

kd-tree construction and analysis with OpenMP and OpenMPI

Imola Fodor SM3500474
Foundations of High Performance Computing
University of Trieste

Deadline 28.02.2022

Contents

1	Introduction	2
2	Algorithm	2
3	Implementation	4
3.1	OpenMP	4
3.2	Hybrid solution	6
4	Performance model and scaling	8
4.1	Hardware	8
4.2	Parallel speed-up	8
4.3	Strong scaling	9
4.4	Weak scaling	9
4.5	Thread equilibrium for Hybrid solution	10
5	Conclusion	11

1 Introduction

A K Dimensional tree (or k-d tree) is a tree data structure that is used to represent points with more than one property in a k-dimensional space. It is a convenient way to organize points by several criteria at once and it provides eg. a convenient way to search, cluster points by their overall similarity.

In this work a way to build and parallelize such a tree is presented.

2 Algorithm

The construction of the tree is done by:

- Finding median/pivot by median of medians of the input array. Each section of constant length is sorted by insertion sort.
- Recursively proceed on left and right portions on the left and right of the found median. Each median is a node.
- Terminate when length of a portion is 0.
- Return root node.

The time complexity of the divide and conquer algorithm is $O(n \log n)$ since the partition process always picks the middle element as pivot. The median of the medians is found in linear time, $O(n)$.

Time complexity for partitioning n datapoints:

$$T(n) = 2T(n/2) + \theta(n) \quad (1)$$

If the partitioning would be done in a non optimal way (choosing as pivot the lowest, largest element of the array), the complexity would be $O(n^2)$.

Algorithm 1 Build kD-tree

Input arrayOfNodes**Output** treeRootNode

```
1: function BUILDKDTree(startNode, length, axis, dim)

2:   if  $length = 0$  then                                     ▷ base case
3:     return 0
4:   end if

5:    $myaxis \leftarrow$  round robin approach between 0 and 1

6:    $medianNode \leftarrow$  MEDIANOFMEDIANS(startNode, startNode + length
    -1, myaxis, len)

7:    $medianNode.left \leftarrow$  BUILDKDTree(startNode, medianNode -
    startNode, myaxis, dim)
8:    $medianNode.right \leftarrow$  BUILDKDTree(startNode, startNode +
    length - (medianNode + 1), myaxis, dim)
9:   return treeNode
10: end function

11: function MEDIANOFMEDIANS(startNode, endNode, myaxis, length)
12:   if  $length < 10$  then                                     ▷ base case
13:     INSERTIONSORT( $startNode, length, myaxis$ )
14:      $median \leftarrow middleElement$ 
15:   else
16:     subarrays  $\leftarrow ceiling(n/5)$ 
17:     allocate array medians of length subarrays
18:     for  $i \leftarrow 1, subarrays$  do
19:       INSERTIONSORT( $startNode, length, myaxis$ )
20:        $medians[i] \leftarrow middleElement$ 
21:     end for

22:     if  $numSubarrays = high$  then
23:        $median \leftarrow$  MEDIANOFMEDIANS( $medians, end, myaxis, length$ )
24:     else
25:       INSERTIONSORT( $medians, num\_subarrays, myaxis$ )
26:        $median \leftarrow medians[middleElement]$ 
27:     end if
28:   end if
29:   return median
30: end function                                     3

31: procedure INSERTIONSORT(startNode, length, axis)
32:   similar to the sorting of playing cards in hands
33: end procedure
```

3 Implementation

Herein the strategy used for OpenMP, OpenMP and their hybrid solution.

3.1 OpenMP

OpenMP is one of the application programming interfaces that facilitates the employment of a shared memory paradigm for parallelization within a node. Below the simplified, decorated Algorithm1 with the instruction read by OpenMP during compilation.

The algorithm has not shown scaling by increasing the number of resources on 1M datapoints. The execution time reported below.

Results 1M datapoints

Number of threads	2	4	8	16	48	96
Time [s]	0.18	0.17	0.18	0.18	0.19	0.19

Algorithm 2 Build kD-tree w/ OpenMP

Input arrayOfNodes**Output** treeRootNode

```
1: function MAIN
2:   #pragma omp parallel
3:   #pragma omp single nowait
4:   initialize random arrayOfNodes
5:    $root \leftarrow \text{BUILDKDTREE}(\text{arrayOfNodes}, \text{length}, 0, 2)$ 
6:   #pragma omp barrier
7:   printroot
8: end function
9: function BUILDKDTREE(startNode, length, axis, dim)

10:   if  $\text{length} = 0$  then ▷ base case
11:     return 0
12:   end if

13:    $\text{myaxis} \leftarrow$  round robin approach between 0 and 1

14:    $\text{medianNode} \leftarrow \text{MEDIANOFMEDIANS}(\text{startNode}, \text{startNode} + \text{length}$ 
     $-1, \text{myaxis}, \text{len})$ 

15:   #pragma omp task
16:    $\text{medianNode.left} \leftarrow \text{BUILDKDTREE}(\text{leftPoints}, \text{length}, \text{myaxis}, \text{dim})$ 

17:   #pragma omp task
18:    $\text{medianNode.right} \leftarrow \text{BUILDKDTREE}(\text{rightPoints}, \text{length}, \text{myaxis},$ 
     $\text{dim})$ 
19:   return treeNode
20: end function

21: function MEDIANOFMEDIANS(startNode, endNode, myaxis, length)
22:   if  $\text{length} < 10$  then ▷ base case
23:      $\text{INSERTIONSORT}(\text{startNode}, \text{length}, \text{myaxis})$ 
24:      $\text{median} \leftarrow \text{middleElement}$ 
25:   else
26:     for  $i \leftarrow 1, \text{subarrays}$  do
27:       #pragma omp parallel for
28:        $\text{INSERTIONSORT}(\text{startNode}, \text{length}, \text{myaxis})$ 
29:        $\text{medians}[i] \leftarrow \text{middleElement}$ 
30:     end for
31:      $\text{INSERTIONSORT}(\text{medians}, \text{num\_subarrays}, \text{myaxis})$ 
32:      $\text{median} \leftarrow \text{medians}[\text{middleElement}]$ 
33:   end if
34:   return median
35: end function
```

Continuing ...

```
36: procedure INSERTIONSORT(startNode, length, axis)
37:     similar to the sorting of playing cards in hands
38: end procedure
```

3.2 Hybrid solution

The hybrid solution, other than leveraging parallelization within a node, leverages also the different nodes that reside in a cluster, scaling up the solution. To achieve this, point-to-point messaging needs to be designed along the NUMA node on disposal, since each node has its own memory space. Once a task receives its message, each one computes the assignment on its portion of data using OpenMP threads, as under subsection OpenMP.

In order for each task to get its portion of data, rank0, the master rank, gets to find the tree nodes up until the level of the kdtree from where it is possible to assign the unique, fairly balanced, and only portions of data for each task (the left-most chunk being processed by rank 0). This approach was decided upon, to ensure load balancing and to easily reconstruct the tree, by sending the subtrees back to rank0.

MPI calls are made inside parallel regions, but there is no restriction when a certain thread receives the message from the master rank, hence the *MPI_THREAD_MULTIPLE* paradigm is flagged.

Below the function that employs tasks, and the main function.

Algorithm 3 Build kD-tree w/ Hybrid

Input arrayOfNodes, numProcs, rank

Output treeRootNode

```
1: function MAIN
2:   MPI_Init_thread(..., MPI_THREAD_MULTIPLE, ...)
3:   #pragma omp parallel
4:   #pragma omp single nowait
5:   initialize random arrayOfNodes
6:   if rank = 0 then
7:     root  $\leftarrow$  FINDFIRSTNODES(arrayOfNodes, length, 0, 2,
                                depth=0, rank = -1)  $\triangleright$  rank is a pointer
8:   else
9:     MPI_Recv(length, 1, MPI_INT, 0, 2, MPI_COMM_WORLD, ...)
10:    MPI_Recv(portion, length, MPI_BYTE, 0, 0, MPI_COMM_WORLD, ...)

11:    toSend  $\leftarrow$  BUILDKD TREE(portion, length, myaxis, 2)  $\triangleright$ 
    ideally the head of the subtree would be sent back to rank 0, for full tree
    construct
12:  end if
13:  MPI_Finalize()
14: end function
15: function FINDFIRSTNODES(startNode, length, axis, dim)

16:   myaxis  $\leftarrow$  round robin approach between 0 and 1

17:   if depth ==  $\log_2(\text{numProcs})$  then  $\triangleright$  base case
18:     return 0
19:   end if

20:   medianNode  $\leftarrow$  MEDIANOFMEDIANS(startNode, startNode + length
    -1, myaxis, len)

21:   #pragma omp task
22:   if depth ==  $\log_2(\text{numProcs}) - 1$  then
23:     rank  $\leftarrow$  rank + 1  $\triangleright$  no round robin needed
24:     if rank = 0 then
25:       toSend  $\leftarrow$  BUILDKD TREE(leftPoints, leftLength, myaxis, 2)
26:     else
27:       MPI_Send(leftLength, 1, MPI_INT, rank, 2, MPI_COMM_WORLD)
28:       MPI_Send(leftPoints, leftLength, MPI_BYTE, rank, 0, MPI_COMM_WORLD)

29:     end if 7
30:     rank  $\leftarrow$  rank + 1  $\triangleright$  no round robin needed
31:     MPI_Send(rightLength, 1, MPI_INT, rank, 2, MPI_COMM_WORLD)
32:     MPI_Send(rightPoints, rightLength, MPI_BYTE, rank, 0, MPI_COMM_WORLD)

33:   end if
```

Continuing ...

```
34:   depth ← depth + 1
35:   medianNode.left ← FINDFIRSTNODES(leftPoints, length, myaxis,
    dim, depth, rank)
36:   #pragma omp task
37:   depth ← depth + 1
38:   medianNode.right ← FINDFIRSTNODES(rightPoints, length, myaxis,
    dim, depth, rank)
39:   return treeNode
40: end function
```

4 Performance model and scaling

Measurements for speedup and efficiency graphs are used to acquire an indication of how well the implementation is performing in regard to some reference implementation. As reference, the corresponding serial execution of the code is used, on the same hardware, namely, for 10k datapoints, the execution time is 0.38s.

4.1 Hardware

- CPU name: Intel(R) Xeon(R) Gold 6226 CPU @ 2.70GHz
- Sockets: 2
- Cores per socket: 12
- Threads per core: 2, HT on

4.2 Parallel speed-up

In order to compute the parallel speed-up of the implementation, the following formula is used:

$$S(P) = T(1)/T(P) \quad (2)$$

, where $T(1)$ is the serial execution time, and $T(P)$ is the parallel execution time of the same problem size, with P tasks.

4.3 Strong scaling

Parallel efficiency, also referred to as strong scaling, is calculated by:

$$E(P) = S(P)/P \quad (3)$$

When $S(P) = P$ it is considered a perfect speed-up. This in "real-life" is not an ultimate goal, since the code usually has parts that need to execute in a serial fashion. Ahmdal hence defined a way to model these more common, realistic implementations:

$$Sahm(P) = 1/(s + (p/N)) \quad (4)$$

In case of the above described implementation, the purely serial part (s) is the initialization of the array of nodes. Furthermore, a show-stopper to a full parallelization is the part where the master rank computes the nodes up until a certain level of the tree.

Results 10k datapoints

Number of procs	2	4	8	16	32
Time [s]	0.39	0.42	0.42	0.42	x

Example run:

`mpicc - fopenmp - lm - gkdtee_build.c - okdtee.x`

`mpirun -np 2 kdtree.x`

Start Time: 0.000000, End Time: 0.001544, Elapsed Time: 0.00154371

Start Time: 0.000000, End Time: 0.432807, Elapsed Time: 0.43280652

No further analysis has been made, since no improvement in time has been noticed as the number of procs increased.

4.4 Weak scaling

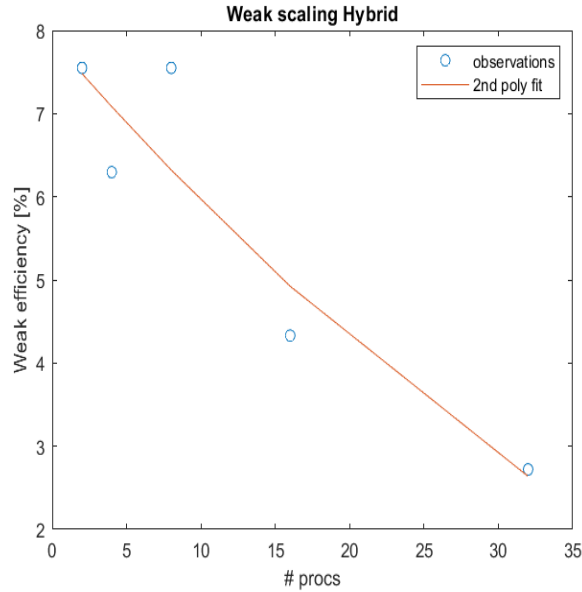
Gustafson instead pointed out, that in practice the sizes of problems scale with the amount of available resources. This is called weak scaling, where the scaled speedup is calculated based on the amount of work done for a scaled problem size (in contrast to Amdahls law which focuses on fixed problem size). Gustafson proposed that with increased resources, the serial part remains the same, doesn't increase, even if the problem size increases.

The implementation of this work does not exactly comply with the latter assumption, since the more tasks we can employ, the master rank will

have a bigger amount of levels (of the tree) to compute. Nevertheless, the weak scaling analysis is done, since the possible serial part is in any case, not utterly expensive, given the max amount of tasks we can employ.

$$S_{gus}(P) = s + pN \quad (5)$$

The weak efficiency can be calculated by eq.(3), changing the problem size with each change of resources. The problem size was increased by a factor of 10 (starting from $1k$), with each scale up in resources.



On *Figure4*, it can be noted that the efficiency starts at 8%, would confirm that the implementation is not very suitable for leveraging large resources.

4.5 Thread equilibrium for Hybrid solution

To understand better the optimal number of threads to be used for the hybrid approach, further analysis is needed. Using for the purpose the OpenMP function *omp_get_wtime*. It was noted that the program, to execute on $10k$ datapoints, took slightly lower time for 2 threads, $0.39s$ whileas for larger threads the time was $> 0.4s$.

5 Conclusion

The logic of the implementation in theory promised scaling with more tasks employed, but the timings extracted did not reflect the same.

The master rank slows the computation, since it may happen that before sending the chunks to other ranks, it computes his subtree creation task. This was addressed by programmatically directing the MPI calls to execute before the rank 0 would execute its task, but no improvement was shown in the performance.

Pinning has not been specified for potential increase in performance. Profiling to better understand the time each section takes neither.

OpenMPI strategy could have been used to find the pivot element. In that case, rank 0 would have scattered the sections to the ranks, the ranks would have computed the insertion sort, to then send back their medians to rank 0, rank 0 gathering them.