

# K-Means Clustering

Enhancing performance and scalability through parallelism using  
MPI

Luís Neto (a77763), Gonçalo Raposo (a77211)

September 5, 2019

## 1 Case of Study

*K-means clustering* is an algorithm commonly used to partition data into  $k$  groups based on their similarities. The algorithm consists in three step:

1. Selection of  $k$  random centroids.
2. Assignment Step: each point is assigned to a cluster based on closest centroid determined with *Euclidean distance*.
3. Update Step: each centroid is updated to be the mean of his cluster.

The latter two phases are repeated until the centroids no longer change.

## 2 Distributed Memory Implementation

### 2.1 Partition

In this development phase, we started by decomposing the algorithm in different phases taking into account the functional decomposition and the data domain of what could be parallelized.

The program begins by finding the maximum value of the data set. Based on the maximum value, we generate  $k$  random centroids, being so, the data domain is the data set points.

Afterward, we calculate the minimum distance between each point and centroids. Subsequently, all the points are assigned to the correspondent cluster. In this phase, it is fairly easy to understand that the data domain are the points from the data set.

In the third phase of the algorithm, we calculate the error and update all centroids. Since each centroid is the average of his cluster, the data domain is the data set's points.

### 2.2 Communication

As said before, this algorithm has 3 phases that need to execute in order, due to the data dependencies. After the first phase, the process with rank 0 will communicate the data set points and centroids to the workers.

In the end of phase 2, each process shares his local clusters and updates them, with the information received from the other processes.

In the last phase, we calculate part of the error and update the centroids, with the local information, afterwards each process shares his centroids and updates them, with the information received from the other processes.

Finally, we calculate the final error and send a final message from the process with rank 0 to the other processes. This message will continue or stop the algorithm. When the algorithm converges all Workers send a message to the process with rank 0 containing the current sets, otherwise, we restart the iterative process.

### 2.3 Agglomeration

In this phase, we consider that the granularity is a chunk of the data set with size equal to  $\frac{datasetSize}{\#Processes}$ . With this granularity, we can use several collective operations to increase performance and scalability.

In the first communication, we choose to use *Scatter* and *Broadcast* to divide the data set by the processes and to send the centroids, respectively.

Between phase 2 and 3, we tested 2 strategies, a *Reduce* followed by a *Broadcast* and a *Allreduce*, this strategies are used to share and update the local size counters of each cluster. The efficiency of these operations depends directly from the OpenMPI implementation, however, the latter strategy should be better than the former. In the next sections we will refer to this communication as comm2.

In order to calculate the final error, we need to send the information of each local updated centroid, therefore, we also tested the 2 strategies described above. In the next sections we will refer to this communication as comm3.

Finally, to stop the algorithm we *Broadcast* a message and if this message is positive we use an *Scatter* operation to get the clusters of process, in contrast, if the message is negative, we continue the algorithm. In the next section we will refer to this communications as error\_check.

### 2.4 Task Mapping

As explained before, in phase 2 we assign each point to a cluster based on closest centroid and in Phase 3 each centroid is updated to the average of his cluster, so, the workload for each element of data, centroid or point, is the same. With this we can do a regular partition of the data set and pass it to each process.

Since the algorithm is iterative, being the local centroids updated with the last iteration cluster's points, we can say that this algorithm fits in the category of the Heartbeat Algorithms.

### 2.5 Algorithm Analysis

In this section, we describe the sequential and parallel algorithm's complexity,  $N$  will represent the data set size,  $K$  the number of clusters and  $P$  the number of processes.

For the sequential version, in the first phase we iterate the data set in order to obtain the maximum value, subsequently, we generate all centroids, therefore, this phase complexity is  $N + K$ .

Phase 2 calculates the minimum distance between each point and centroids, consequentially, this phase complexity is  $N * K$ .

Finally, for phase 3, we iterate the clusters in order to calculate the first part of the error, update the clusters by iterating the data set and caculate the final error by iterating the clusters, hence, this phase complexity is  $2 * K + N$ . With this we can say that the sequential algorithm complexity is:

$$(N + K) + (N * (K + 1) + 2 * K) \quad (1)$$

In the parallel version, we equally divide the data set by the processes, so, we can rewrite the last complexity as:

$$(N + K) + (\frac{N}{P} * (K + 1) + 2 * K) \quad (2)$$

In this version, we also need to take in account the communication complexity, we will consider the following values for the primitives used.

1. Broadcast :  $\log_2(P)$
2. Reduce :  $\log_2(P)$
3. Scatter :  $P$
4. Gather :  $P$

After the first phase, we send the data set chunk to each process and clusters using the primitive Scatter and Broadcast, hence, the communication complexity is  $2 * \log_2(P) + P$ .

At the end of second phase, we reduce the local counters of each cluster to one process and broadcast the result, therefore, the complexity is  $2 * \log_2(P)$ .

After the update the local clusters, we reduce the local clusters to one process and broadcast the result, therefore, the complexity is  $4 * \log_2(P)$ .

Finally, after the final error calculation, we broadcast a message to each process in order to stop or continue the algorithm. If the algorithm continues the complexity is  $\log_2(P)$ , otherwise, the complexity is  $P + \log_2(P)$ , due to the fact that we use a *Gather* primitive to collect the final result.

Having the communication complexity we can write the parallel algorithm complexity as:

$$(N + K) + (\frac{N}{P} * (K + 1) + 2 * K) + 9 * \log_2 P + 2 * P \quad (3)$$

With this the speedup can be obtained by:

$$\frac{(N + K) + (N * (K + 1) + 2 * K)}{(N + K) + (\frac{N}{P} * (K + 1) + 2 * K) + 9 * \log_2(P) + 2 * P} \quad (4)$$

### 3 Input tests description

In order to test the developed algorithms, 2 different data sets were created by a Python script using a real uniform distribution, the first data set has 1966080 points and the last one 62914560. We choose the size of the first data set, in order to, fully occupy the level 3 cache. The other data set size was chosen, so that, it won't fit in any cache level.

## 4 Results Analysis

### 4.1 Communication

In order to choose the best mapping and communication strategy, we built 3 figures where we can observe the overall communication time, the partial communication time and the time spent on CPU and communications.

#### 4.1.1 Overall Communication Time

By observing the figure, is possible to see that, in terms of communication, the mapping by core is better up to 16 core, however, when we use 2 nodes, the time spent of this mapping increases and it is similar when we use 16 cores in both machines. It also possible to identify that *comm1* is the most expensive communication.

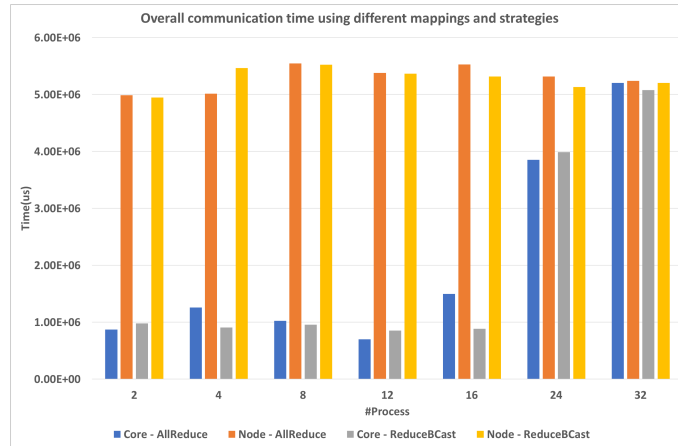


Figure 1: COMM time(us) with different mappings and strategies for the larger data set

#### 4.1.2 Partial Communication Time

In this figure, only the time spent in *comm2* and *comm3* are presented, so that we can compare both strategies and mappings. Between 2 and 32 processes, in both mappings, it is possible to see that for some #processes it is better to use *ReduceBroadcast* and for others to use *Allreduce*, however, the former strategy, in average, is better than the latter.

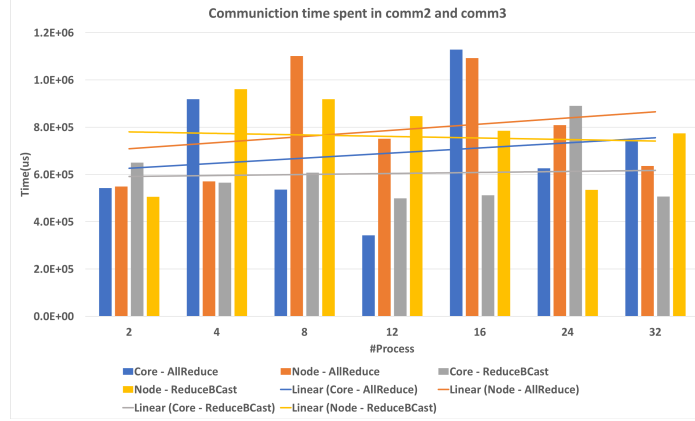


Figure 2: Communication time in *comm2* and *comm3*

## 4.2 Execution Profile

With this section, we will describe the algorithm's execution profile. As referred in section 2, this algorithm executes 4 communications and has 3 computation phases. The time spent in CPU is expected to be greater than the communication time, due to the fact that between iterations we only pass small messages. These messages either are the local sets counters or the centroids, with  $clusters * sizeof(double)$  and  $2 * clusters * sizeof(double)$  bytes, respectively.

Figure 2 represents the median time spent in each phase and in each communication. In order to obtain these results, we used the larger data set, core mapping and *Reduce + Broadcast* strategy.

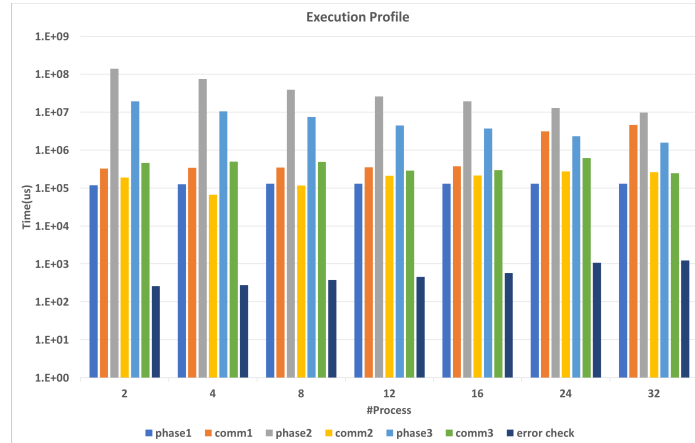


Figure 3: Execution profile for using the larger data set, core mapping and *Reduce + Broadcast* strategy.

As expected, when we use up to 16 processes, phase 2 and 3 times are greater than the communication times. Beyond that number communications will start to increase, since we start to use 2 computing platforms.

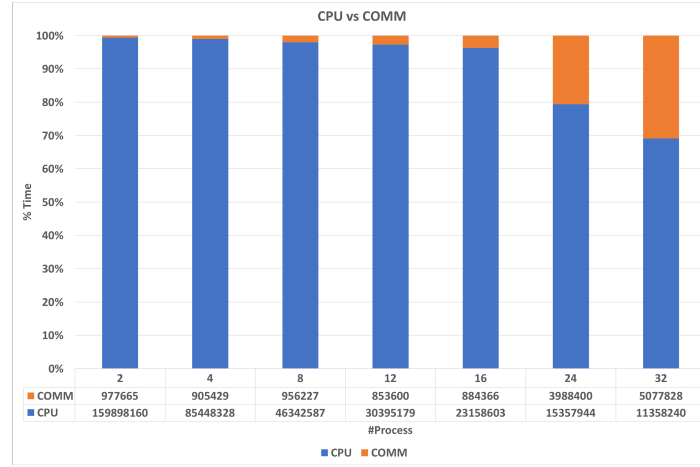


Figure 4: Execution profile for using the larger data set, core mapping and *Reduce + Broadcast* strategy.

### 4.3 Load Balance

In section D of the appendix, there are all the tables with the execution time of each process, for each possible combination of mapping, strategy, and size.

From the data in these tables, we can confirm that our algorithm has a good load balance, being the time of execution in each process approximately equal to  $\frac{Totaltime}{\#process}$ .

### 4.4 Achieved Performance

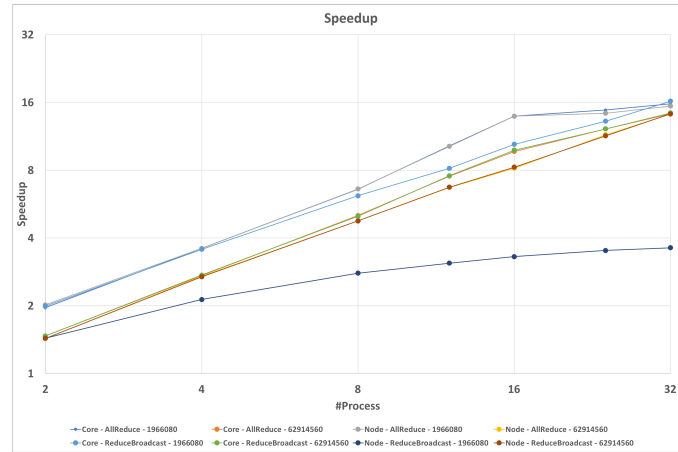


Figure 5: Speedup of each possible combination of size, mapping and strategy.

#### 4.4.1 Data set with 1966080 points

For this data set, it is possible to identify that the worst mapping is to map the processes by node and use the *ReduceBroadcast* method. Meanwhile, all the other strategies implemented have similar behaviors scaling well up to 16 processes and speedups close to the theoretical ones.

#### 4.4.2 Data set with 62914560 points

In respect to this data set, its general behavior is equivalent of the data set previously analyzed. In this case, the strategy of using *ReduceBroadcast* while applying the processes by each node available is the worst among all techniques. All the other, have similar speedups gains by the number of processes. To note that the speedup obtained is almost half by every process added to the task.

## 5 Conclusion

With the performance results obtained, we concluded that the best strategy to implement the *k-means* algorithm using a distributed memory approach is to map all the processes by core available while using the *ReduceBroadcast* collective primitive. We can also conclude that both paradigms have similar speedup, however, if the data set is larger than the RAM of one machine, only the MPI version will allow us to solve the problem.

Finally, taking in account the cost of implementation, the distributed memory paradigm is easier to program and the code is more readable.

# Appendix

## A Computing Platforms Hardware characterization

Manufacturer	Intel Corporation
Architecture	Ivy Bridge
Model	Xeon E5-2650v2
# Cores	16
# Threads	32
Processor Frequency	2.6 GHz
L1 Cache	32 KB (Data + Instruction)
L2 Cache	256 KB
L3 Cache	20 MB
RAM Memory	64Gb
Network	gbe/myri

Table 1: Nodes 641 hardware characterization

## B Node Mappings

### B.1 By Node

### B.1.1 2 Processes

```
[compute-641-8.local:06519] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:
[BB/../.././../.././../.././][../.././../.././../.././]
[compute-641-20.local:16484] MCW rank 1 bound to socket 0[core 0[hwt 0-1]]:
[BB/../.././../.././../.././][../.././../.././../.././]
```

### B.1.2 4 Processes

```
[compute-641-8.local:06526] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:
[BB/../.././../.././../.././][../.././../.././../.././]
[compute-641-8.local:06526] MCW rank 2 bound to socket 0[core 1[hwt 0-1]]:
[../../BB/../.././../.././../.././][../.././../.././../.././../.././]
[compute-641-20.local:16490] MCW rank 1 bound to socket 0[core 0[hwt 0-1]]:
[BB/../.././../.././../.././../.././][../.././../.././../.././../.././]
[compute-641-20.local:16490] MCW rank 3 bound to socket 0[core 1[hwt 0-1]]:
[../../BB/../.././../.././../.././][../.././../.././../.././../.././]
```

### B.1.3 8 Processes

```
[compute-641-8.local:06532] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:
```





```

[compute-641-8.local:06556] MCW rank 4 bound to socket 0[core 2[hwt 0-1]]:
[../BB/../.././../..][.././.././.././../..]
[compute-641-8.local:06556] MCW rank 6 bound to socket 0[core 3[hwt 0-1]]:
[.././../BB/../.././../..][.././.././.././../..]
[compute-641-8.local:06556] MCW rank 8 bound to socket 0[core 4[hwt 0-1]]:
[.././.././BB/../.././../..][.././.././.././../..]
[compute-641-8.local:06556] MCW rank 10 bound to socket 0[core 5[hwt 0-1]]:
[.././.././../BB/../.././../..][.././.././.././../..]
[compute-641-8.local:06556] MCW rank 12 bound to socket 0[core 6[hwt 0-1]]:
[.././.././.././BB/../.././../..][.././.././.././../..]
[compute-641-8.local:06556] MCW rank 14 bound to socket 0[core 7[hwt 0-1]]:
[.././.././.././../BB/../.././../..][.././.././.././../..]
[compute-641-20.local:16519] MCW rank 1 bound to socket 0[core 0[hwt 0-1]]:
[BB/../.././.././../..][.././.././.././../..]
[compute-641-20.local:16519] MCW rank 3 bound to socket 0[core 1[hwt 0-1]]:
[../BB/../.././.././../..][.././.././.././../..]
[compute-641-20.local:16519] MCW rank 5 bound to socket 0[core 2[hwt 0-1]]:
[.././BB/../.././.././../..][.././.././.././../..]
[compute-641-20.local:16519] MCW rank 7 bound to socket 0[core 3[hwt 0-1]]:
[.././../BB/../.././.././../..][.././.././.././../..]
[compute-641-20.local:16519] MCW rank 9 bound to socket 0[core 4[hwt 0-1]]:
[.././.././../BB/../.././../..][.././.././.././../..]
[compute-641-20.local:16519] MCW rank 11 bound to socket 0[core 5[hwt 0-1]]:
[.././.././.././BB/../.././../..][.././.././.././../..]
[compute-641-20.local:16519] MCW rank 13 bound to socket 0[core 6[hwt 0-1]]:
[.././.././.././../BB/../.././../..][.././.././.././../..]
[compute-641-20.local:16519] MCW rank 15 bound to socket 0[core 7[hwt 0-1]]:
[.././.././.././.././BB/../.././../..][.././.././.././../..]

```

#### B.1.6 24 Processes

```

[compute-641-8.local:06575] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:
[BB/../.././.././../..][.././.././.././../..]
[compute-641-8.local:06575] MCW rank 2 bound to socket 0[core 1[hwt 0-1]]:
[../BB/../.././.././../..][.././.././.././../..]
[compute-641-8.local:06575] MCW rank 4 bound to socket 0[core 2[hwt 0-1]]:
[.././BB/../.././.././../..][.././.././.././../..]
[compute-641-8.local:06575] MCW rank 6 bound to socket 0[core 3[hwt 0-1]]:
[.././../BB/../.././.././../..][.././.././.././../..]
[compute-641-8.local:06575] MCW rank 8 bound to socket 0[core 4[hwt 0-1]]:
[.././.././../BB/../.././../..][.././.././.././../..]
[compute-641-8.local:06575] MCW rank 10 bound to socket 0[core 5[hwt 0-1]]:
[.././.././.././BB/../.././../..][.././.././.././../..]
[compute-641-8.local:06575] MCW rank 12 bound to socket 0[core 6[hwt 0-1]]:
[.././.././.././../BB/../.././../..][.././.././.././../..]
[compute-641-8.local:06575] MCW rank 14 bound to socket 0[core 7[hwt 0-1]]:
[.././.././.././.././BB/../.././../..][.././.././.././../..]
[compute-641-8.local:06575] MCW rank 16 bound to socket 1[core 8[hwt 0-1]]:
[.././.././.././.././../BB/../.././../..][BB/../.././.././../..]

```

```
[compute-641-8.local:06575] MCW rank 18 bound to socket 1[core 9[hwt 0-1]]:
[.././.././.././.././..][../BB/./.././.././..]
[compute-641-8.local:06575] MCW rank 20 bound to socket 1[core 10[hwt 0-1]]:
[.././.././.././.././..][.././BB/./.././.././..]
[compute-641-8.local:06575] MCW rank 22 bound to socket 1[core 11[hwt 0-1]]:
[.././.././.././.././..][.././../BB/./.././.././..]
[compute-641-20.local:16537] MCW rank 1 bound to socket 0[core 0[hwt 0-1]]:
BB/./.././.././.././..][.././.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 3 bound to socket 0[core 1[hwt 0-1]]:
[../BB/./.././.././.././..][.././.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 5 bound to socket 0[core 2[hwt 0-1]]:
[.././BB/./.././.././.././..][.././.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 7 bound to socket 0[core 3[hwt 0-1]]:
[.././../BB/./.././.././..][.././.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 9 bound to socket 0[core 4[hwt 0-1]]:
[.././.././../BB/./.././..][.././.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 11 bound to socket 0[core 5[hwt 0-1]]:
[.././.././.././BB/./.././..][.././.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 13 bound to socket 0[core 6[hwt 0-1]]:
[.././.././.././.././BB/./..][.././.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 15 bound to socket 0[core 7[hwt 0-1]]:
[.././.././.././.././.././BB][.././.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 17 bound to socket 1[core 8[hwt 0-1]]:
[.././.././.././.././.././..][BB/./.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 19 bound to socket 1[core 9[hwt 0-1]]:
[.././.././.././.././.././..][../BB/./.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 21 bound to socket 1[core 10[hwt 0-1]]:
[.././.././.././.././.././..][.././BB/./.././.././.././.././..]
[compute-641-20.local:16537] MCW rank 23 bound to socket 1[core 11[hwt 0-1]]:
[.././.././.././.././.././..][.././../BB/./.././.././.././.././..]
```

### B.1.7 32 Processes

```
[compute-641-8.local:06601] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:
```

```
[BB/././././././././././././././././././././././././././././]
```

```
[compute-641-8.local:06601] MCW rank 2 bound to socket 0[core 1[hwt 0-1]]:
```

```
[./BB/././././././././././././././././././././././././././././]
```

```
[compute-641-8.local:06601] MCW rank 4 bound to socket 0[core 2[hwt 0-1]]:
```

```
[././BB/././././././././././././././././././././././././././././]
```

```
[compute-641-8.local:06601] MCW rank 6 bound to socket 0[core 3[hwt 0-1]]:
```

```
[././././BB/././././././././././././././././././././././././]
```

```
[compute-641-8.local:06601] MCW rank 8 bound to socket 0[core 4[hwt 0-1]]:
```

```
[././././././BB/././././././././././././././././././././././]
```

```
[compute-641-8.local:06601] MCW rank 10 bound to socket 0[core 5[hwt 0-1]]:
```

```
[././././././././BB/././././././././././././././././././././]
```

```
[compute-641-8.local:06601] MCW rank 12 bound to socket 0[core 6[hwt 0-1]]:
```

```
[././././././././././BB/./././././././././././././././././././]
```

```
[compute-641-8.local:06601] MCW rank 14 bound to socket 0[core 7[hwt 0-1]]:
```

```
[././././././././././././BB/./././././././././././././././././]
```



## B.2 By Core

### B.2.1 2 Processes

```
[compute-641-8.local:04530] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:
```

```
[BB/././././././././.[././././././././.]
```

```
[compute-641-8.local:04530] MCW rank 1 bound to socket 0[core 1[hwt 0-1]]:
```

```
[./BB/./././././././.[././././././././.]
```

### B.2.2 4 Processes

```
[compute-641-8.local:04536] MCW rank 0 bound to socket 0[core 0[hwt0-1]]:
[BB/../../../../][../../../../]
[compute-641-8.local:04536] MCW rank 1 bound to socket 0[core 1[hwt 0-1]]:
[./BB/../../../../][../../../../]
[compute-641-8.local:04536] MCW rank 2 bound to socket 0[core 2[hwt 0-1]]:
[../../../../BB/../../../../][../../../../]
[compute-641-8.local:04536] MCW rank 3 bound to socket 0[core 3[hwt 0-1]]:
[../../../../BB/../../../../][../../../../]
```

### B.2.3 8 Processes

```
[compute-641-8.local:04546] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:
[BB/../../../../][./../../../../]
[compute-641-8.local:04546] MCW rank 1 bound to socket 0[core 1[hwt 0-1]]:
[./BB/../../../../][./../../../../]
[compute-641-8.local:04546] MCW rank 2 bound to socket 0[core 2[hwt 0-1]]:
[../../../../BB/../../../../][./../../../../]
[compute-641-8.local:04546] MCW rank 3 bound to socket 0[core 3[hwt 0-1]]:
[../../../.././BB/../../../../][./../../../../]
[compute-641-8.local:04546] MCW rank 4 bound to socket 0[core 4[hwt 0-1]]:
[../../../../././BB/../../../../][./../../../../]
[compute-641-8.local:04546] MCW rank 5 bound to socket 0[core 5[hwt 0-1]]:
[../../../.././././BB/../../../../][./../../../../]
[compute-641-8.local:04546] MCW rank 6 bound to socket 0[core 6[hwt 0-1]]:
[../../../../././././BB/../../../../][./../../../../]
[compute-641-8.local:04546] MCW rank 7 bound to socket 0[core 7[hwt 0-1]]:
[../../../.././././././BB/../../../../][./../../../../]
```

### B.2.4 12 Processes

```
[compute-641-8.local:04564] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:  
[BB/././././././././././][././././././././././]  
[compute-641-8.local:04564] MCW rank 1 bound to socket 0[core 1[hwt 0-1]]:  
[./BB/./././././././././][././././././././././]  
[compute-641-8.local:04564] MCW rank 2 bound to socket 0[core 2[hwt 0-1]]:  
[././BB/././././././././][././././././././././]  
[compute-641-8.local:04564] MCW rank 3 bound to socket 0[core 3[hwt 0-1]]:
```

```

[../././BB/././././][.././././././././]
[compute-641-8.local:04564] MCW rank 4 bound to socket 0[core 4[hwt 0-1]]:
[.././././BB/./././][.././././././././]
[compute-641-8.local:04564] MCW rank 5 bound to socket 0[core 5[hwt 0-1]]:
[../././././BB/././][.././././././././]
[compute-641-8.local:04564] MCW rank 6 bound to socket 0[core 6[hwt 0-1]]:
[.././././././BB/./][.././././././././]
[compute-641-8.local:04564] MCW rank 7 bound to socket 0[core 7[hwt 0-1]]:
[../././././././BB][.././././././././]
[compute-641-8.local:04564] MCW rank 8 bound to socket 1[core 8[hwt 0-1]]:
[.././././././././][BB/./././././././]
[compute-641-8.local:04564] MCW rank 9 bound to socket 1[core 9[hwt 0-1]]:
[.././././././././][../BB/././././././]
[compute-641-8.local:04564] MCW rank 10 bound to socket 1[core 10[hwt 0-1]]:
[.././././././././][.././BB/./././././]
[compute-641-8.local:04564] MCW rank 11 bound to socket 1[core 11[hwt 0-1]]:
[.././././././././][../././BB/././././]

```

### B.2.5 16 Processes

```

[compute-641-8.local:04590] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:
[BB/././././././][.././././././././]
[compute-641-8.local:04590] MCW rank 1 bound to socket 0[core 1[hwt 0-1]]:
[../BB/./././././][.././././././././]
[compute-641-8.local:04590] MCW rank 2 bound to socket 0[core 2[hwt 0-1]]:
[.././BB/././././][.././././././././]
[compute-641-8.local:04590] MCW rank 3 bound to socket 0[core 3[hwt 0-1]]:
[../././BB/./././][.././././././././]
[compute-641-8.local:04590] MCW rank 4 bound to socket 0[core 4[hwt 0-1]]:
[.././././BB/././][.././././././././]
[compute-641-8.local:04590] MCW rank 5 bound to socket 0[core 5[hwt 0-1]]:
[../././././BB/./][.././././././././]
[compute-641-8.local:04590] MCW rank 6 bound to socket 0[core 6[hwt 0-1]]:
[.././././././BB/./][.././././././././]
[compute-641-8.local:04590] MCW rank 7 bound to socket 0[core 7[hwt 0-1]]:
[../././././././BB][.././././././././]
[compute-641-8.local:04590] MCW rank 8 bound to socket 1[core 8[hwt 0-1]]:
[../././././././][BB/././././././]
[compute-641-8.local:04590] MCW rank 9 bound to socket 1[core 9[hwt 0-1]]:
[../././././././][../BB/./././././]
[compute-641-8.local:04590] MCW rank 10 bound to socket 1[core 10[hwt 0-1]]:
[../././././././][.././BB/././././]
[compute-641-8.local:04590] MCW rank 11 bound to socket 1[core 11[hwt 0-1]]:
[../././././././][../././BB/./././]
[compute-641-8.local:04590] MCW rank 12 bound to socket 1[core 12[hwt 0-1]]:
[../././././././][.././././BB/././]
[compute-641-8.local:04590] MCW rank 13 bound to socket 1[core 13[hwt 0-1]]:
[../././././././][../././././BB/./]
[compute-641-8.local:04590] MCW rank 14 bound to socket 1[core 14[hwt 0-1]]:

```



```

[../..../BB/..][../..../..../..]
[compute-641-20.local:15351] MCW rank 22 bound to socket 0[core 6[hwt 0-1]]:
[../..../BB/..][../..../..../..]
[compute-641-20.local:15351] MCW rank 23 bound to socket 0[core 7[hwt 0-1]]:
[../..../BB/..][../..../..../..]

```

### B.2.7 32 Processes

```

[compute-641-8.local:04660] MCW rank 0 bound to socket 0[core 0[hwt 0-1]]:
[BB/..../..../..][../..../..../..]
[compute-641-8.local:04660] MCW rank 1 bound to socket 0[core 1[hwt 0-1]]:
[../BB/..../..../..][../..../..../..]
[compute-641-8.local:04660] MCW rank 2 bound to socket 0[core 2[hwt 0-1]]:
[../BB/..../..../..][../..../..../..]
[compute-641-8.local:04660] MCW rank 3 bound to socket 0[core 3[hwt 0-1]]:
[../..../BB/..../..][../..../..../..]
[compute-641-8.local:04660] MCW rank 4 bound to socket 0[core 4[hwt 0-1]]:
[../..../BB/..../..][../..../..../..]
[compute-641-8.local:04660] MCW rank 5 bound to socket 0[core 5[hwt 0-1]]:
[../..../BB/..../..][../..../..../..]
[compute-641-8.local:04660] MCW rank 6 bound to socket 0[core 6[hwt 0-1]]:
[../..../BB/..][../..../..../..]
[compute-641-8.local:04660] MCW rank 7 bound to socket 0[core 7[hwt 0-1]]:
[../..../BB/..][../..../..../..]
[compute-641-8.local:04660] MCW rank 8 bound to socket 1[core 8[hwt 0-1]]:
[../..../BB/..../..][BB/..../..../..]
[compute-641-8.local:04660] MCW rank 9 bound to socket 1[core 9[hwt 0-1]]:
[../..../BB/..../..][../BB/..../..../..]
[compute-641-8.local:04660] MCW rank 10 bound to socket 1[core 10[hwt 0-1]]:
[../..../BB/..../..][../BB/..../..../..]
[compute-641-8.local:04660] MCW rank 11 bound to socket 1[core 11[hwt 0-1]]:
[../..../BB/..../..][../BB/..../..../..]
[compute-641-8.local:04660] MCW rank 12 bound to socket 1[core 12[hwt 0-1]]:
[../..../BB/..../..][../BB/..../..../..]
[compute-641-8.local:04660] MCW rank 13 bound to socket 1[core 13[hwt 0-1]]:
[../..../BB/..../..][../BB/..../..../..]
[compute-641-8.local:04660] MCW rank 14 bound to socket 1[core 14[hwt 0-1]]:
[../..../BB/..][../BB/..../..../..]
[compute-641-8.local:04660] MCW rank 15 bound to socket 1[core 15[hwt 0-1]]:
[../..../BB/..][../BB/..../..../..]
[compute-641-20.local:15368] MCW rank 16 bound to socket 0[core 0[hwt 0-1]]:
[BB/..../..../..][../..../..../..]
[compute-641-20.local:15368] MCW rank 17 bound to socket 0[core 1[hwt 0-1]]:
[../BB/..../..../..][../..../..../..]
[compute-641-20.local:15368] MCW rank 18 bound to socket 0[core 2[hwt 0-1]]:
[../BB/..../..../..][../..../..../..]
[compute-641-20.local:15368] MCW rank 19 bound to socket 0[core 3[hwt 0-1]]:
[../BB/..../..../..][../..../..../..]
[compute-641-20.local:15368] MCW rank 20 bound to socket 0[core 4[hwt 0-1]]:

```



```

[../..../BB/../../..][../..../..../..../..]
[compute-641-20.local:15368] MCW rank 21 bound to socket 0[core 5[hwt 0-1]]:
[../..../..../BB/../../..][../..../..../..../..]
[compute-641-20.local:15368] MCW rank 22 bound to socket 0[core 6[hwt 0-1]]:
[../..../..../BB/../../..][../..../..../..../..]
[compute-641-20.local:15368] MCW rank 23 bound to socket 0[core 7[hwt 0-1]]:
[../..../..../..../BB][../..../..../..../..]
[compute-641-20.local:15368] MCW rank 24 bound to socket 1[core 8[hwt 0-1]]:
[../..../..../..../..][BB/../../..../..../..]
[compute-641-20.local:15368] MCW rank 25 bound to socket 1[core 9[hwt 0-1]]:
[../..../..../..../..][../BB/../../..../..../..]
[compute-641-20.local:15368] MCW rank 26 bound to socket 1[core 10[hwt 0-1]]:
[../..../..../..../..][../..BB/../../..../..../..]
[compute-641-20.local:15368] MCW rank 27 bound to socket 1[core 11[hwt 0-1]]:
[../..../..../..../..][../..../BB/../../..../..]
[compute-641-20.local:15368] MCW rank 28 bound to socket 1[core 12[hwt 0-1]]:
[../..../..../..../..][../..../..BB/../../..../..]
[compute-641-20.local:15368] MCW rank 29 bound to socket 1[core 13[hwt 0-1]]:
[../..../..../..../..][../..../..../BB/../../..]
[compute-641-20.local:15368] MCW rank 30 bound to socket 1[core 14[hwt 0-1]]:
[../..../..../..../..][../..../..../..../BB/..]
[compute-641-20.local:15368] MCW rank 31 bound to socket 1[core 15[hwt 0-1]]:
[../..../..../..../..][../..../..../..../..BB]

```

## C Algorithm Efficiency

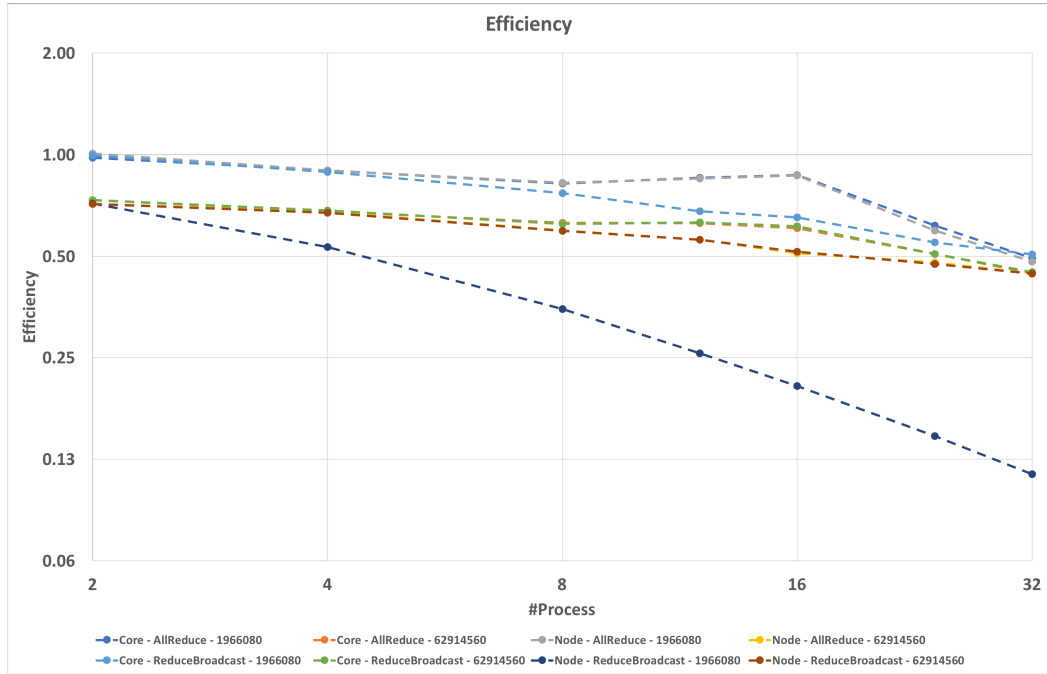


Figure 6: Algorithm's Efficiency per #process

## D Algorithm Load Balance Data

	Core - AllReduce - 1966080						
P1	3250225	1773115	973711	627109	463435	440398	419407
P2	3250228	1772714	973413	626933	462724	441614	419194
P3	X	1773034	973340	626604	462799	440323	419984
P4	X	1772687	973459	626168	462175	442058	418862
P5	X	X	973083	626173	462340	441660	419297
P6	X	X	973139	626069	462194	440175	418218
P7	X	X	972885	626059	462350	441756	416986
P8	X	X	973190	684428	465260	437604	415664
P9	X	X	X	684601	463795	438268	417197
P10	X	X	X	626825	462869	438323	416684
P11	X	X	X	684599	463348	436440	415836
P12	X	X	X	684722	463745	438373	419485
P13	X	X	X	X	463088	438238	419174
P14	X	X	X	X	463060	438374	419421
P15	X	X	X	X	463086	438125	419515
P16	X	X	X	X	463438	451492	430961
P17	X	X	X	X	X	449151	427549
P18	X	X	X	X	X	445565	425934
P19	X	X	X	X	X	444617	423470
P20	X	X	X	X	X	444137	423015
P21	X	X	X	X	X	445233	423530
P22	X	X	X	X	X	445975	424788
P23	X	X	X	X	X	447211	425460
P24	X	X	X	X	X	456848	439437
P25	X	X	X	X	X	X	434067
P26	X	X	X	X	X	X	429440
P27	X	X	X	X	X	X	428112
P28	X	X	X	X	X	X	428037
P29	X	X	X	X	X	X	428467
P30	X	X	X	X	X	X	429372
P31	X	X	X	X	X	X	431790
P32	X	X	X	X	X	X	442305

Figure 7: Load Balance using L3 size data set, Reduce+Broadcast strategy and map by Node , time in microseconds

	Core - AllReduce - 62914560						
P1	161036295	86903623	47109528	31552725	24556621	20064670	17661518
P2	161036272	86903302	47110309	31553172	24555431	20067490	17654873
P3	X	86903602	47109988	31551525	24554672	20062099	17655962
P4	X	86903415	47109397	31551204	24554053	20062325	17654505
P5	X	X	47109253	31551154	24554340	20068430	17663347
P6	X	X	47109643	31550735	24554593	20068649	17663313
P7	X	X	47109673	31550827	24554901	20069067	17663703
P8	X	X	47109670	31614591	24558898	20068447	17655587
P9	X	X	X	31614057	24556431	20064589	17645346
P10	X	X	X	31613559	24554585	20048186	17637659
P11	X	X	X	31613406	24554226	20046218	17649617
P12	X	X	X	31550456	24554414	20061102	17649142
P13	X	X	X	X	24554743	20061019	17649595
P14	X	X	X	X	24554733	20061336	17649584
P15	X	X	X	X	24556432	20061479	17649440
P16	X	X	X	X	24554742	20110732	17729604
P17	X	X	X	X	X	20105073	17731499
P18	X	X	X	X	X	20092806	17729184
P19	X	X	X	X	X	20092534	17724946
P20	X	X	X	X	X	20092555	17739526
P21	X	X	X	X	X	20093093	17729173
P22	X	X	X	X	X	20111303	17726664
P23	X	X	X	X	X	20110182	17725191
P24	X	X	X	X	X	20137692	17715527
P25	X	X	X	X	X	X	17724934
P26	X	X	X	X	X	X	17724229
P27	X	X	X	X	X	X	17719171
P28	X	X	X	X	X	X	17728725
P29	X	X	X	X	X	X	17724437
P30	X	X	X	X	X	X	17717284
P31	X	X	X	X	X	X	17718655
P32	X	X	X	X	X	X	17731028

Figure 8: Load Balance using RAM size data set, Reduce+Broadcast strategy and map by Node, time in microseconds

	Core - ReduceBcast - 1966080						
P1	3155165	1773472	969258	629323	463229	457348	426766
P2	3155159	1773300	969267	629085	463575	455275	425985
P3	X	1773329	968355	628792	463863	455341	428254
P4	X	1773150	967825	628633	463768	453822	428263
P5	X	X	967844	628562	463448	456861	427903
P6	X	X	967869	628682	463236	456830	421403
P7	X	X	967419	628531	463248	456727	425675
P8	X	X	967891	629650	465361	455240	427316
P9	X	X	X	685234	464185	454487	424097
P10	X	X	X	685122	464263	456525	424099
P11	X	X	X	685000	463539	454288	422376
P12	X	X	X	685183	463298	456141	423693
P13	X	X	X	X	463325	455947	423801
P14	X	X	X	X	463330	456260	424189
P15	X	X	X	X	464449	456018	424031
P16	X	X	X	X	463483	460912	434469
P17	X	X	X	X	X	463963	432236
P18	X	X	X	X	X	461663	426899
P19	X	X	X	X	X	457791	426685
P20	X	X	X	X	X	458454	427310
P21	X	X	X	X	X	459891	428651
P22	X	X	X	X	X	458462	429390
P23	X	X	X	X	X	460641	431757
P24	X	X	X	X	X	471793	444969
P25	X	X	X	X	X	X	440127
P26	X	X	X	X	X	X	432119
P27	X	X	X	X	X	X	430071
P28	X	X	X	X	X	X	430724
P29	X	X	X	X	X	X	431578
P30	X	X	X	X	X	X	432446
P31	X	X	X	X	X	X	433717
P32	X	X	X	X	X	X	450662

Figure 9: Load Balance using L3 size data set, Allreduce strategy and map by Node, time in microseconds

	Core - ReduceBcast - 62914560						
P1	161043479	86514848	47501097	31445145	24232983	19922419	17513921
P2	161043479	86515047	47500784	31444821	24231747	19913484	17518168
P3	X	86514986	47500782	31444295	24230803	19923246	17521297
P4	X	86514874	47500174	31443949	24230478	19918035	17522094
P5	X	X	47500072	31443558	24230583	19924570	17521869
P6	X	X	47500157	31444069	24230692	19923958	17522383
P7	X	X	47500080	31443946	24231144	19924761	17522524
P8	X	X	47500116	31501773	24236718	19918781	17514292
P9	X	X	X	31501621	24233419	19898834	17506906
P10	X	X	X	31501374	24231317	19905524	17497493
P11	X	X	X	31501286	24231905	19911501	17515602
P12	X	X	X	31444044	24231744	19911038	17515644
P13	X	X	X	X	24232067	19910336	17516258
P14	X	X	X	X	24232127	19910098	17515932
P15	X	X	X	X	24232183	19910120	17515908
P16	X	X	X	X	24233318	20114950	17593285
P17	X	X	X	X	X	20103959	17591275
P18	X	X	X	X	X	20101632	17590966
P19	X	X	X	X	X	20095424	17603031
P20	X	X	X	X	X	20101151	17592016
P21	X	X	X	X	X	20100805	17573031
P22	X	X	X	X	X	20125474	17584302
P23	X	X	X	X	X	20146462	17592072
P24	X	X	X	X	X	20122313	17600156
P25	X	X	X	X	X	X	17588067
P26	X	X	X	X	X	X	17585896
P27	X	X	X	X	X	X	17588863
P28	X	X	X	X	X	X	17580339
P29	X	X	X	X	X	X	17573667
P30	X	X	X	X	X	X	17577871
P31	X	X	X	X	X	X	17578234
P32	X	X	X	X	X	X	17592524

Figure 10: Load Balance using RAM size data set, Allreduce strategy and map by Node, time in microseconds

	Node - AllReduce - 1966080						
P1	3231829	1821536	1059013	801655	635935	506297	414972
P2	3231527	1816871	1059890	801138	630146	502961	415135
P3	X	1819714	1058749	801263	627646	519818	414517
P4	X	1822751	1061166	801104	629996	503805	408036
P5	X	X	1056222	805210	629514	509146	410984
P6	X	X	1056193	801396	629894	507517	409270
P7	X	X	1058738	801124	663297	500848	406187
P8	X	X	1070490	799820	633851	501307	405575
P9	X	X	X	798066	640142	502101	411021
P10	X	X	X	803408	634748	502972	404724
P11	X	X	X	804454	636100	503457	411568
P12	X	X	X	816201	636252	555165	406693
P13	X	X	X	X	636537	503835	413451
P14	X	X	X	X	639326	503365	408818
P15	X	X	X	X	641984	552707	421303
P16	X	X	X	X	646885	550597	409550
P17	X	X	X	X	X	504925	417559
P18	X	X	X	X	X	507141	420760
P19	X	X	X	X	X	554562	410092
P20	X	X	X	X	X	556582	412163
P21	X	X	X	X	X	555876	409823
P22	X	X	X	X	X	555418	410641
P23	X	X	X	X	X	554854	409554
P24	X	X	X	X	X	517967	408006
P25	X	X	X	X	X	X	410303
P26	X	X	X	X	X	X	407059
P27	X	X	X	X	X	X	411123
P28	X	X	X	X	X	X	412023
P29	X	X	X	X	X	X	407639
P30	X	X	X	X	X	X	408557
P31	X	X	X	X	X	X	410155
P32	X	X	X	X	X	X	428619

Figure 11: Load Balance using L3 size data set, Reduce + Broadcast strategy and map by Core, time in microseconds

	Node - AllReduce - 62914560						
P1	165669824	88673686	50697079	35301892	30006547	21629653	17669833
P2	165668525	88669452	50694332	35306792	29985782	21614941	17678305
P3	X	88672332	50696845	35308445	29978377	21619709	17676724
P4	X	88674666	50690407	35309697	29972569	21606586	17670924
P5	X	X	50689378	35311755	29977817	21616073	17674100
P6	X	X	50685830	35990818	29975799	21617664	17673696
P7	X	X	50688896	35995895	29976594	21624602	17674380
P8	X	X	50711100	36331200	30035776	21674330	17739876
P9	X	X	X	36336452	30022202	21665539	17741001
P10	X	X	X	36343867	30021443	21665492	17735519
P11	X	X	X	36350064	30029102	21679318	17719898
P12	X	X	X	36349503	30017281	21685348	17719166
P13	X	X	X	X	30054683	21684497	17717961
P14	X	X	X	X	30026517	21683520	17718203
P15	X	X	X	X	30027522	21680122	17727836
P16	X	X	X	X	30044368	21681035	17746496
P17	X	X	X	X	X	21684465	17729964
P18	X	X	X	X	X	21663192	17749670
P19	X	X	X	X	X	21686930	17744322
P20	X	X	X	X	X	21678638	17750931
P21	X	X	X	X	X	21677458	17751519
P22	X	X	X	X	X	21672803	17727006
P23	X	X	X	X	X	21672027	17728112
P24	X	X	X	X	X	21682548	17745140
P25	X	X	X	X	X	X	17742407
P26	X	X	X	X	X	X	17715518
P27	X	X	X	X	X	X	17742753
P28	X	X	X	X	X	X	17737345
P29	X	X	X	X	X	X	17743104
P30	X	X	X	X	X	X	17729228
P31	X	X	X	X	X	X	17731266
P32	X	X	X	X	X	X	17745877

Figure 12: Load Balance using RAM size data set, Reduce + Broadcast strategy and map by Core, time in microseconds



	Node - ReduceBcast - 1966080						
P1	4769571	3319116	2852045	2795192	2626630	2590025	3486717
P2	4763137	3218843	2855613	2574420	2628842	2579839	3482236
P3	X	3398992	2665578	2575431	2632587	2567779	3484764
P4	X	3399891	2486828	2614688	2591595	2559504	3486592
P5	X	X	3096521	2527876	2591737	2561690	3482755
P6	X	X	3078220	2870354	2586427	2561645	3484954
P7	X	X	2874649	2603084	2588170	2566731	3930841
P8	X	X	2875481	2861190	2453687	2596578	3931305
P9	X	X	X	2685506	2381884	2145319	3687572
P10	X	X	X	2660669	2607978	2144296	3685494
P11	X	X	X	2571919	2360296	2111990	3688690
P12	X	X	X	2814330	2909931	2155428	3699763
P13	X	X	X	X	2617342	2115365	3239476
P14	X	X	X	X	2835887	2605731	3241517
P15	X	X	X	X	2439940	2547792	3239300
P16	X	X	X	X	2710503	2806348	2802772
P17	X	X	X	X	X	2615661	2803609
P18	X	X	X	X	X	2820762	2561643
P19	X	X	X	X	X	2851583	2563217
P20	X	X	X	X	X	2165686	3683784
P21	X	X	X	X	X	2304310	3678261
P22	X	X	X	X	X	2628671	3681018
P23	X	X	X	X	X	2541785	3683220
P24	X	X	X	X	X	2575586	3239388
P25	X	X	X	X	X	X	3240542
P26	X	X	X	X	X	X	3239705
P27	X	X	X	X	X	X	2802709
P28	X	X	X	X	X	X	2803605
P29	X	X	X	X	X	X	2554067
P30	X	X	X	X	X	X	2554844
P31	X	X	X	X	X	X	2556066
P32	X	X	X	X	X	X	2765386

Figure 13: Load Balance using L3 size data set, Allreduce strategy and map by Core, time in microseconds

	Node - ReduceBcast - 62914560						
P1	166102740	88960723	49696011	36219776	29716549	20897737	17673362
P2	166103715	88958579	49697891	36220140	29701424	20900162	17660936
P3	X	88960066	49699290	36219901	29699057	20903363	17670917
P4	X	88964274	50741800	36217366	29702582	20909084	17673707
P5	X	X	50733685	36220259	29705891	20911422	17672952
P6	X	X	50735062	36258872	29701023	20912716	17673383
P7	X	X	50737021	36253880	29704897	20935616	17675606
P8	X	X	50738769	36256008	29770103	20982735	17738255
P9	X	X	X	36241614	29770668	20986214	17733260
P10	X	X	X	36241149	29751542	21019368	17739864
P11	X	X	X	36258416	29754914	21071065	17739348
P12	X	X	X	36243437	29772451	21408312	17756759
P13	X	X	X	X	29751555	21417829	17761854
P14	X	X	X	X	29748729	21636141	17771440
P15	X	X	X	X	29754814	21712246	17771898
P16	X	X	X	X	29766501	21758064	17708460
P17	X	X	X	X	X	21770251	17690599
P18	X	X	X	X	X	21817116	17782136
P19	X	X	X	X	X	21824528	17730952
P20	X	X	X	X	X	21904103	17732214
P21	X	X	X	X	X	21925369	17739886
P22	X	X	X	X	X	21988943	17728491
P23	X	X	X	X	X	21998948	17740739
P24	X	X	X	X	X	21961429	17717975
P25	X	X	X	X	X	X	17848321
P26	X	X	X	X	X	X	17723676
P27	X	X	X	X	X	X	17791483
P28	X	X	X	X	X	X	17718393
P29	X	X	X	X	X	X	17743480
P30	X	X	X	X	X	X	17727001
P31	X	X	X	X	X	X	17726600
P32	X	X	X	X	X	X	17740825

Figure 14: Load Balance using RAM size data set, Allreduce strategy and map by Core, time in microseconds