# NTNU
## Innovation and Creativity

# Seismic processing using Parallel 3D FMM

**Idar Borlaug**

Master of Science in Computer Science
Submission date:  June 2007
Supervisor:        Anne Cathrine Elster, IDI

Norwegian University of Science and Technology
Department of Computer and Information Science

Problem Description

Finding oil includes seismic processing of data obtained from experiments done in the field. The computations are generally very compute intensive and hence lends themselves to parallel computing, since often large datasets spanning storages on several compute nodes is involved.

In this thesis, a seismic application is considered and fast both serial and parallel algorithms studied and compared. This includes evaluating how to best divide the data between different nodes for a given application. In particular, methods for finding salt formations in seismic data using FMM, will be investigated.

Assignment given: 22. January 2007
Supervisor: Anne Cathrine Elster, IDI

# Abstract

This thesis develops and tests 3D Fast Marching Method (FMM) algorithm and apply these to seismic simulations. The FMM is a general method for monotonically advancing fronts, originally developed by Sethian. It calculates the first arrival time for an advancing front or wave. FMM methods are used for a variety of applications including, fatigue cracks in materials, lymph node segmentation in CT images, computing skeletons and centerlines in 3D objects and for finding salt formations in seismic data.

Finding salt formations in seismic data, is important for the oil industry. Oil often flows towards gaps in the soil below a salt formation. It is therefore, important to map the edges of the salt formation, for this the FMM can be used. This FMM creates a first arrival time map, which makes it easier to see the edges of the salt formation.

Herrmann developed a 3D parallel algorithm of the FMM testing waves of constant velocity. We implemented and tested his algorithm, but since seismic data typically causes a large variation of the velocities, optimizations were needed to make this algorithm scale. By optimising the border exchange and eliminating much of the roll backs, we deleveped and implemented a much improved 3D FMM which achieved close to theoretical performance, for up to at least 256 nodes on the current supercomputer at NTNU.

Other methods like, different domain decompositions for better load balancing and running more FMM picks simultaneous, will also be discussed.

# Acknowledgements

I would like to thank Dr. Anne C. Elster for being my advisor on this Master thesis, and for giving me valuable input throughout the project.

I would also like to thank all my friends and co students in room ITV-458 for helping me with various problems underway.

Thanks to Tore Fevang at Schluberger for providing me with the motivation for this assignment, and for helping me get started and understanding the problem.

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

Finding salt formations in seismic data is very important for the oil industry. Salt formations are mushroom shaped. This creates a cavity below the mushroom head, where oil will flow from the surrounding rock, and be trapped. If one can map the boundaries of the salt formation, it will be made easier for drilling companies to drill in the correct place.

Seismic data are generated by sending low-frequency shock waves into the earth, and reading the wave reflection with sensors. To generate the shock waves, explosives or air guns can be used. The sensors are laid out in a grid, and these sensors listen for waves reflected from the different layers of rock. The amplitude of the reflected waves are then stored. This results in a grid map with wave amplitudes.

Finding salt in these data can be troublesome, because salt is a crystalline form. Therefore, it reflects the waves in many different directions, making it very difficult to read.

The Fast Marching Method (FMM) [11] is a general method for monotonically advancing fronts. It calculates the first arrival time for an advancing front or wave. This can be used to solve a number of problems, follow fatigue cracks in materials [2], lymph node segmentation in CT images [3], computing skeletons and centerlines in 3D objects [4] and for finding salt formations in seismic data.

To solve this problem, a paper was published last year, called Seismic event tracking by global path optimization[5], from Borwn et.al for Amerada Hess Corporation. This paper describes how they tried to map the boundaries of a salt formation. They used the Fast Marching Method for mapping the boundaries, however, they did not use the FMM for the entire domain. They took planes from the seismic 3D data, using Delaunay triangulation. Using these planes or slices from the data, they then tried to map the edges of the salt formations.

In this thesis, the Fast Marching Method will be extended to 3D space. Because solving the FMM for a 3D space will take much longer, a parallel version of the Fast Marching Method will be made.

## 1.1 Contributions

The following highlights the contributions of this master thesis:

- Implementing and testing the Parallel version of the Fast Marching Method developed by Hermann.

- Development of a new parallel Fast Marching Method algorithm, because it will give better results.

- Different domain decompositions have been looked at and tested. To see if they gave a significant difference in speed.

- New computational model for calculating the theoretical execution time for the new parallel Fast Marching Methods have been constructed. Also an speedup equation have been made for PFMM, which roughly indicate the speedup achievable with a given number of nodes.

- Two ideas for improving performance even more have been discussed in Future work.

## 1.2 Outline

- Chapter 1 is this introduction.

- Chapter 2 is background information explaining seismic surveys, super computer, especially Njord, the Fast Marching Method and different variations of the Fast Marching Method.

- Chapter 3 explains my solution to the problem of finding salt formations in seismic data. It goes through the different choices and explains the two different applications that resulted from two parallel variations of the Fast Marching Method

- Chapter 4 describes a model of the runtime of each application

- Chapter 5 shows the different benchmarks and discusses for bout applications

- Chapter 6 is the conclusion

- Bibliography is a list of citations in this thesis, and a description of them. Links are included where available.

# Chapter 2

# Background Material

Finding salt formations in seismic data, explained in Section 2.1, can be difficult. [5] tried to find a solution for this in their article Seismic event tracking by global path optimization discussed in Section 2.2. The Fast Marching Method is a method used in [5] to locate the salt formations. The Fast Marching Method is discussed in Section 2.6. [5] only solved the problem for 2D planes, to solve the whole 3D domain, one needs a lot more computations. To solve that issue, one needs to solve the problem in parallel on many computers. Herrmann made a parallel version of the Fast Marching Method, discussed in Section 2.8. Normal computer could suffice when solving this problem, but it would be faster to run it on a Cluster super computer, discussed in Section 2.3. All the results in this thesis have been created while the program runs at the super computer Njord, discussed in Section 2.4.

## 2.1   Seismic data

One often use seismic data for locating oil reservoirs in the earth. To generate shock waves that traverse the earth they often use explosives or air-guns. Explosives are mostly used on land while air-guns are used on ships. To collect the data from the shock waves a grid with sensors is used. On land, this grid is laid out over the targeted area. On sea, it's pulled after a vessel with several kilometer long cables, carrying the sensors. When the shock waves are reflected back from the earth, the signals are gathered by the sensors and stored. These stored signals are mapped onto a 3D map that models the earth. Out of this 3D map, it's possible to read what lies beneath the surface.

Oil is often located in open spaces created by salt formations. A salt formation looks like a mushroom. On the edges below the top, open space are often formed. In this space, oil can flow from the surrounding rock, making huge areas of easily available oil. Therefore, it's important to find these salt formations.

## 2.2   Seismic event tracking by global path optimization

This is a summary of what is discussed in the article with the same name[5].

Tracking salt boundaries in seismic data is very difficult. It has gaps and other interferences, which makes it difficult for most trackers to function.  Most trackers look a few steps ahead (local optimizations), which makes them more prone to go down gaps in the data. An "ideal" auto tracker would:

- track using a global measure of optimality

- intelligently traverse "holes" in the salt events

- naturally track multi-valued salt boundaries.

To make this work, a positive velocity field is needed. To do this with seismic data, a number to the data can be added, so everything becomes positive. To track the shortest distance between two points one will use the FMM and add the two different points.  This will highlight the shortest path between those points.

To do this they implemented the FMM in a 2D space by taking 2D slices from the 3D space and running the FMM on them. To make this a 3D solution they used Delaunay triangulation [6].  This made for a quick 3D algorithm when new picks were chosen, because not all triangles needs to be updated with each new pick. A pick is a point chosen by a geologist. He picks two or three picks then run the FMM for each pick and add the results.

## 2.3   Cluster super computers

Cluster super computers are typically many computers or nodes that are connected.  These nodes/computers will work on problems, e.g. mathematical problems, and solve them fast. For using such computers, one needs a way to interact with the other computers.

Message passing is a method for passing messages between the same program, running on different nodes in a super computer cluster. This can be messages like "hi I am done", "I got 10 as a result" or "you need this data". Messages can be small or large, they can be a huge array or a flag saying the program is in a specific state.

A well known specification for message passing is MPI, Message Passing Interface [7].  This is a library for message passing, which programs, that run on super computer clusters, can use. It has a variety of different communication styles implemented, that one might need when programming for super

computer clusters. This is normal send/receive operations and more complex operations like broadcast- and reduce operations. It has also different forms of send/receive communication, synchronous, blocking and non-blocking, ready mode, etc. All these are useful for different scenarios.

Different mathematical problems are often solved on super computer clusters. Many of these use a grid or array, describing the problem. When solving a problem on a large array, it's normal to distribute a part of this array to each node in the cluster. This way, every node can solve a small part of the problem, resulting in a speed increase. To solve only one part of such a problem one normally needs data from the nodes that have the neighbouring data. This is often solved by using a border around the node's data, containing some of the neighbours data. In Figure 2.1 the border cells are marked by the nodes number. The border cells are often exchanged between the nodes a number of times when a problem is solved. This will involve communication with four nodes in the super computer cluster.

|   | 1 | 1 | 1 | 1 |   |
|---|---|---|---|---|---|
| 3 | 4 | 4 | 4 | 4 | 5 |
| 3 | 4 | 4 | 4 | 4 | 5 |
| 3 | 4 | 4 | 4 | 4 | 5 |
| 3 | 4 | 4 | 4 | 4 | 5 |
|   | 7 | 7 | 7 | 7 |   |

Figure 2.1: Border cells

## 2.4   Njord

Njord [8] is the current super computer at the Norwegian University of Science and Technology NTNU, delivered by IBM. It is a computer which uses SMP nodes and distributed memory between the nodes. Each node has 8 CPUs with 2 cores each. Each core can run two threads at the same time (SMT). However the threads does share some vital resources like the floating point unit. Vecause of this its not default to run two threads on each CPU on Njord. Each node can run 16 thread/processes at full speed. There are 62 nodes on Njord, making it possible to run jobs with 992 processes in theory. Some nodes are reserved for special programs. The largest size job that can run is 864 processes. Each node with 16 processes have enough memory for 832 MB for each proses. The CPUs used are IBM Power5+ [9], with 36 MB level 3 cache. The interconnect between the processes on each node is shared memory, between the nodes its a very fast network interconnect.

## 2.5 Manhattan distance

Manhattan distance [10] is a concept coming from Manhattan in New York. Because Manhattan is only square blocks with roads in a grid around them, the distance a cab driver has to take from one point to another in Manhattan is the Manhattan distance. In Figure 2.2 the blue, yellow and red lines describe the Manhattan distance. The green line is the Euclidian distance, which is the shortest distance between two points.



Figure 2.2: Manhattan distance

## 2.6 Fast Marching Method (FMM)

The Fast Marching Method is a solution for the eikonal equation.

$$|\Delta u(x)| = F(x), x\epsilon\Omega \tag{2.1}$$

The eikonal equation gives you a travel time field, with first arrival times at all points within the solution. Sethian[11] has developed a fast marching method that solves this equation in a grid.

Fast Marching Method is a method that solves the travel time field without moving over each point more than once, a one pass algorithm. This is achieved by making a narrowband around the starting point or start structures, and moving this narrowband outwards one point at a time.

In Figure 2.3 you can see the narrowband represented by 0, where the already calculated points being marked as -1. Those which are not modified are marked at 1. For not using numbers later i will mark them as the following. Those outside the narrowband will be marked as OUTSIDE, those on the narrowband as BAND and those inside as KNOWN.

To solve Equation 2.1 correctly, the gradient operator has to be approximated by upwind, entropy-satisfying finite differences [12]. The approximation most often used is from [13].

Figure 2.3: Narrowband

$$[max(D_{ijk}^{-x}G, -D_{ijk}^{+x}G, 0)^2 +$$
$$max(D_{ijk}^{-y}G, -D_{ijk}^{+y}G, 0)^2 + \qquad (2.2)$$
$$max(D_{ijk}^{-z}G, -D_{ijk}^{+z}G, 0)^2]^{1/2} = 1.$$

Where

$$
\begin{aligned}
D_{ijk}^{-x}G &= \frac{G_{ijk} - G_{i-1jk}}{\triangle x}, D_{ijk}^{+x}G = \frac{G_{i+1jk} - G_{ijk}}{\triangle x} \\
D_{ijk}^{-y}G &= \frac{G_{ijk} - G_{ij-1k}}{\triangle y}, D_{ijk}^{+y}G = \frac{G_{ij+1k} - G_{ijk}}{\triangle y} \qquad (2.3) \\
D_{ijk}^{-z}G &= \frac{G_{ijk} - G_{ijk-1}}{\triangle z}, D_{ijk}^{+z}G = \frac{G_{ijk+1} - G_{ijk}}{\triangle z}
\end{aligned}
$$

Where G is the arrival time matrix. The simple solution to this problem is to iteratively update all nodes within the array after Equation 2.2 until it finds a stable solution. Because this equation has an upwind property, each point in the array is only dependent on its smaller neighbours, and one can use the faster FMM algorithm.

The fast marching method solves the travel time to different point in a grid. It moves a narrowband outwards from a starting point, which can be a single point or multiple points that are connected. It basically consists of a few simple steps, shown in Figure 2.4.

## 2.6.1 FMM Initial data

Make two arrays, one for storing the band information and one for storing travel times. Set the initial point or points to KNOWN. Add the adjacent points to the narrowband and mark them BAND. The rest must be marked as OUT-SIDE.

1. set up initial data
2. loop begin:  extract the point from the narrowband with lowest travel
   time
3. mark the extracted point as KNOWN
4. add neighbours not in narrowband or KNOWN to narrowband.
5. recalculate all adjacent nodes by Equation 2.2
6. loop end

Figure 2.4: Serial FMM algorithm

### 2.6.2  FMM main loop

First, one needs to extract the lowest travel time from the narrowband.  The
narrowband can be stored efficiently in a min sorted heap[14]. Set this point to
KNOWN. This is now removed from the narrowband and will never be added.
Add adjacent points to the narrowband and calculate their arrival times ac-
cording to Equation 2.2. Repeat the loop until all nodes are marked KNOWN.

| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 2  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 2  | 1  | 2  | -1 | -1 | -1 |
| -1 | -1 | 2  | 1  | 0  | 1  | 2  | -1 | -1 |
| -1 | -1 | -1 | 2  | 1  | 2  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 2  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

Figure 2.5: Traveltime array step n

An ilustration of one loop step can be seen in Figure 2.5 and Figure 2.6.  In
the first figure all nodes marked 2 are in the narrowband. In the second figure
a new step has been made.  The top 2 value has been marked KNOWN and
removed from the narrowband.  Its neighbouring nodes have been added to
the narrowband and its values calculated.  In the next step the lowest value
from the narrowband would be choosen. That would be one of the remaining
2 points.

## 2.7   Spherical vs Cartesian coordinates

To calculate the arrival times, one can use two different coordinate systems,
the Cartesian xyz and Spherical $r\theta\phi$. Where $r$ is the length from origin and $\theta$

| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
|----|----|----|----|----|----|----|----|----|
| -1 | -1 | -1 | -1 | 3  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 3  | 2  | 3  | -1 | -1 | -1 |
| -1 | -1 | -1 | 2  | 1  | 2  | -1 | -1 | -1 |
| -1 | -1 | 2  | 1  | 0  | 1  | 2  | -1 | -1 |
| -1 | -1 | -1 | 2  | 1  | 2  | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 2  | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

Figure 2.6: Traveltime array step n+1

is the angel between the x axis and the $r$ line and $\phi$ is the angel between z-axis and $r$ along the vertical plane.

To calculate arrival time at a specific point, one can use the 1 order approximation 2.4.

$$
\begin{aligned}
t_{ijk} &= t_{i-1jk} + \frac{\triangle x}{v}, \, t_{ijk} = t_{i+1jk} + \frac{\triangle x}{v} \\
t_{ijk} &= t_{ij-1k} + \frac{\triangle y}{v}, \, t_{ijk} = t_{ij+1k} + \frac{\triangle y}{v} \\
t_{ijk} &= t_{ijk-1} + \frac{\triangle z}{v}, \, t_{ijk} = t_{ijk+1} + \frac{\triangle z}{v}
\end{aligned}
\tag{2.4}
$$

Where v is the velocity at point (i,j,k) and t is the arrival time array.

Using Equation 2.4 will give large errors on sparse grid point configurations, because it has problems with waves propagating at 45 degrees. According to [1] this will result in a 20% error near a point source.



Figure 2.7: Cartesian coordinate system. Figure from [1]

If you have three points A B C where A(0,0) is the center, and B(1,0) is one point away on the x axis and C(0,1) is one point away on the y axis. Both B and C have arrival time 1.See Figure 2.7. When the FMM calculates Bs neighbour D (1,1) it will give it arrival time 2 according to Equation 2.4. But is should have $1 + 1/\sqrt{2}$ which will result in the 20% error. One solution to this is too use higher order approximations, seeing further than one point in the grid. Another is to use Spherical coordinates.

$$
[max(D_{ijk}^{-r}t, -D_{ijk}^{+r}t, 0)^2 +
$$
$$
max(D_{ijk}^{-\theta}t, -D_{ijk}^{+\theta}t, 0)^2 + \tag{2.5}
$$
$$
max(D_{ijk}^{-\phi}t, -D_{ijk}^{+\phi}t, 0)^2]^{1/2} = 1
$$

where

$$
D_{ijk}^{-r}t = \frac{t_{i,j,k} - t_{i-1,j,k}}{\triangle r}, D_{ijk}^{+r}t = \frac{t_{i+1,j,k} - t_{i,j,k}}{\triangle r} \tag{2.6}
$$

$$
D_{ijk}^{-\theta}t = \frac{t_{i,j,k} - t_{i,j-1,k}}{r\triangle\theta}, D_{ijk}^{+\theta}t = \frac{t_{i,j+1,k} - t_{i,j,k}}{r\triangle\theta} \tag{2.7}
$$

$$
D_{ijk}^{-\phi}t = \frac{t_{i,j,k} - t_{i,j-1,k}}{rsin\theta\triangle\phi}, D_{ijk}^{+\phi}t = \frac{t_{i,j-1,k} - t_{i,j,k}}{rsin\theta\triangle\phi} \tag{2.8}
$$

For using Spherical coordinates Equation 2.2 must be modified. The modified version is Equation 2.5. To use this function, all values for $\theta$ and $\phi$ where $r$ is 0 are set, and added to the narrowband. Unlike the Cartesian version, the heap tends to be stable in Spherical coordinates, because one usually subtracts one and adds one, when in Cartesian coordinates the heap can become pretty large, especially in 3D grids.

[1] concludes that spherical coordinate systems give more accurate results. Especially with point sources where Cartesian coordinates give a high degree of error. The spherical solution is generally as fast as the Cartesian version, but it will not give the best results for head waves.

## 2.8   Parallel Fast Marching Method

M. Herrmann [15] has written an article on domain decomposition parallelization of the Fast Marching Method. This thesis gives a short brief of some of the different methods he investigated.

There are a few problems with the Fast Marching Method when it comes to parallelization. It has a very serial nature. When the narrowband is moved, one needs to find the lowest travel time point, which will be difficult on a parallel version. How the narrowband will move is also difficult to know before hand, therefore it will not be completely straight forward how to divide the matrix between the nodes. It can even move in a spiral shape.

1. Perform step 1 of serial algorithm 2.4
2. locate the local minimum value in the narrowband.
3. Find the global minimum value by exchanging local minimums.
4. Perform step 3-5 of serial algorithm on the node with global minimum value.
5. If global minimum is border value, exchange with neighbouring node.
6. return to step 2 until all values are calculated.

Figure 2.8: Parallel Fast Marching Method 1.

## 2.8.1   Herrmann's Parallel algorithm 1

When making a parallel version of the Fast Marching Method, it's advantageous to split the domain between the computational nodes. Then a problem arise when one needs to find the smallest value in the narrowband. The smallest value can only be on one node (or a few if there are many equal values). The straight forward solution for this is to calculate the local minimum and use an all reduce function to find which node has the global minimum. When one calculates the border nodes, it's also important to send these changes to the other nodes. This solution results in algorithm 2.8.

This algorithm has an inherit serial part, still only one node can calculate at the same time. One can work on much larger datasets compared to the serial algorithm because the dataset is divided between the different nodes. To calculate the global minimum, one can use an allreduce min function, but it will still contain a global communication point for each node in the dataset.

## 2.8.2   Herrmann's Parallel algorithm 2

The next logical step is to get all nodes to work at the same time. This can be achieved if each nodes propagating narrowband will not interfere with the other nodes. The problem comes if a node receives a border value which will give lower arrival time then what's already calculated, then that node needs to roll back to an earlier state. Herrmann's first attempt at solving this problem resulted in algorithm 2.9

There are a few drawbacks for algorithm 2.9. First, one needs to store each state, which will take a huge amount of storage space. This must be done to be able to roll back. Secondly almost all exchanges to border values will result in a rollback. These problems are corrected in the last algorithm.

## 2.8.3   Herrman's Parallel algorithm 3

The last algorithm will try to correct the problems encountered in algorithm 2. Because of the attributes of Equation 2.2, a point that is larger than the new

1. Perform step 1 of serial algorithm 2.4
2. Check if you received a border value with lower value than the highest in your grid. If so, roll back to a state where the new value is the highest. Add this value to narrowband.
3. locate the locally smallest value in narrowband including new border values.
4. Perform step 3-5 of serial algorithm on the node with local minimum value.
5. If local minimum is border value, exchange with neighbouring node.
6. Store current state
7. return to step 2 until all value are calculated.
8. Wait until all nodes are finished or a new border value is received. Go to step 2.

Figure 2.9: Parallel Fast Marching Method 2.

1. Perform step 1 of serial algorithm 2.4
2. Check if you received a border value with lower value than the highest in your grid. If so, roll back to a state where the new value is the highest. By marking all points with higher value as BAND, add those values to narrowband.
3. locate the locally smallest value in narrowband, including new border values.
4. Perform step 3-5 of serial algorithm on the node with local minimum value.
5. If local minimum is border value, exchange with neighbouring node.
6. return to step 2 until all values are calculated.
7. Wait until all nodes are finished or a new border value is received. Go to step 2.

Figure 2.10: Parallel Fast Marching Method 3.

border value can be set back to BAND and retain its value[15]. So there is no reason to save the complete state, one can only rollback to BAND. This resulted in algorithm 2.10

This new algorithm 2.10 makes it possible to do a roll back without storing each state, one only needs to store the original array. The performance of this algorithm will change greatly depending on how many border exchanges must be made, and the number of rollback operations. The performance will vary a lot, depending on the border exchanges and rollbacks. In some situations it can be very fast, but in others it will require a lot of border exchanges and rollbacks, making it slow. It will still be faster than both algorithm 2.8 and algorithm 2.9.

### 2.8.4   Domain decomposition

The domain decomposition will have a great effect on how fast certain problems will be solved. There are a few ways one can divide the domain, the one most often taken is to divide them by minimising borders. Using a rectangular shape 2:1 it will give the least border area. Another approach is to divide the domain so that all nodes touch the center. This approach is taken by Herrmann [15] in his example.

| 1 | 1 | 2 | 2 | 3 | 3 |
|---|---|---|---|---|---|
| 1 | 1 | 2 | 2 | 3 | 3 |
| 4 | 4 | 5 | 5 | 6 | 6 |
| 4 | 4 | 5 | 5 | 6 | 6 |
| 7 | 7 | 8 | 8 | 9 | 9 |
| 7 | 7 | 8 | 8 | 9 | 9 |

Figure 2.11: Non optimal domain decomposition

Figure 2.11 is a domain decomposition chosen by Herrmann[15] to illustrate that this is not an optimal decomposition. It will not minimise borders and all nodes will not work at the same time. This decomposition can easily be scaled to 27 nodes for 3D space by adding nodes in the z axis.

In Figure 2.12 each node are connected to the middle. If the starting point is in the middle and the velocity field is all 1, then this would result in a circle expanding from the middle with equal work on each node. There would also be no need to do any communication between the nodes. Herrmann used this optimal layout for some of his tests they yielded 0.98 efficiency[15]. The problem is that very few real world cases map to this division. If the velocity field is very varied, one node could end up doing a lot of work. In 2D space this only scales to 4 nodes, and in 3D only 8 nodes. In the optimal case it works very well.

### 2.8.5   Herrmann's results

Hermann showed some graphs, illustrating how each of these domain decompositions would scale. Not surprisingly, the quadratic decomposition worked very well when the start structure was a sphere. This resulted in 0.96 efficiency. The non optimal decomposition didn't do it very well. This is because it will

| 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 |
| 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 |
| 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 |
| 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 |

Figure 2.12: Quadratic grid decomposition

not make all nodes work at the same time, and it is prone to more rollbacks, because of more communication.

# Chapter 3

# Parallel implementation of FMM

Finding salt formations in seismic data, is a challenge. Since this solution is created for solving a 3D domain, it is necessary to use more than one CPU. Therefore, it's important to make a parallel solver. This implicates that one need to change the Fast Marching Method, so it will be parallel. This is not straight forward since the algorithm is inherently serial, only the lowest value in the narrowband can be calculated. In this chapter two solutions are discussed, the second was created because of shortcomings in the first. They are hereafter known as HFMM and PFMM. As HFMM is bassed on Herrmanns algorithm and PFMM is the new algorithm described in this thesis.

## 3.1   Fast Marching Method

The Fast Marching Method is a very fast method for finding arrival times in a velocity field. It is a one pass algorithm and faster than one pass is very hard to make. Unfortunately, it uses a heap for finding the lowest value on the narrowband, giving it $O(\log N)$ runtime. This results in a final algorithm of $O(N \log N)$, which is a fairly fast one. The problem is that it doesn't easily parallelize. With the seismic datasets there is rarely enough ram to run the algorithm on one node. Not to mention the time it will take an $N \log N$ algorithm to pass many gigabytes of data. Therefore, it's important to use a parallel version of the Fast Marching Method.

### 3.1.1   Choosing domain decomposition

Since the seismic datasets are so large, a domain decomposition is essential. There are already two discussed by Herrmann in Section 2.8.4. The quadratic decomposition is the one that will give the best results in the optimal case that Herrmann tried. But there is also another decomposition that Herrmann didn't discuss, rectangular 1:2 scale decomposition. If you divide all your data into 1:2 rectangles you will get the lowest number of border cells.

$$w * C * 2 + h * C * 2 = TC \qquad (3.1)$$

Where w = width, h = height, C = cost per unit and TC = Total Cost.

The minim total cost for Equation 3.1 is when $w = 2h$. This will give the least border size for a given 2D domain. Using this as a domain decomposition, one would end up with the rectangular 1:2 scale decomposition.

I will therefore try to see which of these two domain decompositions yields the best results.

### 3.1.2   Rectangular 1:2 scale decomposition

| 0 | 0 | 1 | 1 | 2 | 2 | 3 | 3 |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 1 | 2 | 2 | 3 | 3 |
| 0 | 0 | 1 | 1 | 2 | 2 | 3 | 3 |
| 0 | 0 | 1 | 1 | 2 | 2 | 3 | 3 |
| 4 | 4 | 5 | 5 | 6 | 6 | 7 | 7 |
| 4 | 4 | 5 | 5 | 6 | 6 | 7 | 7 |
| 4 | 4 | 5 | 5 | 6 | 6 | 7 | 7 |
| 4 | 4 | 5 | 5 | 6 | 6 | 7 | 7 |

Figure 3.1: Rectangular 1:2 scale decomposition

In Figure 3.1 each node has a 1:2 rectangular data area, and the node also map to 1:2 rectangular shape. This will minimise the border cells. Such a grid will be very optimal for iterative solutions where one can run one iteration, then exchange borders and run another. Because each node has the smallest possible border size. However, in a FMM this division might not be optimal because one cannot guarantee that all nodes are on the narrowband. The narrowband could only be using 2 or 3 nodes, then the others would just sit there and wait for data. In seismic data it is, however, not that easy to predict how the narrowband will travel. Another advantage with this model is that it will scale, one can easily add nodes and it's easy to create a minimum border layout. The model will also convert to 3D space by just adding another layer of nodes in the z direction.

### 3.1.3    Choosing a parallel FMM

There are basically two parallel algorithms that one can choose from algorithm 2.10, which is the one Herrmann got the best results from, and iterative domain update according to Equation 2.2. The domain update will use many iterations before it will reach a stable solution.  It is, however, very easy to parallelize, since it's a iterative numerical approach.  But it will not give close to the wall clock time given by algorithm 2.10.  Since the algorithm 2.10 is a one pass algorithm.  If algorithm 2.10 gives no rollbacks.  There is therefore no reason to believe that a iterative approach will beat a one pass with a few rollbacks. The other algorithms presented by Herrmann is worse than algorithm 2.10 so they have not been tested.

I also wanted to test an idea of my own, it's an algorithm that is a cross between Herrmann's 2.10 algorithm and the domain update version.  This is discussed later in Section 3.3.

### 3.1.4    Choosing Spherical vs Cartesian coordinates

I did not choose to use spherical or higher order Cartesian equations.  The reason for this was that there was not enough time and the important part of this assignment is to find a viable parallel algorithm.  If a good parallel FMM is found its very easy to exchange the equation used for finding the new arrival times in each point later.  The first order Cartesian approximation is also very well tested and gives good enough results for most purposes, it's also good enough to check that the application calculates the times correctly. I used Equation 2.4 in my application.

## 3.2    HFMM

In this section, I will explain how my implementation of algorithm 2.10 works, and some minor modifications to the algorithm to avoid some issues that kills performance.

### 3.2.1    Parallel Fast Marching Method

The algorithm I used is an implementation of algorithm 2.10 from Herrmann. I made some minor changes, more implementation specific then actually altering the algorithm. The algorithm is explained step by step in Figure 3.2.

1. Setup initial data
2. Add initial point and set that to BAND
3. (a) Extract lowest value from narrowband and set the point to KNOWN
   (b) Calculate new arrival times for all neighbouring nodes that are not KNOWN
   (c) Add neighbouring nodes not in narrowband or KNOWN to narrowband
   (d) If the extracted point is a border point, send it to neighbours
   (e) check for incoming border values, add those to array and roll back if necessary
   (f) go to step (a) until narrowband is empty
4. Check for incoming border changes, add those and roll back if necessary
5. Check if all nodes are finished, then exit or go to step 3

Figure 3.2: Parallel Fast Marching Method, My version

## 3.2.2  Initial data

There is quite a lot of data that need to be allocated for the FMM. You need a velocity array, this is the data you get from seismic readings. The velocity array / seismic data is the one property thats sets your memory requirements.

The velocity aray is floating point numbers, saved as a float array in my application.

The second array you need is arrival times, where your store the result from your application. This is also a float array.

You also need an int array to store which nodes are KNOWN, OUTSIDE and on the narrowband (BAND). I used to store KNOWN as N where n is number of loops. This was used to test a few optimizations. The optimal would be to use two bits to store each point.

The HEAP also needs a pointer array to store each point that belongs to the narrowband. This heap array need to be big enough to store the largest narrowband size in the entire run. It's not that easy to calculate because it can become very large if there are very varying velocity values. I created the array to the same size as velocity array, it is way too big but I didn't want to get into space problems on the heap, for better memory utilization a vector would be preferable.

I also used a last array to store what value was sent. So if my narrowband came across a value that was already sent, because of a rollback, it wouldn't send it again if it had the same value as last time it was sent. This array will take as much place as the inner boarder but I made it as big as the velocity array. Then I didn't have to make a function that translates from a border point to some point in the boarder array. It would have added more complexity, and only given less memory usage which was not a problem on my test machine.

|                  | My application     | Optimal                       |
|------------------|--------------------|-------------------------------|
| Velocity array   | ARRAYSIZE*4        | ARRAYSIZE*4                   |
| Arrival time array | ARRAYSIZE*4      | ARRAYSIZE*4                   |
| Narrow band array | ARRAYSIZE*4       | ARRAYSIZE*1/4                 |
| HEAP array       | ARRAYSIZE*4+filldata | filldata count*4 + filldata |
| Sent array       | ARRAYSIZE*4        | 0                             |

Table 3.1: Optimal memory usage and my applications memory usage

This results in the memory usage presented in Table 3.2.2. ARRAYSIZE is number of points $x \cdot y \cdot z$ of the velocity array, which setts the problem size.

The total memory usage will be according to Table 3.2.2 on my application (arraysize*4)*5 + heap filldata and optimal memory usage will be (arraysize*4)*2+1/4 + heap filldata count*4 + filldata. So my application uses around twice as much as the optimal case.

### 3.2.3   Narrowband

Using a min sorted heap [14] we would be able to store the narrowband in a very efficient way.  The heap has a few properties that make it ideal for this. It can insert and extract points in O(log N) time, which is very efficient.  For my application, we used the heap from Introduction to Algorithms [14] and changed it, so it became a min heap instead of a max heap, which is implemented in the book.  To store each point, a struct was used, containing x,y,z coordinates and the value. This worked out very well.

### 3.2.4   Calculation of arrival times

To calculate the new arrival times of a specific point I used equation 2.4.  This equation for calculating arrival times is well tested, as its the most straight forward equation.

### 3.2.5   Exchanging border values

The border exchange was synchronous, which gave correct results.  However, having a synchronous communication for each border cell made it very slow. Therefore, it was changed to asynchronous communication.  The program sends border values once they are done and it checks to see if there are any incoming border values and add those.  This mostly eliminates the network delay.  The communication for signaling the nodes when you are done were also made asynchronous.

1. Setup initial data
2. Add initial point and set that to BAND
3.  (a) Extract lowest value from narrowband and set the point to KNOWN
    (b) Calculate new arrival times for all neighbouring nodes that are not KNOWN
    (c) Add neighbouring nodes not in narrowband or KNOWN to narrowband
    (d) go to step (a) until narrowband is empty
4. Exchange borders, reset band array
5. Calculate new values to the border's neighbours (inner border), if these values are smaller than those already there, add them to narrowband
6. Check if all nodes are finished, then exit or go to step 3

Figure 3.3: Parallel Fast Marching Method revised

To avoid unnecessary rollbacks, a few methods was employed. First store the largest value in the arrival time grid, so its fast to check whether the incoming border value is smaller than the largest value. Receive all the pending border values before checking if a rollback is needed, seemed to reduce rollbacks. If there are three incoming border values and two of them require a rollback, one is executed.

## 3.3   PFMM

Because HFMM didn't scale as well as expected, I altered the algorithm a bit, hoping to get better results. The altered/new algorithm is described in Figure 3.3.

Much of the algorithm is exactly the same. The initial data is mostly the same, however, this algorithm doesn't rely on loop numbering, so one can save some space in the band array. The FMM works just like it does in a serial version, but there are some differences.

Instead of sending border values all the time, the program waits until the FMM is done with all the values in one node. Then send all values to neighbouring nodes. This will make the number of network sends quite a lot less. There are no rollbacks, but the system might work a few iterations before it settles and a stable solution is achieved.

If we use the quadratic decomposition , shown in Figure 2.12, which is used by Herrmann, this system will settle in three iterations. Say, the start point is located on node 0. Then node 0 will work its serial FMM until it's done, then exchange values to node 1 and 2. In the second iteration, nodes 1 and 2 will run their serial FMM, node 1 will check its border and find no need to add any new values. When node 1 and 2 are done, node 3 will start its serial FMM in iteration 3. When node 3 is done, all nodes will check their borders and find no values that need updating and all will exit. This example will work if there are

no waves coming back to a node, e.g. in a velocity field with only one value. Given the example Herrmann made by using a sphere located at the center as starting condition, this algorithm will run 1 iteration, it will however need the delay in transferring borders and checking the borders before it can exit so you won't get 100% efficiency compared to the serial version.

If the velocity field is the same value all over, this algorithm will stabilise in the Manhattan distance from the starting point node to the farthest edge iterations. So, if you take the Manhattan distance from the starting node to the farthest edge, you get the number of iterations. However, if you use a velocity field that will give waves that turn back into another node, then you might want to use a few more iterations. Given a 64 node 3D space 4x4x4, the Manhattan distance from the corner to the other will be 8.

# Chapter 4

# Modeling of the FMM algorithm

In this chapter a model for each algorithm is developed. The model will try to give a theoretical equation for runtime.

## 4.1 HFMM

The HFMM algorithm has two important aspects when it comes to calculating the runtime, how far is it from the point to the next node and how many rollbacks will occur. I will now try to explain how one can calculate those aspects and come up with an equation for overall execution time.

$$
\begin{aligned}
arraysize &= n \cdot m \cdot k & (4.1) \\
bordersize &= n \cdot m \cdot 2 + n \cdot k \cdot 2 + m \cdot k \cdot 2 & (4.2)
\end{aligned}
$$

Where n is the length in along x axis, m is length along y axis and k is length along z axis.

*#rollback* is difficult to estimate. It can be affected by several factors. Since HFMM uses asynchronous communication, one node running faster than another can trigger rollbacks in the neighbouring node. The network delay will also affect how many rollbacks that will occur, late border values might trigger rollbacks. If there are several border values coming in at the same time, and many off them will require a rollback, then only one rollback will be executed. The number off rollbacks will vary depending on the system its ran on. Two important factors affect the number of rollbacks, network delay and variations in execution time from node to node.

The time it takes to execute rollbacks is easier to describe. It consists of two factors, *#rollback* and $T_{traverse}$. $T_{traverse}$ is the time it takes to go through the matrix and set the band array back to a previous state. $T_{traverse}$ can be show as

Equation 4.3. *mtime* is the time it takes to modify one point in an array. This results in Equation 4.4 which describes the time it takes to execute all rollbacks on one node. *#rollback* is given by three factors, how high the network delay is, difference in computational speed between the nodes and if the data set require a rollback. This can be estimated through a few runs on the same dataset. But will be affected by borderline and extra border sends, caused by the velocity field.

$$T_{traverse} = arraysize \cdot mtime \qquad (4.3)$$
$$T_{traverse} = n \cdot m \cdot k \cdot mtime$$

$$T_{rollback} = \#rollback \cdot T_{traverse} \qquad (4.4)$$
$$\#rollback = P(rollback) \cdot (bordersize + extraBorderSends)$$

The time it takes for computing a matrix is explained in Equation 4.5. It is affected by two values how long it takes to execute one loop in the FMM and the array size. $T_{loop}$ describes how long one loop takes. However, how long it takes for each loop is difficult to predict, it consists of a few steps, all neighbours that are OUTSIDE will be calculated and added to narrowband. Calculating the new values are a constant time operation, called *calcn*, but adding the points to the narrowband is a log(N) operation. Which means it will be difficult to predict, because N is not known until runtime. It's very difficult to calculate how long it takes to execute $T_{compute}$, but it will be the same for each time the algorithm is run on the same dataset.

$$T_{compute} = T_{loop} \cdot arraysize \qquad (4.5)$$
$$T_{loop} = \#CalcNeighbours \cdot calcn + log(sizeNarrowband)$$

Communication between nodes has a constant time and a not so constant value. The minimum communication time is the time it takes to send all border values and receive them. The difficult part is the extra communication it takes if a rollback occurs. A border value that hasn't changed will not be resent, but if a border value changes because of a rollback, it will be resent. Equation 4.6 show how communication is calculated. $\beta$ is the inverse bandwidth and $\alpha$ is the network delay. For running on njord on one node $\alpha$ is $1.995 \cdot 10^{-6}$ and $\beta$ is $2.143 \cdot 10^{-10}$.

$$T_{comm} = (bordersize + resendCount) \cdot 20 \cdot \beta + \alpha \qquad (4.6)$$
$$T_{comm} = ((n \cdot m \cdot 2 + n \cdot k \cdot 2 + m \cdot k \cdot 2) + resendCount) \cdot 20 \cdot \beta + \alpha$$

To calculate the time it takes for a node to finish, one needs to know when the wave will trigger the node to start. This is not easy to find out, it will not move equally fast in all directions, that is dependant on the velocity array. Where the start point is located is important, if its located close to the border the node besides it will start to work earlier than those on the other side. The time for the wave to hit a node is called $T_{wavetime}$.

The overall time it takes for a node to finish is dependent on a few values, first the $T_{wavetime}$ tells us when it starts. $T_{compute}$ will tell us how long it takes to finish computing that node, $T_{comm}$ the amount of time it takes to send and receive border values, $T_{rollback}$ how long the rollbacks will take. This gives Equation 4.7.

$$T_{node} = T_{wavetime} + T_{compute} + T_{comm} + T_{rollback} \tag{4.7}$$

The overall time it takes to finish the HFMM is given my $T_{overall}$. This is the time it takes for the slowest node to finish. There are a few things that are not accounted for. The time it takes to update all nodes of their running status uses asynchronous communication and is not in $T_{communicate}$. Mostly because it takes too little time to have any affect. $T_{overall}$ is shown in Equation 4.8. P is max number of processors.

$$T_{overall} = MAX(T_{node}0, T_{node}1, ..., T_{node}P - 1) \tag{4.8}$$

Since there are quite a few values that can't be determined before runtime, a few values that vary depending on system and a few values that vary from run to run, this makes it difficult to derive a theoretical model from Equation 4.8.

## 4.2 PFMM

PFMM has one aspect that really affects performance, that's the Manhattan distance from the node which contains the pick to the farthest node.

$T_{compute}$ is the exact same as Equation 4.5 used in HFMM. That means that also $T_{loop}$ is the same. But most of the similarities to HFMMs model ends there.

The border communication for PFMM is given by Equation 4.9. It is only affected by the array size. The border exchange for PFMM is smaller than $T_{comm}$ for HFMM. This is because in PFMM only the arrival time and band array are exchanged.

$$\begin{aligned} T_{comm} &= bordersize \cdot 8 \cdot \beta + \alpha \\ T_{comm} &= (n \cdot m \cdot 2 + n \cdot k \cdot 2 + m \cdot k \cdot 2) \cdot 8 \cdot \beta + \alpha \end{aligned} \tag{4.9}$$

After a border communication each node has to check if its necessary with a rollback. This operation check to see if any values on the border needs to be changed. The time for $T_{rollbackcheck}$ is given by Equation 4.10.

$$\begin{aligned} T_{rollbackcheck} &= T_{loop} \cdot bordersize \\ T_{rollbackcheck} &= T_{loop} \cdot (n \cdot m \cdot 2 + n \cdot k \cdot 2 + m \cdot k \cdot 2) \end{aligned} \tag{4.10}$$

Rollbacks is also an issue with PFMM the penalty for rollbacks are not that severe if the Manhattan distance for the given node is low. Given by Equation 4.11. The rollback will only take as long as $T_{node}$ in worst case, it will stop when nothing else needs to be changed , but if the Manhattan distance is below max one or more rollbacks will be hidden. $M$ is the Manhattan distance to the specific node and $MaxM$ is the max Manhattan distance, given by the distance from the node containing the starting point to the farthest node in the grid.

$$T_{rollback} = T_{node} \cdot (MAX(\#rollbacks0, \#rollbacks1, .., \#rollbacksP - 1) - MaxM + M) \tag{4.11}$$

The time a certain node needs to finish is given by Equation 4.12. See that the starting time is not contained inside this equation as it was for $T_{node}$ in HFMM. The main reason for $T_{node}$ time is the array size.

$$T_{node} = T_{compute} + T_{comm} + T_{rollbackcheck} \tag{4.12}$$

The overall time for PFMM is given by Equation 4.13. The Manhattan distance is the factor that greatly adjusts the $T_{overall}$.

$$T_{overall} = T_{node} \cdot ManhattanDistance + T_{rollback} \tag{4.13}$$

## 4.3 Validation of PFMM

To see if this function works in real life, i tested to see how the Equation corresponds between runs.

As a base for the calclulations the $160x160x160$ matrix was used. Its runtime for the serial code was 8,73 seconds.

$$
\begin{aligned}
n = m = k &= 160 \\
T_{overall} &= T_{node} \cdot ManhattanDistance + T_{rollback} \\
T_{overall} &= 8.73 \\
ManhattanDistance &= 1 \\
T_{rollback} &= 0 \\
8.73 &= T_{node} \cdot 1 + 0
\end{aligned}
$$

When running only one node, the Manhattan distance willbe 1. $T_{rollback}$ is set to 0 as there is no need for a roll back on the serial version.

$$
\begin{aligned}
T_{node} &= T_{compute} + T_{comm} + T_{rollbackcheck} \\
T_{comm} &= 0 \\
T_{rollbackcheck} &= T_{loop} \cdot bordersize \\
T_{rollbackcheck} &= T_{loop} \cdot 160^2 \cdot 6 \\
T_{compute} &= T_{loop} \cdot arraysize \\
T_{compute} &= T_{loop} \cdot 160^3 \\
T_{node} &= T_{loop} \cdot 160^3 + 0 + T_{loop} \cdot 160^2 \cdot 6
\end{aligned}
$$

$T_{comm}$ is set to 0 as there are no communication on the serial version. Since i used the parallel version for checking serial performance there is a roll back check at the end. It could be removed but it wasn't, it will give low impact on speed. The result is that $T_{node}$ is only dependent on $T_{loop}$

$$
\begin{aligned}
8.73 &= T_{loop} \cdot (160^3 + 160^2 \cdot 6) \\
T_{loop} &= 2.054311 \cdot 10^{-06}
\end{aligned}
$$

The result can be seen above. This is the calculated time for $T_{loop}$. The accurate value will vary depending on where in the matrix $T_{loop}$ is executed, but this result is an averge over all points. If the program ran again, $T_{loop}$ should become the same. In a parallel version there will be more $T_{loop}$ at the borders. Those are smaller as there is only one value that needs updating, but it shoudn't affect the model in a significant way.

Now we will use the results in the serial run, to calculate how much time a parallel run will use. This should be close to the measured time. A 8 node configuration will be used, with the same array as the serial run.

This means that each node will have a $80x80x80$ matrix. The $T_{loop}$ value will be used from the previous run. It might be a bit too high because in this run the rollbackcheck will be done more often than on the serial run.

$$
\begin{aligned}
T_{loop} &= 2.054311 \cdot 10^{-06} \\
T_{rollbackcheck} &= T_{loop} \cdot (n \cdot m \cdot 2 + n \cdot k \cdot 2 + m \cdot k \cdot 2) \\
T_{rollbackcheck} &= 2.054311 \cdot 10^{-06} \cdot (80^2 \cdot 6) \\
T_{rollbackcheck} &= 0.078885542
\end{aligned}
$$

For a 8 node configuratino $T_{rollbackcheck}$ was 0.078885542 seconds.

$$
\begin{aligned}
T_{comm} &= bordersize \cdot 8 \cdot \beta + \alpha \\
T_{comm} &= 80^2 \cdot 6 \cdot 8 \cdot 2.143 \cdot 10^{-10} + 1.995 \cdot 10^{-6} \\
T_{comm} &= 0.001336609
\end{aligned}
$$

The communication time for Njord is very low. $T_{comm}$ is very low, only 0.001336609 seconds.

$$
\begin{aligned}
T_{compute} &= T_{loop} \cdot arraysize \\
T_{compute} &= 2.054311 \cdot 10^{-06} \cdot 80^3 \\
T_{compute} &= 1.0518072
\end{aligned}
$$

$T_{compute}$ is 1.0518072 seconds. The serial runtime was 8.73 seconds which divided by 8 is 1.09125. Close to the calculated performance for each node.

$$
\begin{aligned}
T_{node} &= T_{compute} + T_{comm} + T_{rollbackcheck} \\
T_{node} &= 1.0518072 + 0.001336609 + 0.078885542 \\
T_{node} &= 1.1320294
\end{aligned}
$$

A node will complete in 1.1320294 seconds.

$$
\begin{aligned}
ManhattanDistance &= 4 \\
T_{rollback} &= 0 \\
T_{overall} &= T_{node} \cdot ManhattanDistance + T_{rollback} \\
T_{overall} &= 1.1320294 \cdot 4 \\
T_{overall} &= 4.5281175
\end{aligned}
$$

ManhattanDistance is set to 4, this is because that will be the largest distance from the starting point to the fathers edge in a 2x2x2 node grid. $T_{rollback}$ are ignored, we hope there was no rollbacks.

$T_{overall}$ was calculated to 4.52 seconds the measured time was 4.11 seconds. This means that $T_{loop}$ might be a bit incorrect for this run. It is probably a bit high since there are more border checks in a 8 node division than on the serial run. It is still close to the measured performance.

## 4.4 Validation of HFMM

To calculate the performance of the HFMM algorithm, is a challenge. There are many values that only can be estimated. To make this easier I have taked the 8 node run of the HFMM the same setup as in the validation of PFMM.

$$
\begin{aligned}
T_{loop} &= 2.054311 \cdot 10^{-06} \\
T_{compute} &= T_{loop} \cdot arraysize \\
T_{compute} &= 2.054311 \cdot 10^{-06} \cdot 80^3 \\
T_{compute} &= 1.0518072
\end{aligned}
$$

$T_{compute}$ is 1.0518072 seconds. The same as in PFMM.

$$
\begin{aligned}
T_{comm} &= (bordersize + resendCount) \cdot 20 \cdot \beta + \alpha \\
T_{comm} &= (80^2 \cdot 6 + resendCount) \cdot 20 \cdot 2.143 \cdot 10^{-10} + 1.995 \cdot 10^{-6} \\
T_{comm} &= 0.001645824 + 4.286 \cdot 10^{-9} \cdot resendCount + 1.995 \cdot 10^{-6}
\end{aligned}
$$

$T_{comm}$ is dependent on the resendCount, this value is difficult to predict before the application starts.

$$
\begin{aligned}
T_{traverse} &= arraysize \cdot mtime \\
T_{rollback} &= \#rollback \cdot T_{traverse} \\
T_{rollback} &= \#rollback \cdot 80^3 \cdot mtime
\end{aligned}
$$

It is difficult to say how many rollbacks a run will have. This depends on differences in computational speed between the nodes, network latency and the velocity field.

$$
\begin{aligned}
T_{node} &= T_{wavetime} + T_{compute} + T_{comm} + T_{rollback} \\
T_{node} &= T_{wavetime} + 1.0518072 + 0.001645824 + 4.286 \cdot 10^{-9} \cdot resendCount + 1.995 \cdot 10^{-6} + 7 \\
T_{node} &= 1.053455 + T_{wavetime} + T_{rollback} + 4.286 \cdot 10^{-9} \cdot resendCount \\
T_{node} &= 1.053455 + T_{wavetime} + T_{rollback}
\end{aligned}
$$

This is as close to an answer to $T_{node}$ as one can. The resend count is remove becuase Njord has such good interconnect that it will be negliable. The measured time for this run is 9.9 seconds. Each node uses 1.053455. This means that $T_{wavetime}$ and $T_{rollback}$ accounts for almost 9 seconds. The point start in the middle, that means one of the nodes will have the starting point. It will be exchanged to all nodes within the first few loop runs. So $T_{wavetime}$ is small. Therefore, $T_{rollback}$ accounts for most of the 9 seconds, which is very bad.

# Chapter 5

# Performance analysis

This chapter will describe and explain the different results from the different test runs. All speedup measurements are taken compared to PFMM running on one node. This is because PFMM will give mostly the same results as a serial implementation of algorithm 3.3. The measurements are taken before and after the FMM executes. It doesn't take into account the time it takes to distribute the velocity matrix or initialize all the variables.

All the test have been executed on the njord super computer, unless otherwise specified.

## 5.1   HFMM

HFMM is an implementation of Herrmanns algorithm 2.10. The results from this algorithm didn't give the same timings each time as algorithm 2 did. This will most likely be because of the asynchronous nature of the application. All communication was made asynchronous because the usage of synchronous communication for each loop gave too much overhead, which killed performance, even more than the asynchronous version.

From Figure 5.1, one can see that its not much above 1. This means that the application didn't work faster when it ran on more nodes. This can mostly be credited to the rollback function, as you can see from Figure 5.2 the speedup was drastically better without the rollbacks.

The version without rollbacks show that the rollbacks kill performance, but still there are problems with HFMM. There is also a decrease in performance when we use a non optimal domain decomposition. 9 and 27 nodes, give less performance.

There aren't any speed measurements for HFMM above 64 nodes. This is because it has issues with performance when you go from 16 to 32 and even more when you hit 64 nodes. The best performance for HFMM is when it ran on 8-9 processors.

Figure 5.1: Speedup from HFMM

Figure 5.2: Speedup for HFMM with and without rollback

### 5.1.1   Analasis of HFMM

HFMM showed very poor results compared to what one could expect. Herrmann got his algorithm to run with an efficiency of 0.98 with his optimal decomposition of 4 or 8 nodes. He also used a sphere as a starting figure, instead of a point which was used here. His velocity field was all 1. This should, in theory, give an efficiency of 1 without border communication. HFMM uses a lot of communication, it sends each border when it's calculated. It also uses probe functions to check for incoming messages. Since it uses asynchronous communication it should not be hindered as much as it is by communication. It seems the time to send and check for incoming messages uses a lot more resources than one would think.

With some further work it could be possible to remove some of the communication overhead shown in HFMM. That fact doesn't change the fact that it suffered heavily from rollbacks. When it ran with rollback enabled it got so bad that it didn't even beat the serial application.

To avoid rollbacks it would be an idea to use synchronous communication so each node would not get values that were calculated a few loops ago. Then there would be no rollbacks in Herrmanns example and the efficiency described in Herrmann might be achieved, but the overhead of communicating the border nodes when they are discovered makes it impossible for me to achieve these kinds of speedups. It seems theoretically possible to achieve the efficiency Herrmann got, but it's has proven practically difficult. PFMM would most likely get close to that in performance.

## 5.2   PFMM

PFMM worked much better than HFMM.

In Figure 5.3 it is shown that the application increases its speed as the number of nodes increases. As the 80x80x80 matrix reaches 16 nodes, it has a huge performance increase compared to the larger matrices. The reason for this is that mostly all data can be fitted in cache with such a small matrix. However, as it continues form 32 nodes, it drops drastically. This is because the matrix are now so small on each node. Each node will only have a 20x20x20 matrix. Therefore, the time for communicating the borders will be much larger than calculation times. The same can be observed for 160x160x160 matrix. There is also a good increase in performance for 9 and 27 nodes. This can be accredited to the Manhattan distance to all nodes which is 3 on nine nodes, compared to 4 on eight nodes. The same can be said for 27 nodes where the Manhattan distance to the farthest node is 4, but its 5 on 32 nodes. As the matrix size increases, there is work for more nodes.

In Figure 5.4 it is apparent that large matrices gain a lot from more nodes.

Figure 5.3: Speedup for PFMM

Figure 5.4: Speedup for PFMM, 256 nodes

Figure 5.5: Speedup for PFMM, salt data

In Figure 5.6 one can see the difference between using a velocity field of 1 and using real seismic data. The difference is probably because of small differences in the measurements. They are mostly identical, which means that there are no more iterations when using a seismic dataset compared to a velocity field of 1. If there are more iterations, it's carried out by the nodes in the middle.



Figure 5.6: Speedup for PFMM compared to theoretical speedup

In Figure 5.6 the speedup is compared to theoretical speedup. The theoretical speedup is calculated from Equation 5.1.

$$M/n = speedup \qquad (5.1)$$

Where M is the Manhattan distance and n is the number of nodes used.

This is just to get a perspective on how the theoretical max speedup would be if we remove $T_{rollback}$ and $T_{comm}$.

In the two matrix sizes there is a change between 32 and 64 nodes. When the program uses 64 nodes it runs faster than the theoretical speedup. Even the largest matrix runs past the smaller matrix. The reason they beat the theoretical speedup must be because the problem size fits into level three cache, which is very large on Njord. The reason for the largest matrix becoming faster than

| # nodes | Calculated | Measured |
|---------|------------|----------|
| 2       | 8.7296     | 8.79     |
| 4       | 6.6669     | 6.31     |
| 8       | 4.5232     | 4.11     |
| 16      | 2.8226     | 2.4      |

Table 5.1: Theoretical and measured runtime for PFMM

the smaller one, is because there is less work/transfer time ratio than on the smaller matrix. In the smaller matrix, the transfer time for borders dominates more than on the larger matrix.

Another way to compare the runtime to theoretical, is to use the more advanced model, shown in Section 4.2 PFMM. In Table 5.2 one can see the theoretical time compared to the measured time. This is calculated from Equation 4.13. $T_{comm}$ is set to 0, beacuse its so low that it won't affect the results in a significant way, also we assume no rollbacks.



Figure 5.7: PFMM theoretical vs measured time

From Figure 5.7, the difference in theoretical vs measured time can be seen. The reason for its difference can be that $T_{loop}$ are smaller when the matrix size is smaller. Another reason can be better utilization of the cache, and as a result it goes faster.

## 5.2.1    Analasis of PFMM

PFMM gave much better results than HFMM. The reason for this was that it did much less communication and almost avoids rollbacks. In the case presented by Herrmann where we use 4 computing nodes, and start with a sphere in the middle. Then all nodes would have work and no nodes would require a rollback. This would give almost four in speedup compared to the serial version. Though this situation is not interesting for finding salt formations. There the start condition is a single point.

There are a few ways one can calculate arrival times. The normal 1 order approximation. Higher order approximations and using Spherical coordinates. Spherical coordinates gave more correct results but it's hard to follow the head wave. Using Cartesian coordinates approximations doesn't give the correct solution, buts it's not that far from it either. Using a higher order approximation would make it more correct but also take longer to calculate. The applications use 1 order approximation. Mostly because it gives good enough approximations but also because it is faster than the other methods.

When the number of compute nodes increased, the smaller problem sizes had problems maintaining performance. This was because the problem size on each node became so small, that the cost of communicating borders affected the performance. On problem sizes above 320x320x320 could easily run on 256 nodes and still maintain good performance. 160x160x160 had problems above 32 nodes, while 80x80x80 had the same problem. 80x80x80 also gave very good results on 16 to 32 nodes. This would be because the prolbem size then fitted into cache. 160x160x160 might also have such a case between 32 and 64 nodes. But at 64 nodes it gave less performance than at 32 nodes.

The most important factor for speed on PFMM is the Manhattan distance to the farthest node from the starting node. Since all test result was measured by putting the point in the middle of the problem matrix. Distributions which had odd number in each dimension gave the best results. 9 nodes gave much better results than 8 because of the Manhattan distance of 3 vs 4. The same was with 27 and 32, with Manhattan distance of 4 and 5.

Theoretical speedup was calculated by dividing number of nodes by the Manhattan distance. PFMM did beat the theoretical speedup when it reached 64 nodes. This was surprising because $T_{rollback}$ and $T_{comm}$ was ignored in the theoretical speedup. $T_{rollback}$ is probably 0 in all the cases, but $T_{comm}$ is present but probably negligible as Njord has a very fast interconnect. The most probably reason for getting higher than theoretical speedup, is because the problem size must fit in cache. The performance below 64 nodes is still good.

Choosing a decomposition of the problem size that gave the smallest amount or border size, had negligible effect on the speed. It was not possible to distinguish that from normal variations on test runs, which where very small. Normally less than a few percent. The important factor are that the nodes are divided making the Manhattan distance the smallest.

## 5.3 Comparing HFMM and PFMM



Figure 5.8: 320x320x320 matrix comparing HFMM and PFMM

In Figure 5.8, it is obvious that the PFMM works much better. But in theory the HFMM should work very well. But it kneels under the load of sending incremental border exchanges, and rolling back when receiving a value that requires a rollback. In HFMM the other nodes should start working before they start working in PFMM. But because of the penalty involved in sending and receiving values it doesn't beat the performance of PFMM. When HFMM reaches 8 nodes performance starts to drop, this is because the communication and rollbacks start to take much longer because there are much more communication on 16 and 32 nodes than on 8. PFMM scales better, and is not affected by increasing the amount of nodes.

# Chapter 6

# Conclusion

This thesis has focused on developing a parallel method for the Fast Marhcing Method, used in finding salt formations in seismic data. A algorithm from Herrmann was looked at and used as a reference point. Two different approaches of making a parallel Fast Marching Method was tried and tested.

HFMM used the algorithm described be Herrmann. It didn't perform as fast as expected. This was because of too many rollbacks. Because of the asynchroneous communication and different execution times at each node, rollbacks became plentiful. Rollbacks stand for 80-90% of the execution time for 8 node configuration. At first synchroneous communication was tested, this resulted in too much overhead in communication and performed worse than the asynchroneous version.

PFMM had much better performance. The teoretical speedup of this algorithm number of computational nodes divided by the manhattan distance from the starting point to the farthest node. The application ran almost as fast as the theroretical speedup, but with more than 32 nodes it gave faster than theoretical speedup. Because of the one pass nature of the Fast Marching Method it is impossible to get full cpu utilization of all nodes. This can be improved by running more points at the same time, giving only a small percentage penalty, more info in Future work.

An execution time model was developed, this model performed well for PFMM algorithm. It came very close to the actual execution time. If more work had been laid into estimating $T_{loop}$ it would have been even better. The model for HFMM was much more difficult, asynchroneous communication made it very difficult to predict how many rollbacks would be used. The time used in computation for HFMM is very low compared to what is used for rollbacks and communication.

## 6.1 Future work

There are a few optimizations that would be interesting to try, given that the cpu utilization is so low on each node. One optimization would be running more picks at the same time. Also changing the decomposition so the manhattan distance becomes smaller would increase speed.

The speedup from the application was good compared to what the algorithm could theoretical achieve. However it's not a very good utilization of the computing nodes. The start node would be idle while the other nodes calculate their nodes. There are a way to avoid some of this problem, by running more FMM simultaneous. When a geologist tries to locate Salt they pick two to three points at the same time. If two picks could execute at the same time, it would only give a slight increase in computing time when a node had to process bout waves at the same time. This should only occur if the two picks had a node with the same Manhattan distance from each pick. In the case of a 5x5 node grid. There are two collisions, at distance 2 and 3. Figure 6.1 shows the distance for pick 1, while Figure 6.2 shows the distance for pick 2. In Figure 6.3 the colliding nodes are in bold. Since the max distance is 6 for each pick and bouth distance 2 and 3 has colliding nodes. The distance will become 8. Thats a 33% increase in execution time for running two picks at the same time. Which must be a very good increase in performance. It might be even better for 3 picks but more nodes would probably collide. This was not tested and are left as an optimization for later.

| 2 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | 1 | 2 | 3 |
| 2 | 1 | 2 | 3 | 4 |
| 3 | 2 | 3 | 4 | 5 |
| 4 | 3 | 4 | 5 | 6 |

Figure 6.1: Manhattan distance for pick 1

| 6 | 5 | 4 | 3 | 4 |
|---|---|---|---|---|
| 5 | 4 | 3 | 2 | 3 |
| 4 | 3 | 2 | 1 | 2 |
| 3 | 2 | 1 | 0 | 1 |
| 4 | 3 | 2 | 1 | 2 |

Figure 6.2: Manhattan distance for pick 2

In this thesis only rectangular domain decompositions are tested. It might be an idea to test a domain decomposition that use beams from the pick point. In

| 2 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | 1 | **2** | 3 |
| 4 | 1 | **2** | 1 | 2 |
| **3** | **2** | 1 | 0 | 1 |
| 4 | **3** | 2 | 1 | 2 |

Figure 6.3: Manhattan distance for bouth picks, bold where they collide

a 2D case each node could take an equal amount of degrees from the starting point, a cake piece. This will probably give more equal work among the different nodes. The correct result would probably be acheved in 2-3 iterations, beacuse the Manhattan distance would then be 1. The extra iterations are added if the wave moves in and out of different nodes, which is very likely, since the velocity field is not a single value.

Implementing such a division is not straight forward. It will be much more difficult to find out which point in the matrix are on which node. This is because a line at 32 degrees will split many points in the grid. Which one should belong to which node, and how do we add border values. This is solvable but will make the application much more complex. It can also be extended to 3D space by adding another dimension.

# Bibliography

[1] Tariq Alkhalifah and Sergey Fomel. Implementing the fast marching eikonal solver: Spherical versus cartesian coordinates. Available from World Wide Web: `citeseer.ist.psu.edu/543803.html`.

[2] N. Sukumar, D. L. Chopp, and B. Moran. Extended finite element method and fast marching method for three-dimensional fatigue crack propagation. Available from World Wide Web: `citeseer.ist.psu.edu/654120.html`.

[3] Yan J Zhuang TG Zhao B Schwartz LH. Lymph node segmentation from ct images using fast marching method. Available from World Wide Web: `http://www.ncbi.nlm.nih.gov/sites/entrez?cmd=Retrieve\&db=PubMed\&list_uids=15127747\&dopt=Abstract`.

[4] Alexandru Telea and Jarke J. van Wijk. An augmented fast marching method for computing skeletons and centerlines. In *VISSYM '02: Proceedings of the symposium on Data Visualisation 2002*, pages 251–ff, Aire-la-Ville, Switzerland, Switzerland, 2002. Eurographics Association.

[5] Scott A. Morgan Morgan P.Brown and Amerada Hess Corporation Gary Whittle. Seismic event tracking by global path optimization. *SEG/New Orleans 2006 Annual Meeting*, 2006. Article.

[6] Shewchuk. Triangle: Engineering a 2D quality mesh generator and delaunay triangulator. In *WACG: 1st Workshop on Applied Computational Geometry: Towards Geometric Engineering, WACG*. LNCS, 1996. Available from World Wide Web: `citeseer.ist.psu.edu/126089.html`.

[7] MPI Forum. Message passing interface, specification. Available from World Wide Web: `http://www.mpi-forum.org/docs/`.

[8] Notur. Njord. Available from World Wide Web: `http://www.notur.no/hardware/njord/`.

[9] IBM. Power5 system microarchitecture. Available from World Wide Web: `http://researchweb.watson.ibm.com/journal/rd/494/sinharoy.html`.

[10] Eugene F. Krause. Taxicab geometry, 1973.

[11] J. Sethian. A fast marching level set method for monotonically advancing fronts. In *Proc. Nat. Acad. Sci.*, volume 93, pages 1591–1595, 1996. Available from World Wide Web: `citeseer.ist.psu.edu/sethian95fast.html`.

[12] D. Adalsteinsson and J. Sethian. The fast construction of extension velocities in level set methods. *Journal of Computational Physics*, 148:2–22, 1998. Available from World Wide Web: `citeseer.ist.psu.edu/adalsteinsson97fast.html`.

[13] Elisabeth Rouy and Agn&#232;s Tourin. A viscosity solutions approach to shape-from-shading. *SIAM J. Numer. Anal.*, 29(3):867–884, 1992.

[14] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms, Second Edition*. The MIT Press, September 2001. Available from World Wide Web: `http://www.amazon.ca/exec/obidos/redirect?tag=citeulike04-20{\&}path=ASIN/0262531968`.

[15] M. Herrmann. A domain decomposition parallelization of the fast marching method, 2003.

# Appendix A

# Application 1 source code

## A.1   Array

### A.1.1   array_mpi.h

```
1  #ifndef ARRAY_H
2  #define ARRAY_H
3
4  #include "mpi.h"
5  /*
6   * Defines the size of the array
7   */
8  #define SIZEX 160 // 448//8//448
9  #define SIZEY 160 // 704//4//704
10 #define SIZEZ 160 // 1216//4//1216
11 /*
12  * Calculate the strides for global array
13  */
14 #define stride0 ((SIZEZ+2)*(SIZEY+2))
15 #define stride1 (SIZEZ+2)
16 #define stride2 1
17 /*
18  * Gets the index for a position in the global array
19  */
20 #define GETINDEX(i,j,k) stride0*(i+1) + stride1*(j+1) + stride2*(k
      +1)
21 /*
22  * Get the index for a position in the local array
23  */
24 #define GETLINDEX(i,j,k) ((local_z+2)*(local_y+2)) *(i+1) + ((
      local_z+2)*(j+1)) + k+1
25 /*
26  * A large float number, should be larger than anything you
       calculate
27  */
28 #define BIGFLOAT 10000000.0
29 /*
30  * The size of each array, for mallocing memory
31  */
```

```
32  #define  ARRAYSIZE  (SIZEX+2)*(SIZEY+2)*(SIZEZ+2)
33  #define  LOCALARRAYSIZE  (local_x+2)*(local_y+2)*(local_z+2)
34  /*
35   * Rank, MPI rank
36   * cartrank, rank in the cartesian grid
37   * size, number of nodes used
38   */
39  int rank,cartrank, size;
40  /*
41   * Ranks of nodes that are above,below,west,east,north,south
42   */
43  int above, below, west, east, north ,south;
44  /*
45   * Number of nodes in each dimension
46   */
47  int dims[3];
48
49  /*
50   * communicator for the cartesian grid
51   */
52  MPI_Comm gridcomm;
53  /*
54   * MY coordinates in the cartesian node grid
55   */
56  int coords[3];
57  /*
58   * size in each dimension of the global array
59   */
60  int x,y,z;
61  /*
62   * Size in each dimension of the local array
63   */
64  int local_x,local_y,local_z;
65  /*
66   * Print the values of a global array to stdout
67   */
68  void printArray(float* array);
69  /*
70   * Print the local array to stdout
71   */
72  void printLocalArray(float* array);
73  /*
74   * print a local int array to stdout
75   */
76  void printLocalIntArray(int* array);
77  /*
78   * Get what node a global position resides in
79   */
80  int getDest(int x, int y, int z);
81  /*
82   * Get global coordinates from local coordinates
83   */
84  int* getGlobalCord(int x, int y, int z);
85  /*
86   * Get local coordinates from global coordinates
87   */
88  int* getLocalCord(int x, int y, int z);
89  #endif
```

## A.1.2   array_mpi.c

```c
1  #include "array_mpi.h"
2  #include <stdlib.h>
3  #include <stdio.h>
4
5  /*
6   * converts global coord to local coord
7   */
8  int* getLocalCord(int xg, int yg, int zg){
9          int *cr;
10         cr = malloc(sizeof(int)*3);
11         cr[0] = xg;
12         cr[1] = yg;
13         cr[2] = zg;
14   cr[0] -= coords[0] * local_x;
15   cr[1] -= coords[1] * local_y;
16   cr[2] -= coords[2] * local_z;
17         return cr;
18 }
19 /*
20  * convert local coord to global coord
21  */
22 int* getGlobalCord(int x, int y, int z){
23         int *cr;
24         cr = malloc(sizeof(int)*3);
25         cr[0] = x+ local_x*coords[0];
26         cr[1] = y+ local_y*coords[1];
27         cr[2] = z+ local_z*coords[2];
28         return cr;
29         }
30 /*
31  * gets the rank of the node that have input position
32  */
33 int getDest(int x, int y, int z){
34         int cr[3];
35         int value;
36         cr[0] = x/local_x;
37         cr[1] = y/local_y;
38         cr[2] = z/local_z;
39         MPI_Cart_rank(gridcomm,cr,&value);
40         return value;
41 }
42 /*
43  * print the local float array to stdout
44  */
45 void printLocalArray(float* array){
46         int i,j,k;
47         for(k=-1;k<=local_z;k++){
48         printf("%d: array_z_=_%d\n",cartrank,k);
49                 for(i=-1;i<=local_x;i++){
50                         printf("%d:_",cartrank);
51                         for(j=-1;j<=local_y;j++){
52                                 printf("_%8f_",array[GETLINDEX(i,j,k
                                        )]);
53                         }
54                         printf("\n");
55                 }
```

```
56              printf("\n\n");
57              }
58  }
59  /*
60   * print a local integer array to stdout
61   */
62  void printLocalIntArray(int* array){
63          int i,j,k;
64          for(k=-1;k<=local_z;k++){
65          printf("%d:array_z_=_%d\n",cartrank,k);
66                  for(i=-1;i<=local_x;i++){
67                          printf("%d:_",cartrank);
68                          for(j=-1;j<=local_y;j++){
69                                  printf("_%d_",array[GETLINDEX(i,j,k)
                                          ]);
70                          }
71                          printf("\n");
72                  }
73          printf("\n\n");
74          }
75  }
76  /*
77   * print the global array to stdout
78   */
79  void printArray(float* array){
80          int i,j,k;
81
82          for(k=-1;k<=z;k++){
83          printf("array_z_=_%d\n",k);
84                  for(i=-1;i<=x;i++){
85                          for(j=-1;j<=y;j++){
86                                  printf("_%8f_",array[GETINDEX(i,j,k)
                                          ]);
87                          }
88                  printf("\n");
89                  }
90          printf("\n\n");
91          }
92  }
```

# A.2   Fast Marching Method

## A.2.1   fmm_mpi.h

```
1  #include "heap.h"
2
3  /*
4   * a struct to store all variables for a given fmm implementation
5   */
6  typedef
7  struct f {
8    Heap* heap; // the heap
9    float* timearray; // the timearray storing arrivaltimes
10   int* bandarray; // storing band information
11   float* velocityarray; // velocity field
12   float* sentarray; // array containting values of sent points, used
          to avoid sending the same value multiple times
13   int x,y,z,posx,posy,posz; // x,y,z is size of local fmm matrix,
         posx, posy, posz is global coords for the starting point
14   } FmmData;
15
16 typedef
17 /*
18  * Struct used for sending a border cell to another node. String its
         value and posistion with band information
19  */
20 struct me {
21         float value;
22         int x;
23         int y;
24         int z;
25   int n;
26 } MPI_Element;
27 /*
28  * initilalize the FMM set velocity arrat, position of starting
        point and size of array
29  */
30 FmmData* initFMM(float* velarray,int posx,int posy, int posz, int x,
        int y, int z);
31 /*
32  * Free up used variables in FMM
33  */
34 void freeFMM(FmmData* data);
35 /*
36  * Execute the FMM
37  */
38 void executeFMM(FmmData* data);
```

## A.2.2   fmm_mpi.c

```c
1   #include <stdlib.h>
2   #include <string.h>
3   #include <stdio.h>
4   #include <math.h>
5   #include "fmm_mpi.h"
6   #include "heap.h"
7   #include "array_mpi.h"
8
9
10  #define BAND 0
11  #define OUTSIDE −1
12  #define KNOWN n
13
14  //#define DEBUG
15  /*
16   * number of loops
17   */
18  int n;
19
20  /*
21   * largest_solution is the largest value in the array
22   * rollbacksmallest is the smallest value of received border values,
            which dictates the rollback number
23   */
24  float largest_solution , rollbacksmallest;
25
26  /*
27   * the n value which one should rollback to
28   */
29  int rollbackn;
30  /*
31   * mpi datatype for sending border points
32   */
33  MPI_Datatype mpi_element_struct;
34  /*
35   * Used for debug output
36   */
37  float valuemax=0;
38  float valuemin=0;
39  /*
40   * array containing the working status of each node
41   */
42  int* working;
43
44  /*
45   * Add a point to the heap
46   */
47  void addToHeap(FmmData* data , int px , int py , int pz){
48    Element *temp;
49    #ifdef DEBUG
50    printf("%d:_adding_%d,_%d,_%d_to_heap,heap_size_is_%d_maxsize_if_%
          d\n" , cartrank , px , py , pz , data−>heap−>heapsize , data−>heap−>maxsize
          );
51    #endif
52    temp= malloc(sizeof(Element));
53    temp−>value = data−>timearray[GETLINDEX(px , py , pz)];
```

```
54    temp->x = px;
55    temp->y = py;
56    temp->z = pz;
57    heapInsert(data->heap,temp);
58    #ifdef DEBUG
59    printf("%d:_inserted_%d,_%d,_%d_to_heap,heap_size_is_%d\n",
          cartrank,px,py,pz,data->heap->heapsize);
60    #endif
61
62  }
63
64  /*
65   * initialize the FMM
66   * velarray is the velocity field
67   * posx,y,z is position of the starting point
68   * x,y,z is the size of the array, most are read from array_mpi.h
69   */
70  FmmData* initFMM(float* velarray, int posx,int posy, int posz, int x
        , int y, int z){
71    working = malloc(sizeof(int) * size);
72    /* init datatypes */
73    MPI_Element e;
74    MPI_Datatype type[5] = { MPI_FLOAT, MPI_INT, MPI_INT, MPI_INT,
          MPI_INT };
75    int blocklen[5] = { 1, 1, 1, 1 , 1};
76    MPI_Aint disp[5];
77    disp[0] = 0;
78    disp[1] = sizeof(float);
79    disp[2] = sizeof(int) + disp[1];
80    disp[3] = sizeof(int) + disp[2];
81    disp[4] = sizeof(int) + disp[3];
82    MPI_Type_create_struct(5, blocklen, disp, type, &
          mpi_element_struct);
83    MPI_Type_commit(&mpi_element_struct);
84    /* end init datatypes */
85    #ifdef DEBUG
86    printf("%d:_done_init_mpi_datatypes\n",cartrank);
87    #endif
88    int send = 0;
89    FmmData* data;
90    data = malloc(sizeof(FmmData));
91    if(data == 0){
92      printf("%d:_Failed_to_allocate_memory,_exiting\n",cartrank);
93      exit(1);
94    }
95
96    data->bandarray = malloc(sizeof(int)*LOCALARRAYSIZE);
97    if(data->bandarray == 0){
98      printf("%d:_Failed_to_allocate_memory,_exiting\n",cartrank);
99      exit(1);
100   }
101   data->timearray = malloc(sizeof(float)*LOCALARRAYSIZE);
102   if(data->timearray == 0){
103     printf("%d:_Failed_to_allocate_memory,_exiting\n",cartrank);
104     exit(1);
105   }
106
107   data->velocityarray = velarray;
```

```
108    #ifdef DEBUG
109    printf("%d:_initializing_heap\n",cartrank);
110    #endif
111    data->heap = initHeap(LOCALARRAYSIZE);
112    data->x = x;
113    data->y = y;
114    data->z = z;
115    data->posx = posx;
116    data->posy = posy;
117    data->posz = posz;
118
119    #ifdef DEBUG
120    printf("%d:_clearing_memory_bandarray_%p_size_%d\n",cartrank,data
           ->bandarray,sizeof(int)*LOCALARRAYSIZE);
121    fflush(stdout);
122    #endif
123    /*
124     * setting the band array to outside
125     */
126    memset(data->bandarray,OUTSIDE,sizeof(int)*LOCALARRAYSIZE);
127    #ifdef DEBUG
128    printf("%d:_Heap_ok,_clearing_memory_timearray_%p_size_%d\n",
           cartrank,data->timearray,sizeof(float)*LOCALARRAYSIZE);
129    #endif
130    /*
131     * zeroing the arrival time array
132     */
133    int i = 0;
134    for(i=0;i<LOCALARRAYSIZE;i++)
135      data->timearray[i] = 0;
136    //memset(data->timearray,0,sizeof(float)*LOCALARRAYSIZE);
137    //bzero(data->timearray,sizeof(float)*LOCALARRAYSIZE);
138    data->sentarray = malloc(sizeof(float)*LOCALARRAYSIZE);
139    bzero(data->sentarray,sizeof(float)*LOCALARRAYSIZE);
140    #ifdef DEBUG
141    printf("%d:_done_allocating_memory,_setting_starting_point\n",
           cartrank);
142    fflush(stdout);
143    #endif
144    /*
145     * inserting the starting point on the correct node and add it to
            the heap / narrow band
146     */
147    if(getDest(posx,posy,posz) == cartrank){
148      int *cr = getLocalCord(posx,posy,posz);
149      // insert starting point
150      data->timearray[GETLINDEX(cr[0],cr[1],cr[2])] = 0;
151      data->bandarray[GETLINDEX(cr[0],cr[1],cr[2])] = BAND;
152      //printf("%d: sat pos %d %d %d, as known\n",cartrank,cr[0],cr
            [1],cr[2]);
153      Element* element;
154      element = malloc(sizeof(Element));
155      element->x = cr[0];
156      element->y = cr[1];
157      element->z = cr[2];
158      element->value = data->timearray[GETLINDEX(cr[0],cr[1],cr[2])];
159      heapInsert(data->heap,element);
160      send = 1;
```

```
161      }
162      return data;
163    }
164    /*
165     * return min
166     */
167    float min(float per, float truls){
168      if(per > truls)
169        return truls;
170      else return per;
171    }
172    /*
173     * return max
174     */
175    float max(float per, float truls){
176      if(per > truls)
177        return per;
178      else return truls;
179    }
180
181    /*
182     * prints xy plane from a float array
183     */
184    void printFloatArray(int sizex, int sizey, int z, float* array){
185      int i,j;
186      printf("\n_Printing_matrix\n");
187      for(i=0;i<sizex;i++){
188        for(j=0;j<sizey;j++){
189          printf("_%8.2f_",array[GETLINDEX(i,j,z)]);
190        }
191        printf("\n");
192      }
193    }
194
195    /*
196     * calculate the arrival time for a point x,y,z
197     */
198    float calcDistance(FmmData* data, int x, int y, int z){
199      float sol;
200      sol = BIGFLOAT;
201      if(data->bandarray[GETLINDEX(x+1,y,z)] > BAND){
202        sol = min(data->timearray[GETLINDEX(x+1,y,z)] +1/data->
              velocityarray[GETLINDEX(x,y,z)],sol);
203        #ifdef DEBUG
204        printf("%d:_sol_is_%f_for_x+1\n",cartrank,data->timearray[
              GETLINDEX(x+1,y,z)] +1/data->velocityarray[GETLINDEX(x,y,z)])
              ;
205        #endif
206      }
207      if(data->bandarray[GETLINDEX(x-1,y,z)] > BAND){
208        sol = min(data->timearray[GETLINDEX(x-1,y,z)] +1/data->
              velocityarray[GETLINDEX(x,y,z)],sol);
209        #ifdef DEBUG
210        printf("%d:_sol_is_%f_for_x-1\n",cartrank,data->timearray[
              GETLINDEX(x-1,y,z)] +1/data->velocityarray[GETLINDEX(x,y,z)])
              ;
211        #endif
212      }
```

```
213    if(data->bandarray[GETLINDEX(x,y+1,z)] > BAND){
214      sol = min(data->timearray[GETLINDEX(x,y+1,z)] +1/data->
             velocityarray[GETLINDEX(x,y,z)],sol);
215      #ifdef DEBUG
216      printf("%d:_sol_is_%f_for_y+1\n",cartrank,data->timearray[
             GETLINDEX(x,y+1,z)] +1/data->velocityarray[GETLINDEX(x,y,z)])
             ;
217      #endif
218    }
219    if(data->bandarray[GETLINDEX(x,y-1,z)] > BAND){
220      sol = min(data->timearray[GETLINDEX(x,y-1,z)] +1/data->
             velocityarray[GETLINDEX(x,y,z)],sol);
221      #ifdef DEBUG
222      printf("%d:_sol_is_%f_for_y-1\n",cartrank,data->timearray[
             GETLINDEX(x,y-1,z)] +1/data->velocityarray[GETLINDEX(x,y,z)])
             ;
223      #endif
224    }
225    if(data->bandarray[GETLINDEX(x,y,z+1)] > BAND){
226      sol = min(data->timearray[GETLINDEX(x,y,z+1)] +1/data->
             velocityarray[GETLINDEX(x,y,z)],sol);
227      #ifdef DEBUG
228      printf("%d:_sol_is_%f_for_z+1\n",cartrank,data->timearray[
             GETLINDEX(x,y,z+1)] +1/data->velocityarray[GETLINDEX(x,y,z)])
             ;
229      #endif
230    }
231    if(data->bandarray[GETLINDEX(x,y,z-1)] > BAND){
232      sol = min(data->timearray[GETLINDEX(x,y,z-1)] +1/data->
             velocityarray[GETLINDEX(x,y,z)],sol);
233      #ifdef DEBUG
234      printf("%d:_sol_is_%f_for_z-1\n",cartrank,data->timearray[
             GETLINDEX(x,y,z-1)] +1/data->velocityarray[GETLINDEX(x,y,z)])
             ;
235      #endif
236    }
237
238    return sol;
239  }
240
241  /*
242   * calculate a new point the the arrival time array if its inside
          the array and not KNOWN, if its outside add it to the heap
243   */
244  void calcElement(FmmData* data, int px, int py, int pz){
245          int add = 0;
246          float sol;
247          Element *temp;
248          /*
249           * Check if the point is inside the array */
250          if(px >= 0 && px < data->x && py >= 0 && py < data->y && pz
                >= 0 && pz < data->z){
251          /* Make sure the point is not KNOWN */
252          if(data->bandarray[GETLINDEX(px,py,pz)] <= BAND){
253                  sol = calcDistance(data,px,py,pz);
254                  /* Check if the number is not smaller than the one
                         we calculated, when we exit, should never happend
                          in serial version */
```

```
255                    if(data−>timearray[GETLINDEX(px,py,pz)] != 0 && data
                           −>timearray[GETLINDEX(px,py,pz)] <= sol){
256                            return;
257                    }
258                    /* If the point is OUSIDE add it to the heap/
                           narrowband */
259                    if(data−>bandarray[GETLINDEX(px,py,pz)] == OUTSIDE)
                           {
260                            data−>bandarray[GETLINDEX(px,py,pz)] = BAND;
261                            add = 1;
262                    }
263
264                    #ifdef DEBUG
265                    printf("%d:_sol_is_%f,_for_%d_%d_%d\n",cartrank,sol,
                           px,py,pz);
266                    #endif
267                    /* store max and min for debug purposes */
268                    #ifdef DEBUG
269                    if(sol > valuemax && sol != BIGFLOAT){
270                            valuemax=sol;
271                    }
272                    if(sol < valuemin){
273                            valuemin = sol;
274                    }
275                    #endif
276                    /* set the new arrivaltime */
277                    data−>timearray[GETLINDEX(px,py,pz)] = sol;
278                    #ifdef DEBUG
279                    printf("%d:_x_%d,_y_%d,_z_%d,_k[i]_%d_l[i]_%d_m[i]_
                           %d,_value_%f_\n",cartrank,data−>x, data−>y, data
                           −>z,px,py,pz, data−>timearray[GETLINDEX(px,py,pz)
                           ]);
280                    #endif
281                    /* add it to the narrowband if it should be added */
282                    if(add){
283                            addToHeap(data,px,py,pz);
284                            }
285                    }
286            }
287 }
288
289 //remove ???
290 void checkforchange(FmmData* data,int px, int py, int pz){
291    float sol = BIGFLOAT;
292    int i;
293    int o[6] = {px−1,px,px+1,px,px,px};
294    int l[6] = {py,py−1, py,py+1,py,py};
295    int m[6] = {pz,pz,pz,pz,pz−1,pz+1};
296    for(i=0;i<6;i++){
297      if(o[i] >= 0 && o[i] < data−>x && l[i] >= 0 && l[i] < data−>y
            && m[i] >= 0 && m[i] < data−>z && data−>bandarray[GETLINDEX(o
            [i],l[i],m[i])] > BAND){
298
299         #ifdef DEBUG
300         printf("%d:_checking_for_rollback_%d_%d_%d_\n",cartrank,o[i],l
               [i],m[i]);
301         #endif
302         sol = calcDistance(data,px,py,pz);
```

```
303
304          if(sol < data−>timearray[GETLINDEX(o[i],l[i],m[i])] && data−>
                 bandarray[GETLINDEX(o[i],l[i],m[i])] > BAND){
305            // new value is smaller lets add this to our heap
306
307            #ifdef DEBUG
308            printf("%d:_rolling_back_%d_%d_%d_old_value_%f_new_value_%f\
                  n",cartrank,o[i],l[i],m[i],data−>timearray[GETLINDEX(o[i
                  ],l[i],m[i])],sol);
309            #endif
310            data−>timearray[GETLINDEX(o[i],l[i],m[i])] = sol;
311            data−>bandarray[GETLINDEX(o[i],l[i],m[i])] = BAND;
312            addToHeap(data,o[i],l[i],m[i]);
313            checkforchange(data,o[i],l[i],m[i]);
314          }
315        }
316      }
317  }
318
319  /*
320   * Add a new element to the array
321   */
322  void addElement(FmmData* data,MPI_Element e){
323    int* cr;
324    int i,j,k;
325    cr = getLocalCord(e.x,e.y,e.z);
326
327    #ifdef DEBUG
328    printf("%d:_adding_element_%d_%d_%d,_to_local_%d_%d_%d\n",cartrank
            ,e.x,e.y,e.z,cr[0],cr[1],cr[2]);
329    #endif
330    /*
331     * checking if we have added it before
332     * Should be always no since we don't send the same value multiple
             times
333     */
334    if(data−>timearray[GETLINDEX(cr[0],cr[1],cr[2])] == e.value){
335      #ifdef DEBUG
336      printf("%d:_alreaddy_added_%d_%d_%d\n",cartrank,e.x,e.y,e.z);
337      #endif
338      return;
339    }
340    data−>timearray[GETLINDEX(cr[0],cr[1],cr[2])] = e.value;
341    data−>bandarray[GETLINDEX(cr[0],cr[1],cr[2])] = e.n;
342    int o[6] = {cr[0]−1,cr[0],cr[0]+1,cr[0],cr[0],cr[0]};
343    int l[6] = {cr[1],cr[1]−1, cr[1],cr[1]+1,cr[1],cr[1]};
344    int m[6] = {cr[2],cr[2],cr[2],cr[2],cr[2]−1,cr[2]+1};
345    /*
346     * recalculate all neighbours
347     */
348    for(i=0;i<6;i++){
349
350      #ifdef DEBUG
351      printf("%d:_calc_element_%d_%d_%d\n",cartrank,o[i],l[i],m[i]);
352      #endif
353      calcElement(data, o[i], l[i], m[i]);
354    }
355    /*
```

```
356      * set rollback values, so we can check if a rollback is necesarry
357      */
358      /*if(largest_solution > e.value && rollbacksmallest > e.value){
359      rollbacksmallest = e.value;
360      rollbackn = e.n;
361      }*/
362  }
363
364  /*
365   * rollback all values above input value
366   */
367  void rollback(FmmData* data, float value){
368      int i,j,k;
369      for(i=0;i<local_x;i++)
370          for(j=0;j<local_y;j++)
371              for(k=0;k<local_z;k++){
372                  if(data->bandarray[GETLINDEX(i,j,k)] > BAND && data->
                         timearray[GETLINDEX(i,j,k)] > value){
373
374                      data->bandarray[GETLINDEX(i,j,k)] = BAND;
375                      addToHeap(data,i,j,k);
376                  }
377              }
378  }
379
380  /*
381   * send new values and check for incoming border values
382   */
383  void sendRecvBorderChanges(FmmData* data,int x, int y, int z,int
         send){
384      MPI_Element e;
385      int run= 1;
386      int *cr;
387      int reast,rwest,rnorth,rsouth,rabove,rbelow;
388      int seast,swest,snorth,ssouth,sabove,sbelow;
389      reast = rwest = rnorth = rsouth = rabove = rbelow = 0;
390      seast = swest = snorth = ssouth = sabove= sbelow = 0;
391      /*
392       * if we are to send a value
393       */
394      if(send){
395          /*
396           * see where we have to send the value
397           */
398          if(x == 0){
399              snorth = 1;
400          }
401          if(x == local_x -1){
402              ssouth = 1;
403          }
404          if(y == 0){
405              swest = 1;
406          }
407          if(y == local_y -1){
408              seast = 1;
409          }
410          if(z == 0){
411              sbelow = 1;
```

```
412          }
413        if(z == local_z −1){
414          sabove = 1;
415        }
416        #ifdef DEBUG
417        printf("%d:_%d_%d_%d_sending_to_north_%d_%d_south_%d__%d_west_%d
                _%d_east_%d_%d_below_%d_%d_above_%d_%d\n",cartrank,x,y,z,
                snorth,north,ssouth,south,swest,west,seast,east,sbelow,below,
                sabove,above);
418        #endif
419
420
421     if(snorth || ssouth || swest || seast || sabove || sbelow){
422        cr = getGlobalCord(x,y,z);
423        #ifdef DEBUG
424        printf("%d:_sending_element_at_%d_%d_%d_gave_global_coord_%d_%d_
                %d\n",cartrank,x,y,z,cr[0],cr[1],cr[2]);
425        #endif
426        e.x = cr[0];
427        e.y = cr[1];
428        e.z = cr[2];
429        e.value = data−>timearray[GETLINDEX(x,y,z)];
430        data−>sentarray[GETLINDEX(x,y,z)] = e.value;
431        e.n = data−>bandarray[GETLINDEX(x,y,z)];
432
433        if(snorth){
434          MPI_Send(&e,1,mpi_element_struct,north,1,gridcomm);
435        }
436        if(ssouth){
437          MPI_Send(&e,1,mpi_element_struct,south,1,gridcomm);
438        }
439        if(swest){
440          MPI_Send(&e,1,mpi_element_struct,west,1,gridcomm);
441        }
442        if(seast){
443          MPI_Send(&e,1,mpi_element_struct,east,1,gridcomm);
444        }
445        if(sabove){
446          MPI_Send(&e,1,mpi_element_struct,above,1,gridcomm);
447        }
448        if(sbelow){
449          MPI_Send(&e,1,mpi_element_struct,below,1,gridcomm);
450        }
451     }
452     }
453     rollbacksmallest = BIGFLOAT;
454     rollbackn = 0;
455     /*
456      * a loop to receive all incoming border values
457      */
458     while(run){
459     MPI_Iprobe(MPI_ANY_SOURCE,1,gridcomm,&run,MPI_STATUS_IGNORE);
460     if(run){
461       MPI_Iprobe(north,1,gridcomm,&rnorth,MPI_STATUS_IGNORE);
462       MPI_Iprobe(south,1,gridcomm,&rsouth,MPI_STATUS_IGNORE);
463       MPI_Iprobe(east,1,gridcomm,&reast,MPI_STATUS_IGNORE);
464       MPI_Iprobe(west,1,gridcomm,&rwest,MPI_STATUS_IGNORE);
465       MPI_Iprobe(above,1,gridcomm,&rabove,MPI_STATUS_IGNORE);
```

```
466        MPI_Iprobe(below,1,gridcomm,&rbelow,MPI_STATUS_IGNORE);
467        if(rnorth || rsouth || rwest || reast || rabove || rbelow){
468
469          if(rnorth){
470            MPI_Recv(&e,1,mpi_element_struct,north,1,gridcomm,
                    MPI_STATUS_IGNORE);
471            addElement(data,e);
472          }
473          if(rsouth){
474            MPI_Recv(&e,1,mpi_element_struct,south,1,gridcomm,
                    MPI_STATUS_IGNORE);
475            addElement(data,e);
476          }
477          if(rwest){
478            MPI_Recv(&e,1,mpi_element_struct,west,1,gridcomm,
                    MPI_STATUS_IGNORE);
479            addElement(data,e);
480          }
481          if(reast){
482            MPI_Recv(&e,1,mpi_element_struct,east,1,gridcomm,
                    MPI_STATUS_IGNORE);
483            addElement(data,e);
484          }
485          if(rabove){
486            MPI_Recv(&e,1,mpi_element_struct,above,1,gridcomm,
                    MPI_STATUS_IGNORE);
487            addElement(data,e);
488          }
489          if(rbelow){
490            MPI_Recv(&e,1,mpi_element_struct,below,1,gridcomm,
                    MPI_STATUS_IGNORE);
491            addElement(data,e);
492          }
493        }
494      }
495      }
496      /*
497       * rollback if necessary
498       */
499      if(rollbackn){
500        rollback(data,rollbacksmallest);
501        n = rollbackn;
502        rollbackn = 0;
503        rollbacksmallest = BIGFLOAT;
504      }
505 }
506
507 /*
508  * old synchroneous border exhcange
509  */
510 /*
511 void sendRecvBorderChanges(FmmData* data,int x, int y, int z,int
        send){
512    MPI_Element e;
513
514    int *cr;
515    int reast,rwest,rnorth,rsouth,rabove,rbelow;
516    int seast,swest,snorth,ssouth,sabove,sbelow;
```

```
517     reast = rwest = rnorth = rsouth = rabove = rbelow = 0;
518     seast = swest = snorth = ssouth = sabove= sbelow = 0;
519     if(send){
520       if(x == 0){
521         snorth = 1;
522       }
523       if(x == local_x −1){
524         ssouth = 1;
525       }
526       if(y == 0){
527         swest = 1;
528         }
529       if(y == local_y −1){
530         seast = 1;
531       }
532       if(z == 0){
533         sbelow = 1;
534       }
535       if(z == local_z −1){
536         sabove = 1;
537       }
538     #ifdef DEBUG
539     printf("%d: %d %d %d sending to north %d %d south %d  %d west %d
              %d east %d %d below %d %d above %d %d\n",cartrank,x,y,z,
              snorth,north,ssouth,south,swest,west,seast,east,sbelow,below,
              sabove,above);
540     #endif
541     }

543     MPI_Send(&swest,1,MPI_INT,west,0,gridcomm);
544     MPI_Recv(&reast,1,MPI_INT,east,0,gridcomm,MPI_STATUS_IGNORE);
545     MPI_Send(&seast,1,MPI_INT,east,0,gridcomm);
546     MPI_Recv(&rwest,1,MPI_INT,west,0,gridcomm,MPI_STATUS_IGNORE);

548     MPI_Send(&ssouth,1,MPI_INT,south,0,gridcomm);
549     MPI_Recv(&rnorth,1,MPI_INT,north,0,gridcomm,MPI_STATUS_IGNORE);
550     MPI_Send(&snorth,1,MPI_INT,north,0,gridcomm);
551     MPI_Recv(&rsouth,1,MPI_INT,south,0,gridcomm,MPI_STATUS_IGNORE);

553     MPI_Send(&sabove,1,MPI_INT,above,0,gridcomm);
554     MPI_Recv(&rbelow,1,MPI_INT,below,0,gridcomm,MPI_STATUS_IGNORE);
555     MPI_Send(&sbelow,1,MPI_INT,below,0,gridcomm);
556     MPI_Recv(&rabove,1,MPI_INT,above,0,gridcomm,MPI_STATUS_IGNORE);

558     if(snorth || ssouth || swest || seast || sabove || sbelow){
559       cr = getGlobalCord(x,y,z);

561       #ifdef DEBUG
562       printf("%d: sending element at %d %d %d gave global coord %d %d
              %d\n",cartrank,x,y,z,cr[0],cr[1],cr[2]);
563       #endif
564       e.x = cr[0];
565       e.y = cr[1];
566       e.z = cr[2];
567       e.value = data−>timearray[GETLINDEX(x,y,z)];
568       e.n = data−>bandarray[GETLINDEX(x,y,z)];

570       if(snorth){
```

```
571          MPI_Send(&e,1,mpi_element_struct,north,1,gridcomm);
572        }
573      if(ssouth){
574          MPI_Send(&e,1,mpi_element_struct,south,1,gridcomm);
575        }
576      if(swest){
577          MPI_Send(&e,1,mpi_element_struct,west,1,gridcomm);
578        }
579      if(seast){
580          MPI_Send(&e,1,mpi_element_struct,east,1,gridcomm);
581        }
582      if(sabove){
583          MPI_Send(&e,1,mpi_element_struct,above,1,gridcomm);
584        }
585      if(sbelow){
586          MPI_Send(&e,1,mpi_element_struct,below,1,gridcomm);
587        }
588    }
589
590    if(rnorth || rsouth || rwest || reast || rabove || rbelow){
591
592      if(rnorth){
593          MPI_Recv(&e,1,mpi_element_struct,north,1,gridcomm,
               MPI_STATUS_IGNORE);
594          addElement(data,e);
595        }
596      if(rsouth){
597          MPI_Recv(&e,1,mpi_element_struct,south,1,gridcomm,
               MPI_STATUS_IGNORE);
598          addElement(data,e);
599        }
600      if(rwest){
601          MPI_Recv(&e,1,mpi_element_struct,west,1,gridcomm,
               MPI_STATUS_IGNORE);
602          addElement(data,e);
603        }
604      if(reast){
605          MPI_Recv(&e,1,mpi_element_struct,east,1,gridcomm,
               MPI_STATUS_IGNORE);
606          addElement(data,e);
607        }
608      if(rabove){
609          MPI_Recv(&e,1,mpi_element_struct,above,1,gridcomm,
               MPI_STATUS_IGNORE);
610          addElement(data,e);
611        }
612      if(rbelow){
613          MPI_Recv(&e,1,mpi_element_struct,below,1,gridcomm,
               MPI_STATUS_IGNORE);
614          addElement(data,e);
615        }
616    }
617  }
618  */
619
620  /*
621   * check if someone wants to update their working status
622   */
```

```
623  void checkOthers(){
624    int i;
625    int flag;
626    for(i = 0; i<size;i++){
627      MPI_Iprobe(i,9,gridcomm,&flag,MPI_STATUS_IGNORE);
628      if(flag){
629        MPI_Recv(&working[i],1,MPI_INT,i,9,gridcomm,MPI_STATUS_IGNORE)
               ;
630      }
631    }
632  }
633  /*
634   * notify others that my working status is changed
635   */
636  void notifyOthers(int value){
637    int i;
638    for(i = 0; i<size;i++){
639      MPI_Send(&value,1,MPI_INT,i,9,gridcomm);
640    }
641  }
642
643  /*
644   * Execute the FMM
645   */
646  void executeFMM(FmmData* data){
647    int add=0;
648    int posx,posy,posz;
649    int run = 1;
650    int sendrun = 1;
651    int senddata = 0;
652    int end = 0;
653    int sum = 0;
654    int i;
655    #ifdef DEBUG
656    printf("data_size_is_%d_%d_%d\n",data->x, data->y, data->z);
657    #endif
658    for(i = 0; i<size;i++){
659      working[i] = 1;
660    }
661    largest_solution = 0;
662    n = 1;
663    /* loop will run until all nodes are done */
664    while(!end){
665    /* working loop, will run until there are no more work to be done
          */
666    while(run){
667
668
669      if(heapGetMin(data->heap)){
670
671        Element* e,*temp;
672
673        e = heapExtractMin(data->heap);
674        #ifdef DEBUG
675        printf("%d:_setting_%d_%d_%d_to_known\n",cartrank,e->x,e->y,e->
             z);
676        #endif
677        data->bandarray[GETLINDEX(e->x,e->y,e->z)] = KNOWN;
```

```
678          int k[6] = {e->x-1,e->x,e->x+1,e->x,e->x,e->x};
679          int l[6] = {e->y, e->y-1, e->y, e->y+1,e->y,e->y};
680          int m[6] = {e->z, e->z, e->z, e->z,e->z-1,e->z+1};
681          int i;
682        posx = e->x;
683        posy = e->y;
684        posz = e->z;
685        /*
686         * update largest_solution if this solution is the largest
687         */
688        if(data->timearray[GETLINDEX(posx,posy,posz)] >
               largest_solution){
689         largest_solution = data->timearray[GETLINDEX(posx,posy,posz)];
690        }
691        /* check if this point has been sent before */
692        if(data->sentarray[GETLINDEX(posx,posy,posz)] == 0){
693         senddata = 1;
694        }else if(data->sentarray[GETLINDEX(posx,posy,posz)] <= data->
               timearray[GETLINDEX(posx,posy,posz)]){
695         senddata = 0;
696        }
697
698        /*printFloatArray(data->x,data->y,-2,data->timearray);
699        printFloatArray(data->x,data->y,-1,data->timearray);
700        printFloatArray(data->x,data->y,0,data->timearray);
701        printFloatArray(data->x,data->y,1,data->timearray);
702        printFloatArray(data->x,data->y,2,data->timearray);*/
703        for(i=0;i<6;i++){
704
705         if(k[i] >= 0 && k[i] < data->x && l[i] >= 0 && l[i] < data->y
               && m[i] >= 0 && m[i] < data->z){
706           calcElement(data, k[i], l[i], m[i]);
707         }
708        }
709
710    #ifdef DEBUG
711
712        /*printFloatArray(data->x,data->y,-2,data->timearray);
713        printFloatArray(data->x,data->y,-1,data->timearray);
714        printFloatArray(data->x,data->y,0,data->timearray);
715        printFloatArray(data->x,data->y,1,data->timearray);
716        printFloatArray(data->x,data->y,2,data->timearray);*/
717    #endif
718        free(e);
719
720        n++;
721      }
722      /* Send changes to border and look for incoming changes to the
             border */
723      sendRecvBorderChanges(data,posx,posy,posz,senddata);
724
725    #ifdef DEBUG
726      if(cartrank == 0|| cartrank == -2){printLocalArray(data->
             timearray);
727      printLocalIntArray(data->bandarray);
728      }
729    #endif
730      senddata = 0;
```

```
731        /*Do we end the loop?*/
732        if(heapGetMin(data->heap) == 0){
733          run = 0;
734        }else{
735          run = 1;
736        }
737      checkOthers();
738      #ifdef DEBUG
739      if(cartrank == 0 && n%1000 == 0){
740        printf("%d: reached n %d\n",cartrank,n);
741      }
742      #endif
743    }
744      /* see if border changes are coming */
745      sendRecvBorderChanges(data,0,0,0,0);
746      /* if we have work to do lets notify others and start to work
            again, if not let others know we are done */
747      if(heapGetMin(data->heap) != 0){
748        run = 1;
749        // notify other i am still working
750        notifyOthers(1);
751      }else if(working[cartrank] == 1){
752        // notify that i have stopped working
753        notifyOthers(0);
754      }
755      checkOthers();
756      sum = 0;
757      for(i = 0; i<size;i++){
758        sum += working[i];
759      }
760      if(sum == 0){
761        end = 1;
762      }
763    }
764    #ifdef DEBUG
765    printf("valuemax %f valuemin %f\n",valuemax,valuemin);
766    #endif
767 }
768
769
770 /*
771  * free the variables used in the FMM
772  */
773 void freeFMM(FmmData* data){
774    //free(data->timearray);
775    free(data->bandarray);
776    free(data);
777 }
```

# A.3 Application

## A.3.1 mpi_app.c

```c
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include "array_mpi.h"
4  #include "fileio.h"
5  #include <mpi.h>
6  #include <string.h>
7  #include "fmm_mpi.h"
8  #include "time.h"
9
10 #define MPIDEBUG 0
11 #define DEBUG 0
12
13 /* set that no dimensions sould be cyclic */
14 int periods[3] = {0,0,0};
15 /* set the number of dimensions to use in the cartesian node grid */
16 int ndims = 3;
17
18 int div_x, div_y, div_z;
19
20 float* local_array;
21 float* global_array;
22 float* file_array;
23
24 /*
25  * different mpi datatypes for exchanging the borders
26  */
27 MPI_Datatype xz_plane;
28 MPI_Datatype zy_plane;
29 MPI_Datatype xy_plane;
30 MPI_Datatype y_column;
31 MPI_Datatype y_column_resized;
32 /*
33  * A test function to check if the array dosen't contain a value
34  */
35 int checkArray(int x, int y, int z,float value){
36     int i,j,k;
37     int rvalue = 1;
38     for(i= 0; i<x;i++)
39       for(j =0;j<y;j++)
40         for(k=0;k<z;k++){
41           if(local_array[GETLINDEX(i,j,k)] != value){
42             rvalue = 0;
43             return rvalue;
44           }
45         }
46     return rvalue;
47   }
48 /*
49  * Check if a global array dosen't contain a specific value
50  */
51 int checkGArray(int x, int y, int z,float value){
52     int i,j,k;
53     int rvalue = 1;
```

```
54      for(i= 0; i< x;i++)
55        for(j =0;j <y; j++)
56          for(k=0;k<z;k++){
57            if(global_array[GETINDEX(i,j,k)] != value){
58              rvalue = 0;
59              return rvalue;
60            }
61          }
62      return rvalue;
63   }
64
65  /*
66   * Divide the global matrix into smaller matrixes for each node.
67   * This function calculates the local dimensions
68   */
69  void divide_matrix(){
70
71    div_x = x/dims[0];
72    div_y = y/dims[1];
73    div_z = z/dims[2];
74    local_x = div_x;
75    local_y = div_y;
76    local_z = div_z;
77    if(local_x * dims[0] != x){
78      if(coords[0] == dims[0]){
79        local_x = (x-(local_x*dims[0])) + local_x;
80      }
81    }
82    if(local_y * dims[1] != y){
83      if(coords[1] == dims[1]){
84        local_y = (y-(local_y*dims[1])) + local_y;
85      }
86    }
87    if(local_z * dims[2] != z){
88      if(coords[2] == dims[2]){
89        local_z = (z-(local_z*dims[2])) + local_z;
90      }
91    }
92    printf("%d:_local_x_%d_local_y_%d_local_z_%d_coords_%d,%d,%d_\n",
           cartrank, local_x, local_y, local_z,coords[0],coords[1],coords
           [2]);
93  }
94
95  /*
96   * initialize datatypes for border exchange
97   */
98  void initMPIDatatypes(){
99    MPI_Type_vector(local_y,local_z,local_z+2,MPI_FLOAT,&zy_plane);
100   MPI_Type_commit(&zy_plane);
101
102   MPI_Type_vector(local_x,local_z,(local_z+2)*(local_y+2),MPI_FLOAT
          ,&xz_plane);
103   MPI_Type_commit(&xz_plane);
104
105   MPI_Type_vector(local_y,1,local_z+2,MPI_FLOAT,&y_column);
106   MPI_Type_commit(&y_column);
107   MPI_Type_create_resized(y_column,0,(local_z+2)*(local_y+2)*sizeof(
          float),&y_column_resized);
```

```
108    MPI_Type_vector(1,2,1,y_column_resized,&xy_plane);
109    MPI_Type_commit(&xy_plane);
110  }
111  /*
112   * Initialize the program allocating local matrixes and initializing
            data types
113   */
114  void init(){
115    divide_matrix();
116    local_array = malloc(sizeof(float)*LOCALARRAYSIZE);
117    if(local_array == 0){
118      printf("Coudn't allocate enough memory for local_array\n");
119      exit(1);
120    }
121    bzero(local_array, sizeof(float)*LOCALARRAYSIZE);
122    initMPIDatatypes();
123  }
124  /*
125   * a function for reading a file into each node, where the file
            contains a global array, each node will read its respective part
            into their local matrixes
126   */
127  void scatterdata(char* filename){
128    int i,j;
129    int* cr;
130    int value;
131    char* errorstr;
132    int reslen;
133    FILE* f;
134    int offset = 0;
135    printf("%d: openeing file %s\n",cartrank,filename);
136    f = fopen(filename,"rb");
137    if(!f){
138      printf("%d: unable to open file %d \n",cartrank,f);
139      fflush(stdout);
140      return;
141    }
142    if(DEBUG){
143      printf("%d: opened file %d\n",cartrank,f);
144    }
145    for(i = 0; i<local_x;i++){
146      for(j= 0; j< local_y; j++)
147      {
148        cr = getGlobalCord(i,j,0);
149        offset = sizeof(float)*(cr[2]+cr[1]*z+cr[0]*z*y);
150        fseek(f,offset,SEEK_SET);
151        fread(&local_array[GETLINDEX(i,j,0)],sizeof(float),local_z,f);
152        free(cr);
153      }
154    }
155    fclose(f);
156
157  }
158
159  /*
160   * gather all the local matrixes into a global matrix on node 0
161   */
162  float* gatherdata(float* iarray){
```

```
163    int i,j,k,dest,flag,r,t;
164    float* farray;
165    MPI_Request* requests;
166    MPI_Status* status;
167    requests = malloc(sizeof(MPI_Request)*local_y*local_x*2);
168    status = malloc(sizeof(MPI_Status)*local_y*local_x*2);
169    if(cartrank == 0){
170       if(DEBUG){
171          printf("%d:started gathering\n",cartrank);
172       }
173       farray = malloc(sizeof(float)*ARRAYSIZE);
174       bzero(farray,sizeof(float)*ARRAYSIZE);
175
176    }
177    if(cartrank != 0){
178    for(i=0;i<local_x;i++)
179       for(j= 0; j<local_y;j++){
180          MPI_Send(&iarray[GETLINDEX(i,j,0)],local_z,MPI_FLOAT,0,i*
                 local_y+j,gridcomm);//,&requests[(local_y*i)+j]);
181       }
182
183    if(DEBUG){
184       printf("%d: Done sending \n",cartrank);
185    }
186    }
187    if(cartrank==0){
188       if(MPIDEBUG){
189          printf("%d: starting setting recvs\n",cartrank);
190       }
191       for(r=0;r<x;r+=local_x)
192          for(t=0; t<y;t+=local_y)
193             for(k=0;k<dims[2];k++){
194
195                dest = getDest(r,t,k*local_z);
196                if(DEBUG){
197                   printf("%d: receiving from %d\n",cartrank,dest);
198                }
199                for(i=0;i<local_x;i++)
200                   for(j=0;j<local_y;j++){
201                   MPI_Irecv(&farray[GETINDEX(i+r,j+t,k*local_z)],local_z,
                        MPI_FLOAT,dest,i*local_y+j,gridcomm,&requests[(
                        local_y*local_x)+(i*local_y+j)]);
202                   if(dest == 0){
203                      MPI_Isend(&iarray[GETLINDEX(i,j,0)],local_z,MPI_FLOAT
                           ,0,i*local_y+j,gridcomm,&requests[(local_y*i)+j]);
204                   }
205                   }
206                if(dest == 0){
207                   MPI_Waitall(local_y*local_x*2,requests,status);
208                }
209                if(dest != 0){
210                MPI_Waitall(local_y*local_x,&requests[local_y*local_x],
                        status);
211                }
212             }
213    }
214
215    printf("%d: Done gathering \n",cartrank);
```

```
216    return farray;
217 }
218
219
220 /*
221  * check if the global array is the same as the array inside a file
222  */
223 void checkData(char* filename){
224
225    if(cartrank == 0){
226    int i,j,k;
227       printf("%d:_Reading_file\n",cartrank);
228       file_array = malloc(ARRAYSIZE*sizeof(float));
229       bzero(file_array,ARRAYSIZE*sizeof(float));
230       readfile(file_array,filename,x,y,z);
231
232       printArray(file_array);
233       printf("%d:_checking_data_consistency_%f\n",cartrank,file_array[
           GETINDEX(0,0,0)]);
234       for(i=0;i<x;i++)
235         for(j=0;j<y;j++)
236           for(k=0;k<z;k++){
237              if(file_array[GETINDEX(i,j,k)] != global_array[GETINDEX(i,
                 j,k)]){
238                printf("%d:_error_at_%d_%d_%d_file_%lf_global_%lf\n",
                     cartrank,i,j,k,file_array[GETINDEX(i,j,k)],
                     global_array[GETINDEX(i,j,k)]);
239              }
240           }
241    }
242 }
243
244 /*
245  * Exchange borders
246  */
247 void exchangeBorders(){
248
249
250    // sending/recving north south
251    MPI_Send(&local_array[GETLINDEX(0,0,0)],1,zy_plane,north,0,
           gridcomm);
252    MPI_Recv(&local_array[GETLINDEX(local_x,0,0)],1,zy_plane,south,0,
           gridcomm,MPI_STATUS_IGNORE);
253    MPI_Send(&local_array[GETLINDEX(local_x-1,0,0)],1,zy_plane,south
           ,1,gridcomm);
254    MPI_Recv(&local_array[GETLINDEX(-1,0,0)],1,zy_plane,north,1,
           gridcomm,MPI_STATUS_IGNORE);
255
256
257    //sending/recving east, west
258    MPI_Send(&local_array[GETLINDEX(0,0,0)],1,xz_plane,west,2,gridcomm
           );
259    MPI_Recv(&local_array[GETLINDEX(0,local_y,0)],1,xz_plane,east,2,
           gridcomm,MPI_STATUS_IGNORE);
260    MPI_Send(&local_array[GETLINDEX(0,local_y-1,0)],1,xz_plane,east,3,
           gridcomm);
261    MPI_Recv(&local_array[GETLINDEX(0,-1,0)],1,xz_plane,west,3,
           gridcomm,MPI_STATUS_IGNORE);
```

```
262
263      //sending/receiving above, below
264      MPI_Send(&local_array[GETLINDEX(0,0,local_z−1)],1,xy_plane,above
             ,4,gridcomm);
265      MPI_Recv(&local_array[GETLINDEX(0,0,−1)],1,xy_plane,below,4,
             gridcomm,MPI_STATUS_IGNORE);
266      MPI_Send(&local_array[GETLINDEX(0,0,0)],1,xy_plane,below,5,
             gridcomm);
267      MPI_Recv(&local_array[GETLINDEX(0,0,local_z)],1,xy_plane,above,5,
             gridcomm,MPI_STATUS_IGNORE);
268    }
269
270    /*
271     * a test function for writing a file with values
272     */
273    void writeafile(){
274      FILE* f;
275      f = fopen("/work/idarbo/per.conv","w");
276      int i,j,k;
277      float value;
278      for(i=0;i<x;i++)
279        for(j=0;j<y;j++)
280          for(k=0;k<z;k++){
281          //value = i+j+k;
282          value = 1;
283          if(j < 8 && k < 8){
284            value = 9;
285          }
286
287          fwrite(&value,sizeof(float),1,f);
288          }
289      fclose(f);
290    }
291
292    int main(int argc, char** argv){
293      float* array,*time1,*time2;
294      int timeusec,timesec, timeusec2, timesec2,rtimesec, rtimeusec;
295      FmmData* data;
296          x = SIZEX;
297          y = SIZEY;
298          z = SIZEZ;
299          int i, j, k;
300          /* initialize MPI*/
301      MPI_Init(&argc, &argv);
302      MPI_Comm_rank(MPI_COMM_WORLD, &rank);
303      MPI_Comm_size(MPI_COMM_WORLD, &size);
304      MPI_Dims_create(size,ndims,dims);
305
306      MPI_Cart_create(MPI_COMM_WORLD,ndims, dims, periods, 0, &gridcomm)
             ;
307
308      MPI_Cart_shift(gridcomm, 0, 1, &north, &south);
309
310      MPI_Cart_shift(gridcomm, 1, 1, &west, &east);
311      MPI_Cart_shift(gridcomm, 2, 1, &below, &above);
312
313      MPI_Comm_rank(gridcomm, &cartrank);
314
```

```
315    MPI_Cart_coords(gridcomm, cartrank, 3, coords);
316    //if(cartrank ==0)
317    //writeafile ();
318    /* init the program */
319    init ();
320
321    if (MPIDEBUG) {
322      printf("%d:_west_%d_east_%d_south_%d_north_%d_below_%d_above_%d_
              cartrank_%d\n",rank,west,east,south,north,below,above,
              cartrank);
323      fflush(stdout);
324    }
325
326
327
328    if(cartrank == 0){
329      printf("%d:_Init_complete_reading_data_from_file\n",cartrank);
330      fflush(stdout);
331    }
332    /* read data from file */
333    //scatterdata(argv[1]);
334
335    /* instead of reading set the array to 1.0 */
336    for(i=-1;i<=local_x;i++)
337      for(j=-1;j<=local_y;j++)
338        for(k=-1;k<=local_z;k++)
339          local_array[GETLINDEX(i,j,k)] = 1.0;
340    //printLocalArray(local_array);
341    //MPI_Barrier(gridcomm);
342    if(cartrank == 0){
343      printf("%d:_Read_data_exchanging_borders\n",cartrank);
344      fflush(stdout);
345    }
346    /* exchange borders so border values will be correct in the
            velocity array */
347    exchangeBorders();
348
349    if(cartrank == 0){
350      printf("%d:_Initializing_FMM\n",cartrank);
351      fflush(stdout);
352    }
353    /* initialize the FMM */
354    data = initFMM(local_array,x/2,y/2,z/2,local_x,local_y,local_z);
355
356    if(cartrank == 0){
357      printf("%d:_Executing_FMM\n",cartrank);
358      fflush(stdout);
359    }
360    /* take timing and execute the FMM*/
361    timeusec = getTimeInMicroseconds();
362    timesec = getTimeInSeconds();
363    executeFMM(data);
364    timeusec2 = getTimeInMicroseconds();
365    timesec2 = getTimeInSeconds();
366    fixTime(timesec,timeusec,timesec2,timeusec2,&rtimesec,&rtimeusec);
367    if(cartrank == 0){
368      printf("%d:_execute_FMM_took_%d_seconds_and_%d_microseconds\n",
              cartrank,rtimesec,rtimeusec);
```

```
369     }
370     if (DEBUG) {
371         printLocalArray (data −>timearray ) ;
372     }
373     MPI_Barrier (gridcomm ) ;
374
375     if ( cartrank  ==  0) {
376         printf ( "%d:_FMM_done_gathering_data\n" , cartrank ) ;
377     }
378     /* gather  if  you  can  place  all  data  on  one  node  in  memory ,  else
               write  to  file */
379     // global_array  =  gatherdata ( data −>timearray ) ;
380
381     // global_array  =  gatherdata ( local_array ) ;
382
383
384
385
386     if ( cartrank  ==  0) {
387         if (DEBUG) {
388                 // printArray ( global_array ) ;
389         }
390     }
391     // checkData ( argv [ 1 ] ) ;
392
393     printf ( "%d_is_done_ending_gracefully .\n" , cartrank ) ;
394     MPI_Finalize () ;
395
396  }
```

# Appendix B

# Application 2 source code

## B.1 Array

### B.1.1 array_mpi.h

```
1  #ifndef ARRAY_H
2  #define ARRAY_H
3
4  #include "mpi.h"
5  /*
6   * Defines the size of the array
7   */
8  #define SIZEX 448//448//8//448
9  #define SIZEY 704//704//4//704
10 #define SIZEZ 1216//1216//4//1216
11 /*
12  * Calculate the strides for global array
13  */
14 #define stride0 ((SIZEZ+2)*(SIZEY+2))
15 #define stride1 (SIZEZ+2)
16 #define stride2 1
17 /*
18  * Gets the index for a position in the global array
19  */
20 #define GETINDEX(i,j,k) stride0*(i+1) + stride1*(j+1) + stride2*(k
       +1)
21 /*
22  * Get the index for a position in the local array
23  */
24 #define GETLINDEX(i,j,k) ((local_z+2)*(local_y+2)) *(i+1) + ((
       local_z+2)*(j+1)) + k+1
25 /*
26  * A large float number, should be larger than anything you
        calculate
27  */
28 #define BIGFLOAT 10000000.0
29
30 /*
31  * The size of each array, for mallocing memory
```

```
32    */
33    #define ARRAYSIZE (SIZEX+2)*(SIZEY+2)*(SIZEZ+2)
34    #define LOCALARRAYSIZE (local_x+2)*(local_y+2)*(local_z+2)
35
36    /*
37     * Rank, MPI rank
38     * cartrank, rank in the cartesian grid
39     * size, number of nodes used
40     */
41    int rank,cartrank, size;
42    /*
43     * Ranks of nodes that are above,below,west,east,north,south
44     */
45    int above, below, west, east, north ,south;
46    /*
47     * Number of nodes in each dimension
48     */
49    int dims[3];
50    /*
51     * Different datatypes for local arrays
52     */
53    MPI_Datatype xz_plane;
54    MPI_Datatype zy_plane;
55    MPI_Datatype xy_plane;
56    MPI_Datatype y_column;
57    MPI_Datatype y_column_resized;
58
59    MPI_Datatype int_xz_plane;
60    MPI_Datatype int_zy_plane;
61    MPI_Datatype int_xy_plane;
62    MPI_Datatype int_y_column;
63    MPI_Datatype int_y_column_resized;
64    /*
65     * communicator for the cartesian grid
66     */
67    MPI_Comm gridcomm;
68    /*
69     * MY coordinates in the cartesian node grid
70     */
71    int coords[3];
72    /*
73     * size in each dimension of the global array
74     */
75    int x,y,z;
76    /*
77     * Size in each dimension of the local array
78     */
79    int local_x,local_y,local_z;
80
81    /*
82     * Print the values of a global array to stdout
83     */
84    void printArray(float* array);
85    /*
86     * Print the local array to stdout
87     */
88    void printLocalArray(float* array);
89    /*
```

```
90     * print a local int array to stdout
91     */
92   void printLocalIntArray(int* array);
93   /*
94     * Get what node a global position resides in
95     */
96   int getDest(int x, int y, int z);
97   /*
98     * Get global coordinates from local coordinates
99     */
100  int* getGlobalCord(int x, int y, int z);
101  /*
102    * Get local coordinates from global coordinates
103    */
104  int* getLocalCord(int x, int y, int z);
105  /*
106    * Start exchanging of borders
107    */
108  void exchangeBorders(float* array);
109  /*
110    * Initializa mpi data types
111    */
112  void initMPIDatatypes();
113  /*
114    * Start exchanging integer borders
115    */
116  void int_exchangeBorders(int* array);
117  /*
118    * Initialize integer datatypes
119    */
120  void int_initMPIDatatypes();
121  /*
122    * Wait for all int exchange borders
123    */
124  void wait_exchange_int();
125  /*
126    * wait for all float exchange borders
127    */
128  void wait_exchange_float();
129  /*
130    * Wait for all border exchanges
131    */
132  void wait_exchange_all();
133
134  #endif
```

## B.1.2 array_mpi.c

```c
1  #include "array_mpi.h"
2  #include <stdlib.h>
3  #include <stdio.h>
4  /*
5   * stores requests and status from non blocking border exchange
6   */
7  MPI_Request request[24];
8  MPI_Status status[24];
9  /*
10  * converts global coord to local coord
11  */
12 int* getLocalCord(int xg, int yg, int zg){
13         int *cr;
14         cr = malloc(sizeof(int)*3);
15         cr[0] = xg;
16         cr[1] = yg;
17         cr[2] = zg;
18   cr[0] -= coords[0] * local_x;
19   cr[1] -= coords[1] * local_y;
20   cr[2] -= coords[2] * local_z;
21         return cr;
22 }
23
24 /*
25  * convert local coord to global coord
26  */
27 int* getGlobalCord(int x, int y, int z){
28         int *cr;
29         cr = malloc(sizeof(int)*3);
30         cr[0] = x+ local_x*coords[0];
31         cr[1] = y+ local_y*coords[1];
32         cr[2] = z+ local_z*coords[2];
33         return cr;
34         }
35
36 /*
37  * gets the rank of the node that have input position
38  */
39 int getDest(int x, int y, int z){
40         int cr[3];
41         int value;
42         cr[0] = x/local_x;
43         cr[1] = y/local_y;
44         cr[2] = z/local_z;
45         MPI_Cart_rank(gridcomm,cr,&value);
46         return value;
47 }
48
49 /*
50  * print the local float array to stdout
51  */
52 void printLocalArray(float* array){
53         int i,j,k;
54         for(k=-1;k<=local_z;k++){
55         printf("%d: array_z_=_%d\n",cartrank,k);
56                 for(i=-1;i<=local_x;i++){
```

```
57                           printf("%d:␣",cartrank);
58                           for(j=−1;j<=local_y;j++){
59                                   printf("␣%8f␣",array[GETLINDEX(i,j,k
                                       )]);
60                           }
61                           printf("\n");
62                   }
63           printf("\n\n");
64           }
65  }
66  /*
67   * print a local integer array to stdout
68   */
69  void printLocalIntArray(int* array){
70           int i,j,k;
71           for(k=−1;k<=local_z;k++){
72           printf("%d:array␣z␣=␣%d\n",cartrank,k);
73                   for(i=−1;i<=local_x;i++){
74                           printf("%d:␣",cartrank);
75                           for(j=−1;j<=local_y;j++){
76                                   printf("␣%d␣",array[GETLINDEX(i,j,k)
                                       ]);
77                           }
78                           printf("\n");
79                   }
80           printf("\n\n");
81           }
82  }
83
84  /*
85   * print the global array to stdout
86   */
87  void printArray(float* array){
88           int i,j,k;
89
90           for(k=−1;k<=z;k++){
91           printf("array␣z␣=␣%d\n",k);
92                   for(i=−1;i<=x;i++){
93                           for(j=−1;j<=y;j++){
94                                   printf("␣%8f␣",array[GETINDEX(i,j,k)
                                       ]);
95                           }
96                   printf("\n");
97                   }
98           printf("\n\n");
99           }
100 }
101
102 /*
103  * start non blocking exchange border communication
104  */
105 void exchangeBorders(float* array){
106
107
108          // sending /recving north south
109          MPI_Isend(&array[GETLINDEX(0,0,0)],1,zy_plane,north,0,
                 gridcomm,&request[0]);
110          MPI_Irecv(&array[GETLINDEX(local_x,0,0)],1,zy_plane,south,0,
```

```
            gridcomm,&request[1]); //MPI_STATUS_IGNORE);
111         MPI_Isend(&array[GETLINDEX(local_x-1,0,0)],1,zy_plane,south
                ,1,gridcomm,&request[2]);
112         MPI_Irecv(&array[GETLINDEX(-1,0,0)],1,zy_plane,north,1,
                gridcomm,&request[3]); //MPI_STATUS_IGNORE);
113
114
115         //sending/recving east, west
116         MPI_Isend(&array[GETLINDEX(0,0,0)],1,xz_plane,west,2,
                gridcomm,&request[4]);
117         MPI_Irecv(&array[GETLINDEX(0,local_y,0)],1,xz_plane,east,2,
                gridcomm,&request[5]); //MPI_STATUS_IGNORE);
118         MPI_Isend(&array[GETLINDEX(0,local_y-1,0)],1,xz_plane,east
                ,3,gridcomm,&request[6]);
119         MPI_Irecv(&array[GETLINDEX(0,-1,0)],1,xz_plane,west,3,
                gridcomm,&request[7]); //MPI_STATUS_IGNORE);
120
121         //sending/receiving above, below
122         MPI_Isend(&array[GETLINDEX(0,0,local_z-1)],1,xy_plane,above
                ,4,gridcomm,&request[8]);
123         MPI_Irecv(&array[GETLINDEX(0,0,-1)],1,xy_plane,below,4,
                gridcomm,&request[9]); //MPI_STATUS_IGNORE);
124         MPI_Isend(&array[GETLINDEX(0,0,0)],1,xy_plane,below,5,
                gridcomm,&request[10]);
125         MPI_Irecv(&array[GETLINDEX(0,0,local_z)],1,xy_plane,above,5,
                gridcomm,&request[11]); //MPI_STATUS_IGNORE);
126 }
127
128 /*
129  * initialize mpi datatypes
130  */
131 void initMPIDatatypes(){
132         MPI_Type_vector(local_y,local_z,local_z+2,MPI_FLOAT,&
                zy_plane);
133         MPI_Type_commit(&zy_plane);
134
135         MPI_Type_vector(local_x,local_z,(local_z+2)*(local_y+2),
                MPI_FLOAT,&xz_plane);
136         MPI_Type_commit(&xz_plane);
137
138         MPI_Type_vector(local_y,1,local_z+2,MPI_FLOAT,&y_column);
139         MPI_Type_commit(&y_column);
140         MPI_Type_create_resized(y_column,0,(local_z+2)*(local_y+2)*
                sizeof(float),&y_column_resized);
141         MPI_Type_vector(1,2,1,y_column_resized,&xy_plane);
142         MPI_Type_commit(&xy_plane);
143 }
144
145 /*
146  * start exchanging integer border
147  */
148 void int_exchangeBorders(int* array){
149
150
151         // sending/recving north south
152         MPI_Isend(&array[GETLINDEX(0,0,0)],1,int_zy_plane,north,6,
                gridcomm,&request[12]);
153         MPI_Irecv(&array[GETLINDEX(local_x,0,0)],1,int_zy_plane,
```

```
                        south ,6 , gridcomm,& request [13]) ; //MPI_STATUS_IGNORE) ;
154             MPI_Isend(&array [GETLINDEX( local_x −1,0,0)] ,1 , int_zy_plane ,
                        south ,7 , gridcomm,& request [14]) ;
155             MPI_Irecv(&array [GETLINDEX( −1,0,0)] ,1 , int_zy_plane , north ,7 ,
                        gridcomm,& request [15]) ; //MPI_STATUS_IGNORE) ;
156
157
158             // printf("%d: sending to %d reciving from %d\n", cartrank ,
                        west , east ) ;
159             // sending / recving east , west
160
161
162             MPI_Isend(&array [GETLINDEX(0 ,0 ,0)] ,1 , int_xz_plane , west ,8 ,
                        gridcomm,& request [16]) ;
163             MPI_Irecv(&array [GETLINDEX(0 , local_y ,0)] ,1 , int_xz_plane , east
                        ,8 , gridcomm,& request [17]) ; //MPI_STATUS_IGNORE) ;
164             MPI_Isend(&array [GETLINDEX(0 , local_y −1,0)] ,1 , int_xz_plane ,
                        east ,9 , gridcomm,& request [18]) ;
165             MPI_Irecv(&array [GETLINDEX(0 , −1,0)] ,1 , int_xz_plane , west ,9 ,
                        gridcomm,& request [19]) ; //MPI_STATUS_IGNORE) ;
166
167             // sending / receiving above , below
168
169             MPI_Isend(&array [GETLINDEX(0 ,0 , local_z −1)] ,1 , int_xy_plane ,
                        above ,10 , gridcomm,& request [20]) ;
170             MPI_Irecv(&array [GETLINDEX(0 ,0 , −1)] ,1 , int_xy_plane , below ,10 ,
                        gridcomm,& request [21]) ; //MPI_STATUS_IGNORE) ;
171             MPI_Isend(&array [GETLINDEX(0 ,0 ,0)] ,1 , int_xy_plane , below ,11 ,
                        gridcomm,& request [22]) ;
172             MPI_Irecv(&array [GETLINDEX(0 ,0 , local_z )] ,1 , int_xy_plane ,
                        above ,11 , gridcomm,& request [23]) ; //MPI_STATUS_IGNORE) ;
173     }
174
175     /*
176      * initialise integer mpi datatypes
177      */
178     void int_initMPIDatatypes (){
179             MPI_Type_vector( local_y , local_z , local_z +2,MPI_INT,&
                        int_zy_plane ) ;
180             MPI_Type_commit(&int_zy_plane ) ;
181
182             MPI_Type_vector( local_x , local_z ,( local_z +2)*( local_y +2) ,
                        MPI_INT,&int_xz_plane ) ;
183             MPI_Type_commit(&int_xz_plane ) ;
184
185             MPI_Type_vector( local_y ,1 , local_z +2,MPI_INT,&int_y_column ) ;
186             MPI_Type_commit(&int_y_column ) ;
187             MPI_Type_create_resized ( int_y_column ,0 ,( local_z +2)*( local_y
                        +2)*sizeof( int ),&int_y_column_resized ) ;
188             MPI_Type_vector(1 ,2 ,1 , int_y_column_resized ,&int_xy_plane ) ;
189             MPI_Type_commit(&int_xy_plane ) ;
190     }
191
192     /*
193      * Wait for int border exchange to finish
194      */
195     void wait_exchange_int (){
196       MPI_Waitall(12 ,& request [12] ,& status [12]) ;
```

```
197  }
198  /*
199   * wait for float border exchange to finish
200   */
201  void wait_exchange_float(){
202    MPI_Waitall(12,request,status);
203  }
204
205  /*
206   * wait for all border exchanges to finish
207   */
208  void wait_exchange_all(){
209    MPI_Waitall(24,request,status);
210  }
```

## B.2 Fast Marching Method

### B.2.1 fmm_mpi.h

```
1   #include "heap.h"
2
3   typedef
4   /*
5    * a struct to store all variables for a given fmm implementation
6    */
7   struct f {
8     Heap* heap; // the heap
9     float* timearray; // the timearray storing arrivaltimes
10    int* bandarray; // storing band information
11    float* velocityarray; // velocity field
12    int x,y,z,posx,posy,posz; // x,y,z is size of local fmm matrix,
          posx, posy, posz is global coords for the starting point
13    } FmmData;
14
15  /*
16   * initilalize the FMM set velocity arrat, position of starting
          point and size of array
17   */
18  FmmData* initFMM(float* velarray,int posx,int posy, int posz, int x,
          int y, int z);
19  /*
20   * Free up used variables in FMM
21   */
22  void freeFMM(FmmData* data);
23  /*
24   * Execute the FMM
25   */
26  void executeFMM(FmmData* data);
```

## B.2.2   fmm_mpi.c

```
1   #include <stdlib.h>
2   #include <string.h>
3   #include <stdio.h>
4   #include <math.h>
5   #include "fmm_mpi.h"
6   #include "heap.h"
7   #include "array_mpi.h"
8
9
10  #define BAND 0
11  #define OUTSIDE −1
12  #define KNOWN n
13
14  //#define DEBUG
15  /*
16   * number of loops
17   */
18  int n;
19
20  /*
21   * largest_solution is the largest value in the array
22   * rollbacksmallest is the smallest value of received border values,
            which dictates the rollback number
23   */
24  float largest_solution , rollbacksmallest;
25
26  /*
27   * the n value which one should rollback to
28   */
29  int rollbackn;
30  /*
31   * mpi datatype for sending border points
32   */
33  MPI_Datatype mpi_element_struct;
34  /*
35   * Used for debug output
36   */
37  float valuemax=0;
38  float valuemin=0;
39  /*
40   * array containing the working status of each node
41   */
42  int* working;
43
44  /*
45   * Add a point to the heap
46   */
47  void addToHeap(FmmData* data , int px, int py, int pz){
48      Element *temp;
49      #ifdef DEBUG
50      printf ("%d:_adding_%d,_%d,_%d_to_heap,heap_size_is_%d_maxsize_if_%
            d\n" ,cartrank ,px,py,pz,data−>heap−>heapsize ,data−>heap−>maxsize
            );
51      #endif
52      temp= malloc (sizeof(Element));
53      temp−>value = data−>timearray [GETLINDEX(px,py,pz) ];
```

```
54    temp->x = px;
55    temp->y = py;
56    temp->z = pz;
57    heapInsert(data->heap,temp);
58    #ifdef DEBUG
59    printf("%d:_inserted_%d,_%d,_%d_to_heap,heap_size_is_%d\n",
          cartrank,px,py,pz,data->heap->heapsize);
60    #endif
61
62  }
63
64  /*
65   * initialize the FMM
66   * velarray is the velocity field
67   * posx,y,z is position of the starting point
68   * x,y,z is the size of the array, most are read from array_mpi.h
69   */
70  FmmData* initFMM(float* velarray, int posx, int posy, int posz, int x
        , int y, int z){
71    working = malloc(sizeof(int) * size);
72    /* init datatypes */
73    MPI_Element e;
74    MPI_Datatype type[5] = { MPI_FLOAT, MPI_INT, MPI_INT, MPI_INT,
          MPI_INT };
75    int blocklen[5] = { 1, 1, 1, 1 , 1};
76    MPI_Aint disp[5];
77    disp[0] = 0;
78    disp[1] = sizeof(float);
79    disp[2] = sizeof(int) + disp[1];
80    disp[3] = sizeof(int) + disp[2];
81    disp[4] = sizeof(int) + disp[3];
82    MPI_Type_create_struct(5, blocklen, disp, type, &
          mpi_element_struct);
83    MPI_Type_commit(&mpi_element_struct);
84    /* end init datatypes */
85    #ifdef DEBUG
86    printf("%d:_done_init_mpi_datatypes\n",cartrank);
87    #endif
88    int send = 0;
89    FmmData* data;
90    data = malloc(sizeof(FmmData));
91    if(data == 0){
92      printf("%d:_Failed_to_allocate_memory,_exiting\n",cartrank);
93      exit(1);
94    }
95
96    data->bandarray = malloc(sizeof(int)*LOCALARRAYSIZE);
97    if(data->bandarray == 0){
98      printf("%d:_Failed_to_allocate_memory,_exiting\n",cartrank);
99      exit(1);
100   }
101   data->timearray = malloc(sizeof(float)*LOCALARRAYSIZE);
102   if(data->timearray == 0){
103     printf("%d:_Failed_to_allocate_memory,_exiting\n",cartrank);
104     exit(1);
105   }
106
107   data->velocityarray = velarray;
```

```
108    #ifdef DEBUG
109    printf("%d:_initializing_heap\n",cartrank);
110    #endif
111    data->heap = initHeap(LOCALARRAYSIZE);
112    data->x = x;
113    data->y = y;
114    data->z = z;
115    data->posx = posx;
116    data->posy = posy;
117    data->posz = posz;
118
119    #ifdef DEBUG
120    printf("%d:_clearing_memory_bandarray_%p_size_%d\n",cartrank,data
           ->bandarray,sizeof(int)*LOCALARRAYSIZE);
121    fflush(stdout);
122    #endif
123    /*
124     * setting the band array to outside
125     */
126    memset(data->bandarray,OUTSIDE,sizeof(int)*LOCALARRAYSIZE);
127    #ifdef DEBUG
128    printf("%d:_Heap_ok,_clearing_memory_timearray_%p_size_%d\n",
           cartrank,data->timearray,sizeof(float)*LOCALARRAYSIZE);
129    #endif
130    /*
131     * zeroing the arrival time array
132     */
133    int i = 0;
134    for(i=0;i<LOCALARRAYSIZE;i++)
135      data->timearray[i] = 0;
136    //memset(data->timearray,0,sizeof(float)*LOCALARRAYSIZE);
137    //bzero(data->timearray,sizeof(float)*LOCALARRAYSIZE);
138    data->sentarray = malloc(sizeof(float)*LOCALARRAYSIZE);
139    bzero(data->sentarray,sizeof(float)*LOCALARRAYSIZE);
140    #ifdef DEBUG
141    printf("%d:_done_allocating_memory,_setting_starting_point\n",
           cartrank);
142    fflush(stdout);
143    #endif
144    /*
145     * inserting the starting point on the correct node and add it to
            the heap / narrow band
146     */
147    if(getDest(posx,posy,posz) == cartrank){
148      int *cr = getLocalCord(posx,posy,posz);
149      // insert starting point
150      data->timearray[GETLINDEX(cr[0],cr[1],cr[2])] = 0;
151      data->bandarray[GETLINDEX(cr[0],cr[1],cr[2])] = BAND;
152      //printf("%d: sat pos %d %d %d, as known\n",cartrank,cr[0],cr
              [1],cr[2]);
153      Element* element;
154      element = malloc(sizeof(Element));
155      element->x = cr[0];
156      element->y = cr[1];
157      element->z = cr[2];
158      element->value = data->timearray[GETLINDEX(cr[0],cr[1],cr[2])];
159      heapInsert(data->heap,element);
160      send = 1;
```

```
161      }
162      return data;
163  }
164  /*
165   * return min
166   */
167  float min(float per, float truls){
168      if(per > truls)
169          return truls;
170      else return per;
171  }
172  /*
173   * return max
174   */
175  float max(float per, float truls){
176      if(per > truls)
177          return per;
178      else return truls;
179  }
180
181  /*
182   * prints xy plane from a float array
183   */
184  void printFloatArray(int sizex, int sizey, int z,float* array){
185      int i,j;
186      printf("\n_Printing_matrix\n");
187      for(i=0;i<sizex;i++){
188          for(j=0;j<sizey;j++){
189              printf("_%8.2f_",array[GETLINDEX(i,j,z)]);
190          }
191          printf("\n");
192      }
193  }
194
195  /*
196   * calculate the arrival time for a point x,y,z
197   */
198  float calcDistance(FmmData* data, int x, int y, int z){
199      float sol;
200      sol = BIGFLOAT;
201      if(data->bandarray[GETLINDEX(x+1,y,z)] > BAND){
202          sol = min(data->timearray[GETLINDEX(x+1,y,z)] +1/data->
                  velocityarray[GETLINDEX(x,y,z)],sol);
203          #ifdef DEBUG
204          printf("%d:_sol_is_%f_for_x+1\n",cartrank,data->timearray[
                  GETLINDEX(x+1,y,z)] +1/data->velocityarray[GETLINDEX(x,y,z)])
                  ;
205          #endif
206      }
207      if(data->bandarray[GETLINDEX(x-1,y,z)] > BAND){
208          sol = min(data->timearray[GETLINDEX(x-1,y,z)] +1/data->
                  velocityarray[GETLINDEX(x,y,z)],sol);
209          #ifdef DEBUG
210          printf("%d:_sol_is_%f_for_x-1\n",cartrank,data->timearray[
                  GETLINDEX(x-1,y,z)] +1/data->velocityarray[GETLINDEX(x,y,z)])
                  ;
211          #endif
212      }
```

```
213       if (data->bandarray [GETLINDEX(x,y+1,z)] > BAND){
214          sol = min(data->timearray [GETLINDEX(x,y+1,z)] +1/data->
                 velocityarray [GETLINDEX(x,y,z)],sol);
215          #ifdef DEBUG
216          printf ("%d:_sol_is_%f_for_y+1\n",cartrank,data->timearray [
                 GETLINDEX(x,y+1,z)] +1/data->velocityarray [GETLINDEX(x,y,z)])
                 ;
217          #endif
218       }
219       if (data->bandarray [GETLINDEX(x,y-1,z)] > BAND){
220          sol = min(data->timearray [GETLINDEX(x,y-1,z)] +1/data->
                 velocityarray [GETLINDEX(x,y,z)],sol);
221          #ifdef DEBUG
222          printf ("%d:_sol_is_%f_for_y-1\n",cartrank,data->timearray [
                 GETLINDEX(x,y-1,z)] +1/data->velocityarray [GETLINDEX(x,y,z)])
                 ;
223          #endif
224       }
225       if (data->bandarray [GETLINDEX(x,y,z+1)] > BAND){
226          sol = min(data->timearray [GETLINDEX(x,y,z+1)] +1/data->
                 velocityarray [GETLINDEX(x,y,z)],sol);
227          #ifdef DEBUG
228          printf ("%d:_sol_is_%f_for_z+1\n",cartrank,data->timearray [
                 GETLINDEX(x,y,z+1)] +1/data->velocityarray [GETLINDEX(x,y,z)])
                 ;
229          #endif
230       }
231       if (data->bandarray [GETLINDEX(x,y,z-1)] > BAND){
232          sol = min(data->timearray [GETLINDEX(x,y,z-1)] +1/data->
                 velocityarray [GETLINDEX(x,y,z)],sol);
233          #ifdef DEBUG
234          printf ("%d:_sol_is_%f_for_z-1\n",cartrank,data->timearray [
                 GETLINDEX(x,y,z-1)] +1/data->velocityarray [GETLINDEX(x,y,z)])
                 ;
235          #endif
236       }
237
238       return sol;
239    }
240
241    /*
242     * calculate a new point the the arrival time array if its inside
            the array and not KNOWN, if its outside add it to the heap
243     */
244    void calcElement(FmmData* data, int px, int py, int pz){
245             int add = 0;
246             float sol;
247             Element *temp;
248             /*
249              * Check if the point is inside the array */
250             if (px >= 0 && px < data->x && py >= 0 && py < data->y  && pz
                     >= 0 && pz < data->z){
251             /* Make sure the point is not KNOWN */
252             if (data->bandarray [GETLINDEX(px,py,pz)] <= BAND){
253                     sol = calcDistance (data,px,py,pz);
254                     /* Check if the number is not smaller than the one
                             we calculated , when we exit , should never happend
                             in serial version */
```

```
255                        if (data−>timearray [GETLINDEX(px,py,pz)] != 0 && data
                               −>timearray [GETLINDEX(px,py,pz)] <= sol){
256                                return;
257                        }
258                        /* If the point is OUSIDE add it to the heap/
                               narrowband */
259                        if (data−>bandarray [GETLINDEX(px,py,pz)] == OUTSIDE)
                               {
260                                data−>bandarray [GETLINDEX(px,py,pz)] = BAND;
261                                add = 1;
262                        }
263
264                        #ifdef DEBUG
265                        printf("%d:_sol_is_%f,_for_%d_%d_%d\n",cartrank,sol,
                               px,py,pz);
266                        #endif
267                        /* store max and min for debug purposes */
268                        #ifdef DEBUG
269                        if (sol > valuemax && sol != BIGFLOAT){
270                                valuemax=sol;
271                        }
272                        if (sol < valuemin){
273                                valuemin = sol;
274                        }
275                        #endif
276                        /* set the new arrivaltime */
277                        data−>timearray [GETLINDEX(px,py,pz)] = sol;
278                        #ifdef DEBUG
279                        printf("%d:_x_%d,_,y_%d,_z_%d,_k[i]_%d_l[i]_%d_m[i]_
                               %d,_value_%f_\n",cartrank,data−>x, data−>y, data
                               −>z,px,py,pz, data−>timearray [GETLINDEX(px,py,pz)
                               ]);
280                        #endif
281                        /* add it to the narrowband if it should be added */
282                        if (add){
283                                addToHeap(data,px,py,pz);
284                                }
285                        }
286                }
287 }
288
289 //remove ???
290 void checkforchange (FmmData∗ data ,int px, int py, int pz){
291    float sol = BIGFLOAT;
292    int i;
293    int o[6] = {px−1,px,px+1,px,px,px};
294    int l[6] = {py,py−1, py,py+1,py,py};
295    int m[6] = {pz,pz,pz,pz,pz−1,pz+1};
296    for(i=0;i<6;i++){
297      if(o[i] >= 0 && o[i] < data−>x && l[i] >= 0 && l[i] < data−>y
             && m[i] >= 0 && m[i] < data−>z && data−>bandarray [GETLINDEX(o
             [i],l[i],m[i])] > BAND){
298
299          #ifdef DEBUG
300          printf("%d:_checking_for_rollback_%d_%d_%d_\n",cartrank,o[i],l
                 [i],m[i]);
301          #endif
302          sol = calcDistance (data,px,py,pz);
```

```
303
304        if(sol < data->timearray[GETLINDEX(o[i],l[i],m[i])] && data->
               bandarray[GETLINDEX(o[i],l[i],m[i])] > BAND){
305            // new value is smaller lets add this to our heap
306
307            #ifdef DEBUG
308            printf("%d:_rolling_back_%d_%d_%d_old_value_%f_new_value_%f\
                   n",cartrank,o[i],l[i],m[i],data->timearray[GETLINDEX(o[i
                   ],l[i],m[i])],sol);
309            #endif
310            data->timearray[GETLINDEX(o[i],l[i],m[i])] = sol;
311            data->bandarray[GETLINDEX(o[i],l[i],m[i])] = BAND;
312            addToHeap(data,o[i],l[i],m[i]);
313            checkforchange(data,o[i],l[i],m[i]);
314        }
315      }
316    }
317 }
318
319 /*
320  * Add a new element to the array
321  */
322 void addElement(FmmData* data,MPI_Element e){
323    int* cr;
324    int i,j,k;
325    cr = getLocalCord(e.x,e.y,e.z);
326
327    #ifdef DEBUG
328    printf("%d:_adding_element_%d_%d_%d,_to_local_%d_%d_%d\n",cartrank
               ,e.x,e.y,e.z,cr[0],cr[1],cr[2]);
329    #endif
330    /*
331     * checking if we have added it before
332     * Should be always no since we don't send the same value multiple
               times
333     */
334    if(data->timearray[GETLINDEX(cr[0],cr[1],cr[2])] == e.value){
335      #ifdef DEBUG
336      printf("%d:_alreaddy_added_%d_%d_%d\n",cartrank,e.x,e.y,e.z);
337      #endif
338      return;
339    }
340    data->timearray[GETLINDEX(cr[0],cr[1],cr[2])] = e.value;
341    data->bandarray[GETLINDEX(cr[0],cr[1],cr[2])] = e.n;
342    int o[6] = {cr[0]-1,cr[0],cr[0]+1,cr[0],cr[0],cr[0]};
343    int l[6] = {cr[1],cr[1]-1, cr[1],cr[1]+1,cr[1],cr[1]};
344    int m[6] = {cr[2],cr[2],cr[2],cr[2],cr[2]-1,cr[2]+1};
345    /*
346     * recalculate all neighbours
347     */
348    for(i=0;i<6;i++){
349
350      #ifdef DEBUG
351      printf("%d:_calc_element_%d_%d_%d\n",cartrank,o[i],l[i],m[i]);
352      #endif
353      calcElement(data, o[i], l[i], m[i]);
354    }
355    /*
```

```
356        * set rollback values, so we can check if a rollback is necesarry
357        */
358       /*if(largest_solution > e.value && rollbacksmallest > e.value){
359       rollbacksmallest = e.value;
360       rollbackn = e.n;
361       }*/
362    }
363
364    /*
365     * rollback all values above input value
366     */
367    void rollback(FmmData* data, float value){
368      int i,j,k;
369      for(i=0;i<local_x;i++)
370        for(j=0;j<local_y;j++)
371          for(k=0;k<local_z;k++){
372            if(data->bandarray[GETLINDEX(i,j,k)] > BAND && data->
                  timearray[GETLINDEX(i,j,k)] > value){
373
374              data->bandarray[GETLINDEX(i,j,k)] = BAND;
375              addToHeap(data,i,j,k);
376            }
377          }
378    }
379
380    /*
381     * send new values and check for incoming border values
382     */
383    void sendRecvBorderChanges(FmmData* data,int x, int y, int z,int
        send){
384      MPI_Element e;
385      int run= 1;
386      int *cr;
387      int reast,rwest,rnorth,rsouth,rabove,rbelow;
388      int seast,swest,snorth,ssouth,sabove,sbelow;
389      reast = rwest = rnorth = rsouth = rabove = rbelow = 0;
390      seast = swest = snorth = ssouth = sabove= sbelow = 0;
391      /*
392       * if we are to send a value
393       */
394      if(send){
395        /*
396         * see where we have to send the value
397         */
398        if(x == 0){
399          snorth = 1;
400        }
401        if(x == local_x -1){
402          ssouth = 1;
403        }
404        if(y == 0){
405          swest = 1;
406        }
407        if(y == local_y -1){
408          seast = 1;
409        }
410        if(z == 0){
411          sbelow = 1;
```

```
412        }
413        if(z == local_z −1){
414          sabove = 1;
415        }
416        #ifdef DEBUG
417        printf("%d:_%d_%d_%d_sending_to_north_%d_%d_south_%d__%d_west_%d
                _%d_east_%d_%d_below_%d_%d_above_%d_%d\n",cartrank,x,y,z,
                snorth,north,ssouth,south,swest,west,seast,east,sbelow,below,
                sabove,above);
418        #endif
419
420
421      if(snorth || ssouth || swest || seast || sabove || sbelow){
422        cr = getGlobalCord(x,y,z);
423        #ifdef DEBUG
424        printf("%d:_sending_element_at_%d_%d_%d_gave_global_coord_%d_%d_
                %d\n",cartrank,x,y,z,cr[0],cr[1],cr[2]);
425        #endif
426        e.x = cr[0];
427        e.y = cr[1];
428        e.z = cr[2];
429        e.value = data−>timearray[GETLINDEX(x,y,z)];
430        data−>sentarray[GETLINDEX(x,y,z)] = e.value;
431        e.n = data−>bandarray[GETLINDEX(x,y,z)];
432
433        if(snorth){
434          MPI_Send(&e,1,mpi_element_struct,north,1,gridcomm);
435        }
436        if(ssouth){
437          MPI_Send(&e,1,mpi_element_struct,south,1,gridcomm);
438        }
439        if(swest){
440          MPI_Send(&e,1,mpi_element_struct,west,1,gridcomm);
441        }
442        if(seast){
443          MPI_Send(&e,1,mpi_element_struct,east,1,gridcomm);
444        }
445        if(sabove){
446          MPI_Send(&e,1,mpi_element_struct,above,1,gridcomm);
447        }
448        if(sbelow){
449          MPI_Send(&e,1,mpi_element_struct,below,1,gridcomm);
450        }
451      }
452      }
453      rollbacksmallest = BIGFLOAT;
454      rollbackn = 0;
455      /*
456       * a loop to receive all incoming border values
457       */
458      while(run){
459      MPI_Iprobe(MPI_ANY_SOURCE,1,gridcomm,&run,MPI_STATUS_IGNORE);
460      if(run){
461        MPI_Iprobe(north,1,gridcomm,&rnorth,MPI_STATUS_IGNORE);
462        MPI_Iprobe(south,1,gridcomm,&rsouth,MPI_STATUS_IGNORE);
463        MPI_Iprobe(east,1,gridcomm,&reast,MPI_STATUS_IGNORE);
464        MPI_Iprobe(west,1,gridcomm,&rwest,MPI_STATUS_IGNORE);
465        MPI_Iprobe(above,1,gridcomm,&rabove,MPI_STATUS_IGNORE);
```

```
466        MPI_Iprobe(below,1,gridcomm,&rbelow,MPI_STATUS_IGNORE);
467        if(rnorth || rsouth || rwest || reast || rabove || rbelow){
468
469          if(rnorth){
470            MPI_Recv(&e,1,mpi_element_struct,north,1,gridcomm,
                   MPI_STATUS_IGNORE);
471            addElement(data,e);
472          }
473          if(rsouth){
474            MPI_Recv(&e,1,mpi_element_struct,south,1,gridcomm,
                   MPI_STATUS_IGNORE);
475            addElement(data,e);
476          }
477          if(rwest){
478            MPI_Recv(&e,1,mpi_element_struct,west,1,gridcomm,
                   MPI_STATUS_IGNORE);
479            addElement(data,e);
480          }
481          if(reast){
482            MPI_Recv(&e,1,mpi_element_struct,east,1,gridcomm,
                   MPI_STATUS_IGNORE);
483            addElement(data,e);
484          }
485          if(rabove){
486            MPI_Recv(&e,1,mpi_element_struct,above,1,gridcomm,
                   MPI_STATUS_IGNORE);
487            addElement(data,e);
488          }
489          if(rbelow){
490            MPI_Recv(&e,1,mpi_element_struct,below,1,gridcomm,
                   MPI_STATUS_IGNORE);
491            addElement(data,e);
492          }
493        }
494      }
495      }
496      /*
497       * rollback if necessary
498       */
499      if(rollbackn){
500        rollback(data,rollbacksmallest);
501        n = rollbackn;
502        rollbackn = 0;
503        rollbacksmallest = BIGFLOAT;
504      }
505 }
506
507 /*
508  * old synchroneous border exhcange
509  */
510 /*
511 void sendRecvBorderChanges(FmmData* data,int x, int y, int z,int
       send){
512   MPI_Element e;
513
514   int *cr;
515   int reast,rwest,rnorth,rsouth,rabove,rbelow;
516   int seast,swest,snorth,ssouth,sabove,sbelow;
```

```
517    reast = rwest = rnorth = rsouth = rabove = rbelow = 0;
518    seast = swest = snorth = ssouth = sabove= sbelow = 0;
519    if(send){
520      if(x == 0){
521        snorth = 1;
522      }
523      if(x == local_x −1){
524        ssouth = 1;
525      }
526      if(y == 0){
527        swest = 1;
528        }
529      if(y == local_y −1){
530        seast = 1;
531      }
532      if(z == 0){
533        sbelow = 1;
534      }
535      if(z == local_z −1){
536        sabove = 1;
537      }
538      #ifdef DEBUG
539      printf("%d: %d %d %d sending to north %d %d south %d  %d west %d
             %d east %d %d below %d %d above %d %d\n",cartrank,x,y,z,
             snorth,north,ssouth,south,swest,west,seast,east,sbelow,below,
             sabove,above);
540      #endif
541    }
542
543    MPI_Send(&swest,1,MPI_INT,west,0,gridcomm);
544    MPI_Recv(&reast,1,MPI_INT,east,0,gridcomm,MPI_STATUS_IGNORE);
545    MPI_Send(&seast,1,MPI_INT,east,0,gridcomm);
546    MPI_Recv(&rwest,1,MPI_INT,west,0,gridcomm,MPI_STATUS_IGNORE);
547
548    MPI_Send(&ssouth,1,MPI_INT,south,0,gridcomm);
549    MPI_Recv(&rnorth,1,MPI_INT,north,0,gridcomm,MPI_STATUS_IGNORE);
550    MPI_Send(&snorth,1,MPI_INT,north,0,gridcomm);
551    MPI_Recv(&rsouth,1,MPI_INT,south,0,gridcomm,MPI_STATUS_IGNORE);
552
553    MPI_Send(&sabove,1,MPI_INT,above,0,gridcomm);
554    MPI_Recv(&rbelow,1,MPI_INT,below,0,gridcomm,MPI_STATUS_IGNORE);
555    MPI_Send(&sbelow,1,MPI_INT,below,0,gridcomm);
556    MPI_Recv(&rabove,1,MPI_INT,above,0,gridcomm,MPI_STATUS_IGNORE);
557
558    if(snorth || ssouth || swest || seast || sabove || sbelow){
559      cr = getGlobalCord(x,y,z);
560
561      #ifdef DEBUG
562      printf("%d: sending element at %d %d %d gave global coord %d %d
             %d\n",cartrank,x,y,z,cr[0],cr[1],cr[2]);
563      #endif
564      e.x = cr[0];
565      e.y = cr[1];
566      e.z = cr[2];
567      e.value = data−>timearray[GETLINDEX(x,y,z)];
568      e.n = data−>bandarray[GETLINDEX(x,y,z)];
569
570      if(snorth){
```

```
571          MPI_Send(&e,1,mpi_element_struct,north,1,gridcomm);
572        }
573      if(ssouth){
574          MPI_Send(&e,1,mpi_element_struct,south,1,gridcomm);
575        }
576      if(swest){
577          MPI_Send(&e,1,mpi_element_struct,west,1,gridcomm);
578        }
579      if(seast){
580          MPI_Send(&e,1,mpi_element_struct,east,1,gridcomm);
581        }
582      if(sabove){
583          MPI_Send(&e,1,mpi_element_struct,above,1,gridcomm);
584        }
585      if(sbelow){
586          MPI_Send(&e,1,mpi_element_struct,below,1,gridcomm);
587        }
588    }
589
590    if(rnorth || rsouth || rwest || reast || rabove || rbelow){
591
592      if(rnorth){
593          MPI_Recv(&e,1,mpi_element_struct,north,1,gridcomm,
              MPI_STATUS_IGNORE);
594          addElement(data,e);
595        }
596      if(rsouth){
597          MPI_Recv(&e,1,mpi_element_struct,south,1,gridcomm,
              MPI_STATUS_IGNORE);
598          addElement(data,e);
599        }
600      if(rwest){
601          MPI_Recv(&e,1,mpi_element_struct,west,1,gridcomm,
              MPI_STATUS_IGNORE);
602          addElement(data,e);
603        }
604      if(reast){
605          MPI_Recv(&e,1,mpi_element_struct,east,1,gridcomm,
              MPI_STATUS_IGNORE);
606          addElement(data,e);
607        }
608      if(rabove){
609          MPI_Recv(&e,1,mpi_element_struct,above,1,gridcomm,
              MPI_STATUS_IGNORE);
610          addElement(data,e);
611        }
612      if(rbelow){
613          MPI_Recv(&e,1,mpi_element_struct,below,1,gridcomm,
              MPI_STATUS_IGNORE);
614          addElement(data,e);
615        }
616    }
617 }
618 */
619
620 /*
621  * check if someone wants to update their working status
622  */
```

```
623  void checkOthers (){
624    int i;
625    int flag;
626    for(i = 0; i<size;i++){
627      MPI_Iprobe(i,9,gridcomm,&flag,MPI_STATUS_IGNORE);
628      if(flag){
629        MPI_Recv(&working[i],1,MPI_INT,i,9,gridcomm,MPI_STATUS_IGNORE)
               ;
630      }
631    }
632  }
633  /*
634   * notify others that my working status is changed
635   */
636  void notifyOthers(int value){
637    int i;
638    for(i = 0; i<size;i++){
639      MPI_Send(&value,1,MPI_INT,i,9,gridcomm);
640    }
641  }
642
643  /*
644   * Execute the FMM
645   */
646  void executeFMM(FmmData* data){
647    int add=0;
648    int posx,posy,posz;
649    int run = 1;
650    int sendrun = 1;
651    int senddata = 0;
652    int end = 0;
653    int sum = 0;
654    int i;
655    #ifdef DEBUG
656    printf("data_size_is_%d_%d_%d\n",data->x, data->y, data->z);
657    #endif
658    for(i = 0; i<size;i++){
659      working[i] = 1;
660    }
661    largest_solution = 0;
662    n = 1;
663    /* loop will run until all nodes are done */
664    while(!end){
665    /* working loop, will run until there are no more work to be done
           */
666    while(run){
667
668
669      if(heapGetMin(data->heap)){
670
671        Element* e,*temp;
672
673        e = heapExtractMin(data->heap);
674        #ifdef DEBUG
675        printf("%d:_setting_%d_%d_%d_to_known\n",cartrank,e->x,e->y,e->
             z);
676        #endif
677        data->bandarray[GETLINDEX(e->x,e->y,e->z)] = KNOWN;
```

```
678          int k[6] = {e->x-1,e->x,e->x+1,e->x,e->x,e->x};
679          int l[6] = {e->y, e->y-1, e->y, e->y+1,e->y,e->y};
680          int m[6] = {e->z, e->z, e->z, e->z,e->z-1,e->z+1};
681          int i;
682         posx = e->x;
683         posy = e->y;
684         posz = e->z;
685         /*
686          * update largest_solution if this solution is the largest
687          */
688         if(data->timearray[GETLINDEX(posx,posy,posz)] >
                 largest_solution){
689          largest_solution = data->timearray[GETLINDEX(posx,posy,posz)];
690         }
691         /* check if this point has been sent before */
692         if(data->sentarray[GETLINDEX(posx,posy,posz)] == 0){
693          senddata = 1;
694         }else if(data->sentarray[GETLINDEX(posx,posy,posz)] <= data->
                 timearray[GETLINDEX(posx,posy,posz)]){
695          senddata = 0;
696         }
697
698         /* printFloatArray(data->x,data->y,-2,data->timearray);
699         printFloatArray(data->x,data->y,-1,data->timearray);
700         printFloatArray(data->x,data->y,0,data->timearray);
701         printFloatArray(data->x,data->y,1,data->timearray);
702         printFloatArray(data->x,data->y,2,data->timearray);*/
703         for(i=0;i<6;i++){
704
705          if(k[i] >= 0 && k[i] < data->x && l[i] >= 0 && l[i] < data->y
                 && m[i] >= 0 && m[i] < data->z){
706           calcElement(data, k[i], l[i], m[i]);
707          }
708         }
709
710     #ifdef DEBUG
711
712         /* printFloatArray(data->x,data->y,-2,data->timearray);
713         printFloatArray(data->x,data->y,-1,data->timearray);
714         printFloatArray(data->x,data->y,0,data->timearray);
715         printFloatArray(data->x,data->y,1,data->timearray);
716         printFloatArray(data->x,data->y,2,data->timearray);*/
717     #endif
718         free(e);
719
720         n++;
721         }
722         /* Send changes to border and look for incoming changes to the
                 border */
723         sendRecvBorderChanges(data,posx,posy,posz,senddata);
724
725     #ifdef DEBUG
726         if(cartrank == 0|| cartrank == -2){printLocalArray(data->
                 timearray);
727         printLocalIntArray(data->bandarray);
728         }
729     #endif
730         senddata = 0;
```

```
731        /*Do we end the loop?*/
732        if(heapGetMin(data->heap) == 0){
733          run = 0;
734        }else{
735          run = 1;
736        }
737        checkOthers();
738        #ifdef DEBUG
739        if(cartrank == 0 && n%1000 == 0){
740          printf("%d:_reached_n_%d\n",cartrank,n);
741        }
742        #endif
743      }
744        /* see if border changes are coming */
745        sendRecvBorderChanges(data,0,0,0,0);
746        /* if we have work to do lets notify others and start to work
                again, if not let others know we are done */
747        if(heapGetMin(data->heap) != 0){
748          run = 1;
749          // notify other i am still working
750          notifyOthers(1);
751        }else if(working[cartrank] == 1){
752          // notify that i have stopped working
753          notifyOthers(0);
754        }
755        checkOthers();
756        sum = 0;
757        for(i = 0; i<size;i++){
758          sum += working[i];
759        }
760        if(sum == 0){
761          end = 1;
762        }
763      }
764      #ifdef DEBUG
765      printf("valuemax_%f_valuemin_%f\n",valuemax,valuemin);
766      #endif
767  }
768
769
770  /*
771   * free the variables used in the FMM
772   */
773  void freeFMM(FmmData* data){
774    //free(data->timearray);
775    free(data->bandarray);
776    free(data);
777  }
```

## B.3   Application

### B.3.1   mpi_app.c

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include "array_mpi.h"
4   #include "fileio.h"
5   #include <mpi.h>
6   #include <string.h>
7   #include "fmm_mpi.h"
8   #include "time.h"
9
10  #define MPIDEBUG 0
11  #define DEBUG 0
12
13  /* set that no dimensions sould be cyclic */
14  int periods[3] = {0,0,0};
15  /* set the number of dimensions to use in the cartesian node grid*/
16  int ndims = 3;
17
18  int div_x, div_y, div_z;
19
20  float* local_array;
21  float* global_array;
22  float* file_array;
23
24  /*
25   * change velocity matrix so that all values are positive and use
        sqrt to reduce the difference between the values
26   */
27  void fixMatrix(){
28    float min,max;
29    float rmin,rmax;
30    int i,j,k;
31    min = BIGFLOAT;
32    max = BIGFLOAT*-1;
33    for(i=0;i<local_x;i++)
34      for(j=0;j<local_y;j++)
35        for(k=0;k<local_z;k++){
36          if(local_array[GETLINDEX(i,j,k)] > max)
37            max = local_array[GETLINDEX(i,j,k)];
38          if(local_array[GETLINDEX(i,j,k)] < min)
39            min = local_array[GETLINDEX(i,j,k)];
40        }
41    MPI_Allreduce(&min,&rmin,1,MPI_FLOAT,MPI_MIN,gridcomm);
42    MPI_Allreduce(&max,&rmax,1,MPI_FLOAT,MPI_MAX,gridcomm);
43    min = rmin +1;
44    max = rmax;
45    for(i=0;i<local_x;i++)
46      for(j=0;j<local_y;j++)
47        for(k=0;k<local_z;k++){
48          local_array[GETLINDEX(i,j,k)] = sqrt(local_array[GETLINDEX(i
              ,j,k)] + min);
49        }
50  }
51  /*
```

```
52     * A test function to check if the array dosen't contain a value
53     */
54   int checkArray(int x, int y, int z, float value){
55       int i,j,k;
56       int rvalue = 1;
57       for(i= 0; i<x;i++)
58         for(j =0;j<y;j++)
59           for(k=0;k<z;k++){
60             if(local_array[GETLINDEX(i,j,k)] != value){
61               rvalue = 0;
62               return rvalue;
63             }
64           }
65       return rvalue;
66   }
67
68   /*
69    * Check if a global array dosen't contain a specific value
70    */
71   int checkGArray(int x, int y, int z, float value){
72       int i,j,k;
73       int rvalue = 1;
74       for(i= 0; i< x;i++)
75         for(j =0;j<y;j++)
76           for(k=0;k<z;k++){
77             if(global_array[GETINDEX(i,j,k)] != value){
78               rvalue = 0;
79               return rvalue;
80             }
81           }
82       return rvalue;
83   }
84
85   /*
86    * Divide the global matrix into smaller matrixes for each node.
87    * This function calculates the local dimensions
88    */
89   void divide_matrix(){
90
91     div_x = x/dims[0];
92     div_y = y/dims[1];
93     div_z = z/dims[2];
94     local_x = div_x;
95     local_y = div_y;
96     local_z = div_z;
97     if(local_x * dims[0] != x){
98       if(coords[0] == dims[0]){
99         local_x = (x-(local_x*dims[0])) + local_x;
100      }
101    }
102    if(local_y * dims[1] != y){
103      if(coords[1] == dims[1]){
104        local_y = (y-(local_y*dims[1])) + local_y;
105      }
106    }
107    if(local_z * dims[2] != z){
108      if(coords[2] == dims[2]){
109        local_z = (z-(local_z*dims[2])) + local_z;
```

```
110        }
111      }
112      printf("%d:_local_x_%d_local_y_%d_local_z_%d_coords_%d,%d,%d_\n",
             cartrank, local_x, local_y, local_z, coords[0], coords[1], coords
             [2]);
113 }
114
115
116 /*
117  * Initialize  the  program  allocating  local  matrixes  and  initializing
             data  types
118  */
119 void init(){
120     divide_matrix();
121     local_array = malloc(sizeof(float)*LOCALARRAYSIZE);
122     if(local_array == 0){
123        printf("Coudn't_allocate_enough_memory_for_local_array\n");
124        exit(1);
125     }
126     bzero(local_array, sizeof(float)*LOCALARRAYSIZE);
127     initMPIDatatypes();
128 }
129
130 /*
131  * a  function  for  reading  a  file  into  each  node ,  where  the  file
             contains  a  global  array ,  each  node  will  read  its  respective  part
             into  their  local  matrixes
132  */
133 void scatterdata(char* filename){
134     int i,j;
135     int* cr;
136     int value;
137     char* errorstr;
138     int reslen;
139     FILE* f;
140     int offset = 0;
141     printf("%d:_openeing_file_%s\n", cartrank, filename);
142     f = fopen(filename, "rb");
143     if(!f){
144        printf("%d:_unable_to_open_file_%d_\n", cartrank, f);
145        fflush(stdout);
146        return;
147     }
148     if(DEBUG){
149        printf("%d:_opened_file_%d\n", cartrank, f);
150     }
151     for(i = 0; i<local_x; i++){
152        for(j= 0; j< local_y; j++)
153        {
154           cr = getGlobalCord(i,j,0);
155           offset = sizeof(float)*(cr[2]+cr[1]*z+cr[0]*z*y);
156           fseek(f, offset, SEEK_SET);
157           fread(&local_array[GETLINDEX(i,j,0)], sizeof(float), local_z, f);
158           free(cr);
159        }
160     }
161     fclose(f);
162
```

```
163  }
164
165  /*
166   * gather all the local matrixes into a global matrix on node 0
167   */
168  float* gatherdata(float* iarray){
169    int i,j,k,dest,flag,r,t;
170    float* farray;
171    MPI_Request* requests;
172    MPI_Status* status;
173    requests = malloc(sizeof(MPI_Request)*local_y*local_x*2);
174    status = malloc(sizeof(MPI_Status)*local_y*local_x*2);
175    if(cartrank == 0){
176      if(DEBUG){
177        printf("%d:started gathering\n",cartrank);
178      }
179      farray = malloc(sizeof(float)*ARRAYSIZE);
180      bzero(farray,sizeof(float)*ARRAYSIZE);
181
182    }
183    if(cartrank != 0){
184    for(i=0;i<local_x;i++)
185      for(j= 0; j<local_y;j++){
186        MPI_Send(&iarray[GETLINDEX(i,j,0)],local_z,MPI_FLOAT,0,i*
                 local_y+j,gridcomm);//,&requests[(local_y*i)+j]);
187      }
188    if(DEBUG){
189      printf("%d: Done sending \n",cartrank);
190    }
191    }
192    if(cartrank==0){
193      if(MPIDEBUG){
194        printf("%d: starting setting recvs\n",cartrank);
195      }
196      for(r=0;r<x;r+=local_x)
197        for(t=0; t<y;t+=local_y)
198          for(k=0;k<dims[2];k++){
199
200            dest = getDest(r,t,k*local_z);
201            if(DEBUG){
202              printf("%d: receiving from %d\n",cartrank,dest);
203            }
204            for(i=0;i<local_x;i++)
205              for(j=0;j<local_y;j++){
206              MPI_Irecv(&farray[GETINDEX(i+r,j+t,k*local_z)],local_z,
                     MPI_FLOAT,dest,i*local_y+j,gridcomm,&requests[(
                     local_y*local_x)+(i*local_y+j)]);
207                if(dest == 0){
208                  MPI_Isend(&iarray[GETLINDEX(i,j,0)],local_z,MPI_FLOAT
                       ,0,i*local_y+j,gridcomm,&requests[(local_y*i)+j]);
209                }
210                }
211              if(dest == 0){
212                MPI_Waitall(local_y*local_x*2,requests,status);
213              }
214              if(dest != 0){
215                MPI_Waitall(local_y*local_x,&requests[local_y*local_x],
                     status);
```

```
216                 }
217             }
218       }
219
220      printf("%d:_Done_gathering_\n",cartrank);
221      return farray;
222  }
223
224  /*
225   * check if the global array is the same as the array inside a file
226   */
227  void checkData(char* filename){
228
229      if(cartrank == 0){
230      int i,j,k;
231          printf("%d:_Reading_file\n",cartrank);
232          file_array = malloc(ARRAYSIZE*sizeof(float));
233          bzero(file_array,ARRAYSIZE*sizeof(float));
234          readfile(file_array,filename,x,y,z);
235
236          printArray(file_array);
237          printf("%d:_checking_data_consistency_%f\n",cartrank,file_array[
                GETINDEX(0,0,0)]);
238          for(i=0;i<x;i++)
239            for(j=0;j<y;j++)
240              for(k=0;k<z;k++){
241                if(file_array[GETINDEX(i,j,k)] != global_array[GETINDEX(i,
                    j,k)]){
242                  printf("%d:_error_at_%d_%d_%d_file_%lf_global_%lf\n",
                      cartrank,i,j,k,file_array[GETINDEX(i,j,k)],
                      global_array[GETINDEX(i,j,k)]);
243                }
244              }
245      }
246  }
247
248  /*
249   * a test function for writing a file with values
250   */
251  void writeafile(){
252      FILE* f;
253      f = fopen("/work/idarbo/per.conv","w");
254      int i,j,k;
255      float value;
256      for(i=0;i<x;i++)
257        for(j=0;j<y;j++)
258          for(k=0;k<z;k++){
259          //value = i+j+k;
260          value = 1;
261          if(j < 8 && k < 8){
262            value = 9;
263          }
264
265          fwrite(&value,sizeof(float),1,f);
266          }
267      fclose(f);
268  }
269
```

```
270  int main(int argc, char** argv){
271     float* array ,*time1 ,*time2 ;
272     int timeusec , timesec , timeusec2 , timesec2 ,rtimesec , rtimeusec ;
273  FmmData* data ;
274           x = SIZEX ;
275           y = SIZEY ;
276           z = SIZEZ ;
277           int i , j , k ;
278           /* initialize MPI*/
279     MPI_Init(&argc , &argv ) ;
280     MPI_Comm_rank (MPI_COMM_WORLD, &rank ) ;
281     MPI_Comm_size (MPI_COMM_WORLD, &size ) ;
282     MPI_Dims_create ( size ,ndims ,dims ) ;
283
284     MPI_Cart_create (MPI_COMM_WORLD,ndims , dims , periods , 0, &gridcomm)
              ;
285
286     MPI_Cart_shift (gridcomm , 0, 1, &north , &south ) ;
287
288     MPI_Cart_shift (gridcomm , 1, 1, &west , &east ) ;
289     MPI_Cart_shift (gridcomm , 2, 1, &below , &above ) ;
290
291     MPI_Comm_rank ( gridcomm , &cartrank ) ;
292
293     MPI_Cart_coords (gridcomm , cartrank , 3, coords ) ;
294     // if ( cartrank  ==0)
295     // writeafile ();
296
297     /* init the program */
298     init () ;
299
300     if (MPIDEBUG) {
301        printf ("%d:_west_%d_east_%d_south_%d_north_%d_below_%d_above_%d_
              cartrank_%d\n" ,rank ,west ,east ,south ,north ,below ,above ,
              cartrank ) ;
302        fflush ( stdout ) ;
303     }
304
305
306
307     if ( cartrank == 0){
308        printf ("%d:_Init_complete_reading_data_from_file\n" ,cartrank ) ;
309        fflush ( stdout ) ;
310     }
311     /* read data from file */
312     scatterdata ( argv [1]) ;
313     fixMatrix () ;
314     /* instead of reading set the array to 1.0 */
315     /* for ( i=−1;i<=local_x ; i++)
316        for ( j=−1;j<=local_y ; j++)
317           for ( k=−1;k<=local_z ; k++)
318              local_array [GETLINDEX( i , j , k ) ] = 1.0;
319     */
320     // printLocalArray ( local_array ) ;
321     if ( cartrank == 0){
322        printf ("%d:_Read_data_exchanging_borders\n" ,cartrank ) ;
323        fflush ( stdout ) ;
324     }
```

```
325     /* exchange borders so border values will be correct in the
            velocity array */
326     exchangeBorders(local_array);
327     wait_exchange_float();
328     if(cartrank == 0){
329       printf("%d: Initializing FMM point is %d %d %d \n",cartrank,x/2,
            y/2,z/2);
330       fflush(stdout);
331     }
332     /* initialize the FMM */
333     data = initFMM(local_array,x/2,y/2,z/2,local_x,local_y,local_z);
334
335     if(cartrank == 0){
336       printf("%d: Executing FMM\n",cartrank);
337       fflush(stdout);
338     }
339     /* take timing and execute the FMM*/
340     MPI_Barrier(gridcomm);
341     timeusec = getTimeInMicroseconds();
342     timesec = getTimeInSeconds();
343     executeFMM(data);
344     timeusec2 = getTimeInMicroseconds();
345     timesec2 = getTimeInSeconds();
346     fixTime(timesec,timeusec,timesec2,timeusec2,&rtimesec,&rtimeusec);
347     if(cartrank == 0){
348       printf("%d: execute FMM took %d seconds and %d microseconds\n",
            cartrank,rtimesec,rtimeusec);
349     }
350
351     /*
352      * Code for checking output, in production this is the part one
            should store the arrival time array
353      */
354       if(DEBUG){
355         printLocalArray(data->timearray);
356       }
357     MPI_Barrier(gridcomm);
358
359     if(cartrank == 0){
360       printf("%d: FMM done gathering data\n",cartrank);
361     }
362     //global_array = gatherdata(data->timearray);
363
364     //global_array = gatherdata(local_array);
365
366
367
368
369     if(cartrank == 0){
370
371             //printArray(global_array);
372       if(DEBUG){
373               //printArray(global_array);
374       }
375       /*if(checkGArray(x,y,z,0.0)){
376         printf("%d: global_array is all 0.0\n",cartrank);
377       } else {
378         printf("%d: global_array is good\n",cartrank);
```

```
379        } */
380     }
381
382     printf ("%d is done ending gracefully.\n",cartrank);
383     MPI_Finalize ();
384
385   }
```

# Appendix C

# Common files

## C.1 Heap

### C.1.1 heap.h

```
1  #ifndef __HEAP_H
2  #define __HEAP_H
3  typedef
4  /*
5   * a element in the heap
6   */
7  struct e {
8          float value;
9          int x;
10         int y;
11         int z;
12  } Element;
13
14  /*
15   * heap storage struct
16   */
17  typedef struct h {
18    Element** array;
19    int maxsize;
20    int heapsize;
21  } Heap;
22  /*
23   * Initialize the heap to a specific max size.
24   * The heap should not exceed the size, it will then fail
25   */
26  Heap* initHeap(int size);
27  /*
28   * Return and remove the smallest element in the heap
29   */
30  Element* heapExtractMin(Heap* heap);
31  /*
32   * Returns the smallest element in the heap without removing it
33   */
34  Element* heapGetMin(Heap* heap);
```

```
35  /*
36   * Insert a new element to the heap
37   */
38  void heapInsert(Heap* heap,Element* key);
39  #endif
```

## C.1.2 heap.c

```c
1  #include <stdlib.h>
2  #include <stdio.h>
3  #include "heap.h"
4
5  //#define DEBUG
6  /*
7   * Initialize the heap with size as max heap size */
8  Heap* initHeap(int size){
9    Element **array;
10   Heap* heap;
11   #ifdef DEBUG
12   printf("Ready_to_allocate_heap\n");
13   #endif
14   heap = malloc(sizeof(Heap));
15   array = malloc(sizeof(Element*)*size);
16
17   #ifdef DEBUG
18   printf("Allocated_memory_for_heap\n");
19   #endif
20   heap->array = array;
21   heap->maxsize = size;
22   heap->heapsize=0;
23
24   #ifdef DEBUG
25   printf("returning_heap\n");
26   #endif
27   return heap;
28  }
29
30  /*
31   * A support function that checks that the heap is realy min sorted
32   * Only used in testing if the heap works correctly, shoudn't be
33      used in production code
34   */
34  float checkHeapConsistensy(Heap* heap){
35          int i;
36          float min = 10000000000.0;
37    if(heap->heapsize > 1){
38
39    #ifdef DEBUG
40    printf("checking_consistency,_size_is_%d\n", heap->heapsize);
41    #endif
42          for(i=1;i<=heap->heapsize;i++){
43                  if(min > heap->array[i]->value)
44                  min = heap->array[i]->value;
45          }
46    #ifdef DEBUG
47    printf("root_is_%f,_smallest_is_%f\n", heap->array[1]->value, min)
             ;
48    #endif
49    }
50    return min;
51  }
52
53  /*
54   * Returns the parten to a specific position in the heap
```

```
55    */
56  int inline getParent(int pos){
57      return pos>>1; //pos/2;
58  }
59  /*
60   * Gets the left child of a specific position in the heap
61   */
62  int inline getLeft(int pos){
63      return pos <<1; // pos * 2
64  }
65  /*
66   * Gets the right child of a specific position in the heap
67   */
68  int inline getRight(int pos){
69      return (pos <<1) +1; // pos*2+1
70  }
71
72  /*
73   * Move one value to its correct position in the min sorted heap
74   * Run once for each new element in the heap
75   */
76  void MinHeapify(Heap* heap,int pos){
77      int l,r;
78      int smallest = 0;
79      l = getLeft(pos);
80      r = getRight(pos);
81      if(l <= heap->heapsize && heap->array[l]->value < heap->array[pos
            ]->value){
82          smallest = l;
83      }else{
84          smallest = pos;
85      }
86      if(r<= heap->heapsize && heap->array[r]->value < heap->array[
            smallest]->value)
87          smallest = r;
88      if(smallest != pos){
89          Element* temp;
90          temp = heap->array[pos];
91          heap->array[pos] = heap->array[smallest];
92          heap->array[smallest] = temp;
93          MinHeapify(heap,smallest);
94      }
95  }
96
97  /*
98   * return and remove the smallest element from the heap
99   */
100 Element* heapExtractMin(Heap* heap){
101     if(heap->heapsize <1) return NULL;
102     Element* temp = heap->array[1];
103     heap->array[1] = heap->array[heap->heapsize];
104     heap->heapsize --;
105     #ifdef DEBUG
106     printf("New_heapsize_%d\n",heap->heapsize);
107     #endif
108
109     MinHeapify(heap,1);
110     return temp;
```

```
111  }
112
113  /*
114   * Return the smallest element from the heap
115   * Does not remove anything from the heap
116   */
117  Element* heapGetMin(Heap* heap){
118    if(heap->heapsize <1) return NULL;
119    Element* temp = heap->array[1];
120    return temp;
121  }
122
123  /*
124   * Add a new element to the heap
125   * this function adds a element to a position and moves it up until
          its found its place in the heap.
126   */
127  void heapIncreaseKey(Heap* heap,int pos, Element* key){
128    Element* temp;
129    //if(key->value > array[pos]->value){
130    if(heap->array[pos] != 0){
131      exit(2);
132    }
133    #ifdef DEBUG
134    printf("Inserting_into_%d\n",pos);
135    #endif
136    heap->array[pos] = key;
137    while (pos > 1 && heap->array[getParent(pos)]->value > heap->array
          [pos]->value){
138      #ifdef DEBUG
139      printf("chaging_position_between_%d_and_%d\n",pos,getParent(pos)
            );
140      #endif
141      temp = heap->array[pos];
142      heap->array[pos] = heap->array[getParent(pos)];
143      heap->array[getParent(pos)] = temp;
144      pos = getParent(pos);
145    }
146    //checkHeapConsistensy(heap);
147  }
148
149  /*
150   * add a new element to the heap
151   */
152  void heapInsert(Heap* heap,Element* key){
153    heap->heapsize++;
154    if(heap->heapsize > heap->maxsize)
155    {
156      printf("Exceeded_heap_max_size_of_%d\n",heap->maxsize);
157      exit(5);
158    }
159    #ifdef DEBUG
160    printf("New_heapsize_%d_array_at_%p\n",heap->heapsize,heap->array)
          ;
161    #endif
162    heap->array[heap->heapsize] = 0;
163    heapIncreaseKey(heap,heap->heapsize,key);
164  }
```

### C.1.3 testheap.c

```c
#include <stdio.h>
#include <stdlib.h>
#include "heap.h"

Heap* heap;

int main(int argc, char** argv){
  int i;
  heap = initHeap(500);
  Element *temp;
  Element en;
  Element to;
  Element tre;
  Element fire;
  Element fem;
  Element seks;
  Element sju;
  Element atte;
  Element ni;
  Element ti;
  en.value = 1.0;
  to.value = 2.0;
  tre.value = 3.0;
  fire.value = 4.0;
  fem.value = 5.0;
  seks.value = 6.0;
  sju.value = 7.0;
  atte.value = 8.0;
  ni.value = 9.0;
  ti.value = 10.0;

  heapInsert(heap,&en);
  printf("Inserted en\n");
  temp = heapExtractMin(heap);
  printf("extracted en with value %f\n",temp->value);

  heapInsert(heap,&fem);
  heapInsert(heap,&fire);
  heapInsert(heap,&tre);
  heapInsert(heap,&to);
  temp = heapExtractMin(heap);
  printf("extracted en with value %f\n",temp->value);

  temp = heapExtractMin(heap);
  printf("extracted en with value %f\n",temp->value);

  temp = heapExtractMin(heap);
  printf("extracted en with value %f\n",temp->value);

  temp = heapExtractMin(heap);
  printf("extracted en with value %f\n",temp->value);

  for(i=0;i<500;i++){
    heapInsert(heap,&ni);
  }
}
```

# C.2  Time

## C.2.1  time.h

```
1  int getTimeInMicroseconds(void);
2  int getTimeInSeconds(void);
3  /*
4   * Check time in sec and usec so that they are correct
5   */
6  void fixTime(int sec, int usec, int sec2, int usec2,int *rsec, int *
       rusec);
```

## C.2.2   time.c

```
1  #include <sys/time.h>
2
3  int getTimeInMicroseconds(void)
4  {
5          struct timeval tv;
6          struct timezone tz;
7          gettimeofday(&tv,&tz);
8          return tv.tv_usec;
9  }
10
11 int getTimeInSeconds(void)
12 {
13         struct timeval tv;
14         struct timezone tz;
15         gettimeofday(&tv,&tz);
16         return tv.tv_sec;
17 }
18 /*
19  * Change sec and usec so they are correct
20  */
21 void fixTime(int sec, int usec, int sec2, int usec2,int *rsec, int *
      rusec){
22   *rsec = sec2-sec;
23   if(usec2-usec < 0){
24     *rsec--;
25     *rusec = (usec2-usec) + 1000000;
26   }else{
27     *rusec = usec2-usec;
28   }
29 }
```

## C.3   Fileio

### C.3.1   fileio.h

```
1  /*
2   * stores an array in a text file
3   */
4  void printdatafile(float* array, int x, int y, int z);
5  /*
6   * Reads a binary float file into an array
7   */
8  void readfile(float* array, char* file, int x, int y, int z);
9  /*
10  * saves a float array as a png image
11  */
12 void printFloatImage(float* array, char* file, int sizex, int sizey,
       int z);
13 /*
14  * Reads float array from a text file
15  */
16 void readTextFile(float * array, char* file, int x, int y, int z);
```

## C.3.2   fileio.c

```
1   //#include <gd.h>
2   #include "/opt/freeware/include/gd.h"
3   #include <stdio.h>
4   #include <stdlib.h>
5   #include "array_mpi.h"
6   #include <string.h>
7   #define MAXFLOAT ((float)3.40282347e+38)
8   #define DEBUG
9   /*
10   * stores an array in a text file
11   */
12
13  void printdatafile(float* array, int x, int y, int z){
14    FILE* f;
15    int i, j ,k;
16    f = fopen("out.data", "w");
17    if(!f){
18      printf("Error opening file out.data ");
19      return;
20      }
21
22    for(k=0;k<z;k++){
23      for(j=0;j<y;j++){
24        for(i=0;i<x;i++){
25          fprintf(f,"%d.0 %d.0 %d.0 %lf \n",i,j,k,array[GETINDEX(i,j,k
                )]);
26        }
27        fprintf(f,"\n");
28      }
29    fprintf(f,"\n");
30    }
31    fclose(f);
32  }
33  /*
34   * Reads a binary float file into an array
35   */
36  void readfile(float * array, char* file ,int x, int y, int z){
37    FILE* f;
38    int read = 0;
39    int i,j,k;
40    f = fopen(file,"rb");
41    printf("x = %d, y=%d, z=%d\n",x,y,z);
42    fflush(stdout);
43    bzero(array,sizeof(float)*x+2*y+2*z+2);
44    printf("done zeroing array, starting read\n");
45    fflush(stdout);
46    for(i=0;i<x;i++)
47      for(j=0;j<y;j++){
48        read = fread(&array[GETINDEX(i,j,0)],sizeof(float),z,f);
49        if(read != z){
50          printf("Coudn't read hole file exiting, read %d\n",read);
51          fclose(f);
52          exit(1);
53        }
54    }
55    printf("done read, closing file\n");
```

```
56    fflush(stdout);
57    fclose(f);
58  }
59
60  /*
61   * Reads float array from a text file
62   */
63  void readTextFile(float* array, char* file, int x, int y, int z){
64    FILE *f;
65    float temp;
66    char ctemp = 'h';
67    int i,j,t;
68    f = fopen(file,"r");
69    if(!f){
70      printf("Error opening file %s",file);
71      return;
72    }
73    bzero(array,sizeof(float)*ARRAYSIZE);
74    memset(array,1,sizeof(float)*ARRAYSIZE);
75    fread(&ctemp,sizeof(char),1,f);
76    for(i=0;i<2;i++){
77      while(ctemp != '\n'){
78        printf(".");
79        fread(&ctemp,sizeof(char),1,f);
80      }
81      printf("\n");
82    }
83  // for(t=0;t<y*x;t++){
84    t=0;
85    while(fscanf(f,"%d %d %f\n",&i,&j,&temp) != EOF){
86      array[GETINDEX(i-1,j-1,z)] = temp;
87      t++;
88    }
89    printf("read %d values\n",t);
90    fclose(f);
91  }
92
93  /*
94   * saves a float array as a png image
95   */
96  void printFloatImage(float* array, char* file, int sizex, int sizey,
        int z){
97    /* Declare the image */
98    gdImagePtr im;
99    /* Declare output files */
100   FILE *pngout;
101   #ifdef DEBUG
102   printf("printing image to file \n");
103   #endif
104   int maxt, mint;
105   float max = 0,min = MAXFLOAT;
106   int i,j,r,g,b,t;
107     for(i=0;i<sizex;i++){
108         for(j=0;j<sizey;j++){
109     if(array[GETINDEX(i,j,z)] < min && array[GETINDEX(i,j,z)] !=
          BIGFLOAT){
110       min = array[GETINDEX(i,j,z)];
111       printf("new min value at %d %d %d value is %f\n",i,j,z,min);
```

```
112        } else if(array[GETINDEX(i,j,z)] > max&& array[GETINDEX(i,j,z)]
               != BIGFLOAT){
113          max = array[GETINDEX(i,j,z)];
114        }
115        }
116      }
117
118      for(i=0;i<sizex;i++){
119                for(j=0;j<sizey;j++){
120           if(array[GETINDEX(i,j,z)] == BIGFLOAT|| array[GETINDEX(i,j,z
               )] == BIGFLOAT*2){
121          array[GETINDEX(i,j,z)] = max;
122        }
123        }
124      }
125      #ifdef DEBUG
126      printf("Maxvalue_is_%f,_min_value_is_%f\n",max,min);
127      #endif
128      /* Allocate the image: 64 pixels across by 64 pixels tall */
129      im = gdImageCreate(sizex, sizey);
130
131      /* Allocate the color black (red, green and blue all minimum).
132         Since this is the first color in a new image, it will
133         be the background color. */
134      maxt = 0;
135      mint = 9999999;
136      for(i=0;i<sizex;i++){
137      for(j=0;j<sizey;j++){
138      //printf("color is %d\n",(int) (((float)array[GETINDEX(i,j,z)] / (
           float)(max−min)) *255));
139
140      t = (int) ((((float)array[GETINDEX(i,j,z)]−min) / (float)(max−min)
           ) *(255*4));
141      if(maxt < t)
142        maxt = t;
143      if(mint > t)
144        mint = t;
145      //printf("%d\n",t);
146      /* red− yellow −green − cyan − blue */
147      /* if(t<255){
148        r= 255;
149        g= t;
150        b=0;
151      } else if(t<255*2){
152        r= 255*2 −t;
153        g = 255;
154        b = 0;
155      } else if(t<255*3){
156        r = 0;
157        g = 255;
158        b = t − 255*2;
159      } else{
160        r = 0;
161        g = 255*4 − t;
162        b = 255;
163      }*//* blue − cyan − green − yeallow − red */
164      t = 255*4 −t;
165      if(t<255){
```

```
166        r= 255;
167        g= t;
168        b=0;
169     } else if (t<255*2){
170        r= 255*2 -t;
171        g = 255;
172        b = 0;
173     } else if (t<255*3){
174        r = 0;
175        g = 255;
176        b = t - 255*2;
177     } else{
178        r = 0;
179        g = 255*4 - t;
180        b = 255;
181     }
182        gdImageSetPixel(im, i, j, gdImageColorResolve(im, r/*red*/,g/*
              green*/, b/*blue*/ ));
183     }
184  }
185  #ifdef DEBUG
186  printf("Maxvalue is %d, min value is %d\n",maxt,mint);
187  #endif
188
189  /* Open a file for writing. "wb" means "write binary", important
190     under MSDOS, harmless under Unix. */
191  pngout = fopen(file, "wb");
192
193  /* Output the image to the disk file in PNG format. */
194  gdImagePng(im, pngout);
195
196  /* Close the files. */
197  fclose(pngout);
198
199  /* Destroy the image in memory. */
200  gdImageDestroy(im);
201 }
```

# C.4   Convert

## C.4.1   convert.c

```c
 1  #include <stdio.h>
 2  #include <stdlib.h>
 3
 4  int main(int argc, char** argv){
 5    FILE* filer;
 6    FILE* filew;
 7    int read, write,i;
 8    float data;
 9    filer = fopen(argv[1],"rb");
10    filew = fopen(argv[2],"wb");
11    char* cdata;
12
13    read = fread(&data,sizeof(float),1,filer);
14      if(read!= 1){
15        printf("Read_failed\n");
16      }
17
18    while(read){
19  //   for(i=0;i<1000;i++){
20        cdata = &data;
21        write = fwrite(&cdata[3],sizeof(char),1,filew);
22        if(write != 1){
23          printf("Coudn't_write\n");
24        }
25        write = fwrite(&cdata[2],sizeof(char),1,filew);
26        if(write != 1){
27          printf("Coudn't_write\n");
28        }
29        write = fwrite(&cdata[1],sizeof(char),1,filew);
30        if(write != 1){
31          printf("Coudn't_write\n");
32        }
33        write = fwrite(&cdata[0],sizeof(char),1,filew);
34        if(write != 1){
35          printf("Coudn't_write\n");
36        }
37        read = fread(&data,sizeof(float),1,filer);
38        if(read!= 1){
39          printf("Read_failed\n");
40        }
41    }
42    fclose(filer);
43    fclose(filew);
44  }
```