

Parallel and Multicore Computing Project 2

Shuang Li 1044137

Introduction

This report outlines and analyses an implementation of a message passing parallel algorithm for k-means++ clustering and the n-body simulation.

The algorithm is implemented with Open MPI where a number of nodes collaborate and passes messages to each other to carry out the task required.

The task is described as follows (n denotes the total number of points and m denotes the total number of nodes):

1. n points are sampled by all m nodes from a Gaussian mixture model (GMM).
2. All nodes collaborate to cluster the points into m clusters using k-means++.
3. Each cluster is moved into one single node so that a single node contains all points of a cluster.
4. The nodes collaborate with each other to perform n-body simulation where each point represents a body with a mass of 1 and bodies are attracted to each other via gravity.

The implementation for sampling all n points and clustering them with k-means++ are all rather naive and straightforward. The interesting part is the n-body gravity simulation.

To reduce the simulation time, the Barnes-Hut simulation approximation algorithm is used so that points over a certain distance away are simplified into a single center of mass when calculating the force they exert on another point. To reduce the communication overhead, each node constructs a Barnes-Hut tree and then prune the tree for each other nodes so that the tree only contains the nodes that are actually needed for the calculation of points in each other node. This pruned tree will from now on be referred to as the partial tree.

For node A to construct a partial tree to be sent to node B, we would like to make sure that any node of the tree where $s/d < \theta$ is kept, where s is the width of the region the node of the tree represent, d is the distance between the center of mass of the node and the closest point in cluster B to the center of mass to the node, and θ is a predefined constant. However, we do not want to calculate the closest points in cluster B to every node in the Barnes-Hut tree for cluster A and that would be slow in both computation time and communication overhead. Hence, we instead take the closest point in cluster B with respect to the center of mass of all points in cluster A, and calculate a hyperplane that is perpendicular to the line through the center of mass and the closest point, that goes through the closest point. Assuming that all points in each cluster are closer to their center of mass than any points in other clusters, the distance from any point within cluster A to this hyperplane should be shorter than the distance from point within cluster A to any point within cluster B.

The pseudocode describing the algorithm is shown in Algorithm 1, with the function that it calls shown in Algorithm 2 to 7. Note that the pseudocode is written in the perspective on one computing node.

While in theory, this algorithm should be faster than naively parallelizing Barnes-Hut algorithm by transferring all points to a root cluster and calculating one single Barnes-Hut tree to be sent to all clusters. In practice, it turns out the algorithm has a number of other overheads.

These overheads include calculating m partial trees (again here m is the total number of computing nodes), encoding them into an array of bytes to be sent over to each node and decoding the bytes into a partial tree, plus the extra computing time of creating m^2 partial trees. Furthermore, because each simulation step relies on the assumption of points of each all points in each cluster are closer to their center of mass than any points in other clusters being true, the clusters have to be recomputed and points that changed clusters sent over the corresponding nodes. This adds additional overheads that,

Algorithm 1

```
1:  $\theta \leftarrow 0.1$ 
2:  $\Delta Time \leftarrow 0.01$ 
3:  $m \leftarrow$  number of compute nodes.
4:  $NodeIndex \leftarrow$  the index for the compute node.  $\triangleright$  This is equivalent to MPI_Comm_rank but 1 indexed.
5:  $N \leftarrow$  number of points
6:  $D \leftarrow$  number of dimensions
7:  $c \leftarrow$  number of GMM components
8: for  $i \leftarrow 1$  to  $c$  do
9:    $prob[i] \leftarrow$  the probability of the  $i^{th}$  GMM component
10:   $gmmc[i] \leftarrow$  the  $i^{th}$  GMM component
11: end for
12: for  $i \leftarrow 1$  to  $N/m$  do
13:    $points[i] \leftarrow \text{SAMPLE}(gmmc[\text{SAMPLEWEIGHTEDDISCRETEDISTRIBUTION}(prob)])$ 
14:    $\triangleright \text{SAMPLEWEIGHTEDDISCRETEDISTRIBUTION}(weights)$  generate a random number between 1 and the length of  $weights$  with the number  $i$  having a probability of  $weights[i]/\text{SUM}(weights)$ 
15: end for
16:  $centroids[1] \leftarrow \text{CHOOSERANDOMONE}(points)$   $\triangleright$  choose centroid for k-means++
17: for  $i \leftarrow 2$  to  $m$  do  $\triangleright$  choose  $m$  centroids for  $m$  clusters
18:   for  $j \leftarrow 1$  to  $\text{LEN}(points)$  do
19:      $distances[i] \leftarrow$  distance of  $points[i]$  to the closest centroid
20:   end for
21:    $DistSum \leftarrow$  sum of distances
22:   if  $NodeIndex = 1$  then
23:     for  $i \leftarrow 1$  to  $m$  do
24:        $AllSum[i] \leftarrow DistSum$  from node  $i$ 
25:     end for
26:      $NextCentroidNode \leftarrow \text{SAMPLEWEIGHTEDDISCRETEDISTRIBUTION}(AllSum)$ 
27:   end if
28:    $NextCentroidNode \leftarrow NextCentroidNode$  from node 1
29:   if  $NodeIndex = NextCentroidNode$  then
30:      $NextCentroid \leftarrow points[\text{SAMPLEWEIGHTEDDISCRETEDISTRIBUTION}(distances)]$ 
31:   end if
32:    $centroids[i] \leftarrow NextCentroid$  from node  $NextCentroidNode$ 
33: end for
34:  $ClusterIndices \leftarrow \text{KMEANS}(centroids, points)$ 
35:  $points \leftarrow$  points in  $NodeIndex^{th}$  cluster as indicated by  $ClusterIndices$  in all nodes  $\triangleright$  Significant implementation detail of message passing between nodes omitted.
36: for  $i \leftarrow 1$  to  $\text{LEN}(points)$  do
37:    $velocities[i] \leftarrow 0$ 
38: end for
39:  $InitVariance \leftarrow \text{GETVARIANCE}(points)$ 
40: while  $True$  do
41:    $\text{SIMULATE}(points, velocities)$ 
42:    $CurrentVariance \leftarrow \text{GETVARIANCE}(points)$ 
43:   if  $CurrentVariance < InitVariance/2$  then
44:     break while loop
45:   end if
46:    $ClusterIndices \leftarrow \text{KMEANS}(centroids, points)$   $\triangleright$  one iteration only, implemented separately in real code but reusing pseudocode here.
47:    $points \leftarrow$  points in  $NodeIndex^{th}$  cluster as indicated by  $ClusterIndices$  in all nodes
48: end while
```

Algorithm 2

```
1: function KMEANS(centroids, points)
2:   while centroids changed since last iteration do
3:     for  $i \leftarrow 1$  to LEN(points) do
4:        $ClusterIndices[i] \leftarrow$  index of centroid closest to points
5:     end for
6:     for  $i \leftarrow 1$  to LEN(points) do
7:        $sums[ClusterIndices[i]] \leftarrow sums[ClusterIndices[i]] + points[i]$ 
8:        $counts[ClusterIndices[i]] \leftarrow counts[ClusterIndices[i]] + 1$ 
9:     end for
10:    if NodeIndex == 1 then
11:      for  $i \leftarrow 1$  to LEN(centroids) do
12:         $allSums \leftarrow$  sum of  $sums[i]$  from all nodes
13:         $allCounts \leftarrow$  sum of  $counts[i]$  from all nodes
14:         $centroids[i] \leftarrow allSums/allCounts$   $\triangleright$  Here it is actually each dimension of allSums
        divided by allCounts.
15:      end for
16:    end if
17:    centroids  $\leftarrow$  centroids from node 1
18:  end while
19:  return ClusterIndices
20: end function
```

Algorithm 3

```
1: function GETVARIANCE(points)
2:   mean  $\leftarrow$  average of all points
3:   if NodeIndex = 1 then
4:     sum  $\leftarrow$  0
5:     allSize  $\leftarrow$  0
6:     for  $i \leftarrow 1$  to m do
7:        $means[i] \leftarrow$  mean from node i
8:        $sizes[i] \leftarrow$  LEN(points) from node i
9:        $sum \leftarrow sum + means[i] * sizes[i]$ 
10:       $allSize \leftarrow allSize + sizes[i]$ 
11:    end for
12:     $allMean \leftarrow sum/allSize$ 
13:  end if
14:  mean  $\leftarrow$  allMean from node 1
15:  variance  $\leftarrow$  0
16:  for  $i \leftarrow 1$  to LEN(points) do
17:     $variance \leftarrow variance + \text{DISTANCE}(points[i], mean)$ 
18:  end for
19:   $allVariance \leftarrow$  sum of all variance from all nodes
20:  return allVariance
21: end function
```

Algorithm 4

```
1: function SIMULATE(points, velocities)
2:    $b1 \leftarrow points[1]$ 
3:    $b2 \leftarrow points[1]$ 
4:   for  $i \leftarrow 2$  to  $LEN(points)$  do
5:     for each dimension do
6:        $b1[dimension] \leftarrow MIN(b1[dimension], points[i][dimension])$ 
7:        $b2[dimension] \leftarrow MAX(b2[dimension], points[i][dimension])$ 
8:     end for
9:   end for
10:   $MaxDiff \leftarrow 0$ 
11:  for each dimension do  $MaxDiff \leftarrow MAX(MaxDiff, b2[dimension] - b1[dimension])$ 
12:  end for
13:  for each dimension do
14:     $CurDiff \leftarrow b2[dimension] - b1[dimension]$ 
15:     $Expand \leftarrow (MaxDiff - CurDiff)/2$ 
16:     $b2[dimension] \leftarrow b2[dimension] + Expand$ 
17:     $b1[dimension] \leftarrow b1[dimension] - Expand$ 
18:  end for  $\triangleright$  Now  $b1$  and  $b2$  encapsulate a hypercube where all  $points$  are in it.
19:   $BHTree \leftarrow BARNESHUTTREE(points, b1, b2)$ 
20:  for  $i \leftarrow 1$  to  $m$  do
21:     $means[i] \leftarrow centerOfMass$  from BarnesHutTree of node  $i$ 
22:     $ClosestPoints[i] \leftarrow$  points with the shortest distance to  $means[i]$ 
23:  end for
24:  for  $i \leftarrow 1$  to  $m$  do
25:     $hyperplanes[i] \leftarrow$  hyperplane that go through  $ClosestPoints[i]$  which is perpendicular to the
    line through  $ClosestPoints[i]$  and  $centerOfMass$  of local cluster
26:  end for
27:  for  $i \leftarrow 1$  to  $m$  do
28:     $PartialTrees[i] \leftarrow GETPARTIALTREE(hyperplane[i], BHTree)$ 
29:  end for
30:  for  $i \leftarrow 1$  to  $m$  do
31:    if  $i \neq NodeIndex$  then
32:       $BHTrees[i] \leftarrow PartialTrees[NodeIndex]$  from node  $i$ 
33:    end if
34:  end for
35:   $BHTrees[NodeIndex] \leftarrow BHTree$ 
36:  for  $i \leftarrow 1$  to  $LEN(points)$  do
37:     $acceleration \leftarrow 0$ 
38:    for  $j \leftarrow 1$  to  $LEN(BHTrees)$  do
39:       $acceleration \leftarrow acceleration + GETACCELERATION(points[i], BHTrees[j])$ 
40:    end for
41:     $velocities[i] \leftarrow velocities[i] + acceleration * DeltaTime$ 
42:     $points[i] \leftarrow points[i] + velocities[i] * DeltaTime$ 
43:  end for
44: end function
```

Algorithm 5

```
1: function BARNESHUTTREE(points, b1, b2)
2:   if All points are identical then
3:     NUMCHILDREN(BHTreeNode)  $\leftarrow$  0
4:     MASS(BHTreeNode)  $\leftarrow$  LEN(points)
5:     CENTEROFMASS(BHTreeNode)  $\leftarrow$  points[1]
6:   else
7:     CHILDREN(BHTreeNode)  $\leftarrow$  empty array
8:     for nb1, nb2  $\leftarrow$  every hyperoctant of the hypercube b1, b2 represent do
9:       npoints  $\leftarrow$  empty array
10:      for point  $\leftarrow$  points do
11:        if point in hypercube represented by nb1, nb2 then PUSHBACK(npoints, point)
12:      end if
13:    end for
14:    if LEN(npoints) > 0 then
15:      PUSHBACK(CHILDREN(BHTreeNode), BARNESHUTTREE(npoints, nb1, nb2))
16:    end if
17:    NUMCHILDREN(BHTreeNode)  $\leftarrow$  LEN(CHILDREN(BHTreeNode))
18:    CENTEROFMASS(BHTreeNode)  $\leftarrow$  0
19:    MASS(BHTreeNode)  $\leftarrow$  0
20:    for child  $\leftarrow$  each CHILDREN(BHTreeNode) do
21:      ChildMass  $\leftarrow$  CENTEROFMASS(child) * MASS(child)
22:      CENTEROFMASS(BHTreeNode)  $\leftarrow$  CENTEROFMASS(BHTreeNode) + ChildMass
23:      MASS(BHTreeNode)  $\leftarrow$  MASS(BHTreeNode) + MASS(child)
24:    end for
25:    CENTEROFMASS(BHTreeNode)  $\leftarrow$   $\frac{\text{CENTEROFMASS}(\text{BHTreeNode})}{\text{MASS}(\text{BHTreeNode})}$ 
26:  end for
27: end if
28: B1(BHTreeNode)  $\leftarrow$  b1
29: B2(BHTreeNode)  $\leftarrow$  b2
30: return BHTreeNode
31: end function
```

Algorithm 6

```
1: function GETPARTIALTREE(hyperplane, BHTreeNode)
2:   CENTEROFMASS(PartialNode)  $\leftarrow$  CENTEROFMASS(BHTreeNode)
3:   MASS(PartialNode)  $\leftarrow$  MASS(BHTreeNode)
4:   B1(PartialNode)  $\leftarrow$  B1(BHTreeNode)
5:   B2(PartialNode)  $\leftarrow$  B2(BHTreeNode)
6:   if  $\frac{B1(\text{BHTreeNode}) - B2(\text{BHTreeNode})}{\text{DISTANCE}(\text{hyperplane}, \text{CENTEROFMASS}(\text{BHTreeNode}))} < \theta$  then
7:     NUMCHILDREN(PartialNode)  $\leftarrow$  0
8:   else
9:     NUMCHILDREN(PartialNode)  $\leftarrow$  NUMCHILDREN(BHTreeNode)
10:    CHILDREN(PartialNode)  $\leftarrow$  empty array
11:    for child  $\leftarrow$  each CHILDREN(BHTreeNode) do
12:      PUSHBACK(CHILDREN(PartialNode), GETPARTIALTREE(hyperplane, child))
13:    end for
14:  end if
15:  return PartialNode
16: end function
```

Algorithm 7

```
1: function GETACCELERATION(point, BHTreeNode)
2:   acceleration  $\leftarrow$  0
3:   if  $\frac{B1(BHTreeNode) - B2(BHTreeNode)}{DISTANCE(point, CENTEROFMASS(BHTreeNode))} < \theta \vee \text{NUMCHILDREN}(BHTreeNode) = 0$  then
4:     com  $\leftarrow$  CENTEROFMASS(BHTreeNode)
5:     acceleration  $\leftarrow$  ACCELERATIONFROMCOM(point, com, MASS(BHTreeNode))
6:   else
7:     for child  $\leftarrow$  every CHILDREN(BHTreeNode) do
8:       acceleration  $\leftarrow$  acceleration + GETACCELERATION(point, child)
9:     end for
10:  end if
11:  return acceleration
12: end function
```

when combined with all the other overheads mentioned above, might outweigh the benefit of not having to send all points to all nodes.

In the future, it would be interesting to compare the performance of the current implementation against the naively parallelized Barnes-Hut algorithm described above.

Due to time constraint, the implementation is also not parallelized within a node. As a result, this does not take advantage of the multicore nature of each node. In the future, it'll be interesting to see how much extra performance can be extracted via parallelization within each node.

Methodology

To measure the performance of the implementation and whether parallelization sped up the simulation or not and by how much, the program was run on spartan with 1, 2, 4, 8, 16, 32 and 64 nodes. Because the implementation doesn't utilize OpenMP for inner process / inner node parallelization, each node only have one core allocated to run the simulation. Furthermore, to investigate how the communication overhead between the different nodes effect the overall runtime, the program is also run on a single spartan compute node with 1, 2, 4, 8, 16, 32 and 64 threads. This should significantly reduce the communication overhead between the Open MPI processes and give us more insight to how the extra communication overhead between different spartan compute nodes effect the performance of the program.

To allow more fine-grained observation of the program's performance, the runtime measurement is done on each process for the following stages of the program:

- Generating the points (by sampling from the GMM components).
- Choosing centers as needed by the k-means++ algorithm.
- Collaborating to cluster points using k-means algorithm after the centers are chosen.
- Sending points across the nodes so that all points of each cluster are on a single node.
- The N-Body simulation.

Experiments

Analysis & Discussion