graph - load_one

my cluster load_one

my cluster **load_one** last **custom** sorted **by name**

Metric  load_one ⌄   Show Hosts Scaled:  Auto  Same  None   Size  small ⌄   Columns  4 ⌄  (0 = metric + reports)

Show only nodes matching [          ]  Filter   Max graphs to show  all ⌄   Sorted  ascending  descending  by name

master | worker1 | worker2

(Nodes colored by 1-minute load) | Legend