

Intro to distributed programming

Mohamad Ibrahim Cheikh

November 30, 2020

1 Introduction

Parallel computing refers to simulations in which multiple computational resources are implemented to solve a problem at the same time. Parallelism is usually achieved by dividing the problem into many smaller ones, in which each one is solved separately and also at the same time. Different forms of parallel programming has been developed like bit-level, instruction-level, data, and task parallelism.

The current report will analyze the performance of a distributed parallel execution model, that give the software the capability to control multiple different processes on the same or different machine which is *MPI*.

2 Experimental Configuration

The supercomputer on which we are going to conduct our performance analysis is Beocat, which is a high performance computer cluster at the Kansas State University [1]. The performance analysis will be performed on the Elves nodes of Beocat [2]. Table 1 gives a brief description of the nodes configuration.

Elve Nodes				
Nodes	1-56	57-72,77	73-76,78, 79	80-85
Processors	2x 8-Core Xeon E5-2690	2x 10-Core Xeon E5-2690 v2	2x 10-Core Xeon E5-2690 v2	2x 10-Core Xeon E5-2690v2
Ram	64GB	96GB	384GB	64GB
Hard Drive	1x 250GB 7,200 RPM SATA	1x 250GB 7,200 RPM SATA	1x 250GB 7,200 RPM SATA	1x 250GB 7,200 RPM SATA
NICs	4x Intel I350	4x Intel I350	4x Intel I350	4x Intel I350
10GbE and QDR Infini-band	MT27500 Family (ConnectX-3)	MT27500 Family (ConnectX-3)	MT27500 Family (ConnectX-3)	MT27500 Family (ConnectX-3)

Table 1: Elves node configuration taken from [2]

3 Code

This section will detail how the distributed programming code was implemented. Before discussing how the mpi-parallelized version of the code was implemented, we are going to walk over the steps that the code should implement:

- The first thing, the code will create an empty array for the keywords called `word`, and an array for the text that you want to look up for the keywords in called `line`.

- The second thing the code should do is read the **keyword.txt** file containing the list of keywords, and store them in a word, and read the **wiki_dump.txt** and store them in the array line.
- The third task is to look for the line at which the keyword in the word array occur and add it to a text.
- The final task is to write the out put on a text file.

Following the serial steps detailed above are shown in figure 1.

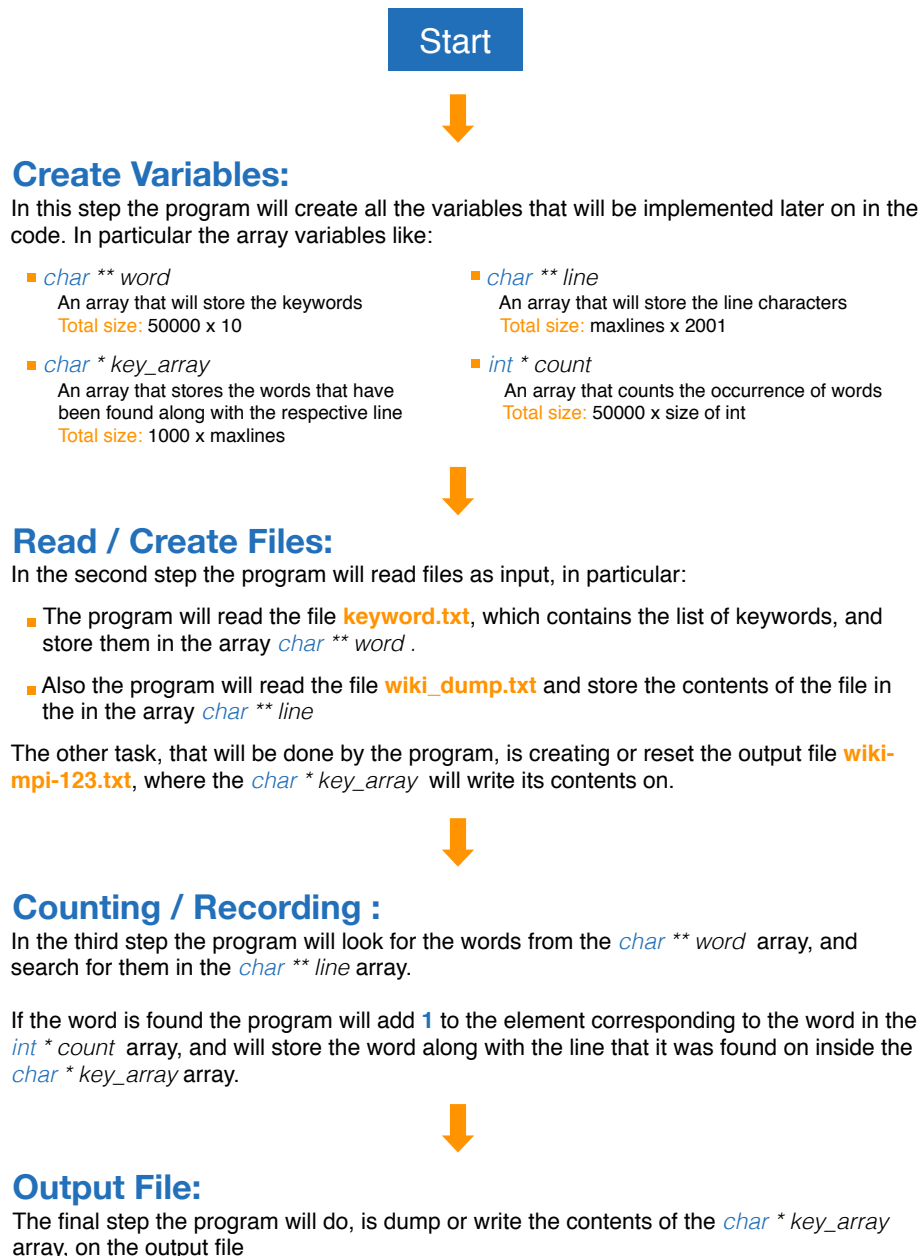


Figure 1: A list of tasks the program will do in serial

3.1 MPI Implementation

The current code was implemented in a way to minimize the amount of communication between the cores, to allow each core to work on its own without needing to depend on another core. Figure 2 shows how the code was serialized for a 4 core simulation. From the figure it is clear that all tasks, starting from the initialization, to the reading, and ending with the counting are being done in parallel manner with no core depending on the other. The only time a processor has to wait is at the end when they need to write the output file.

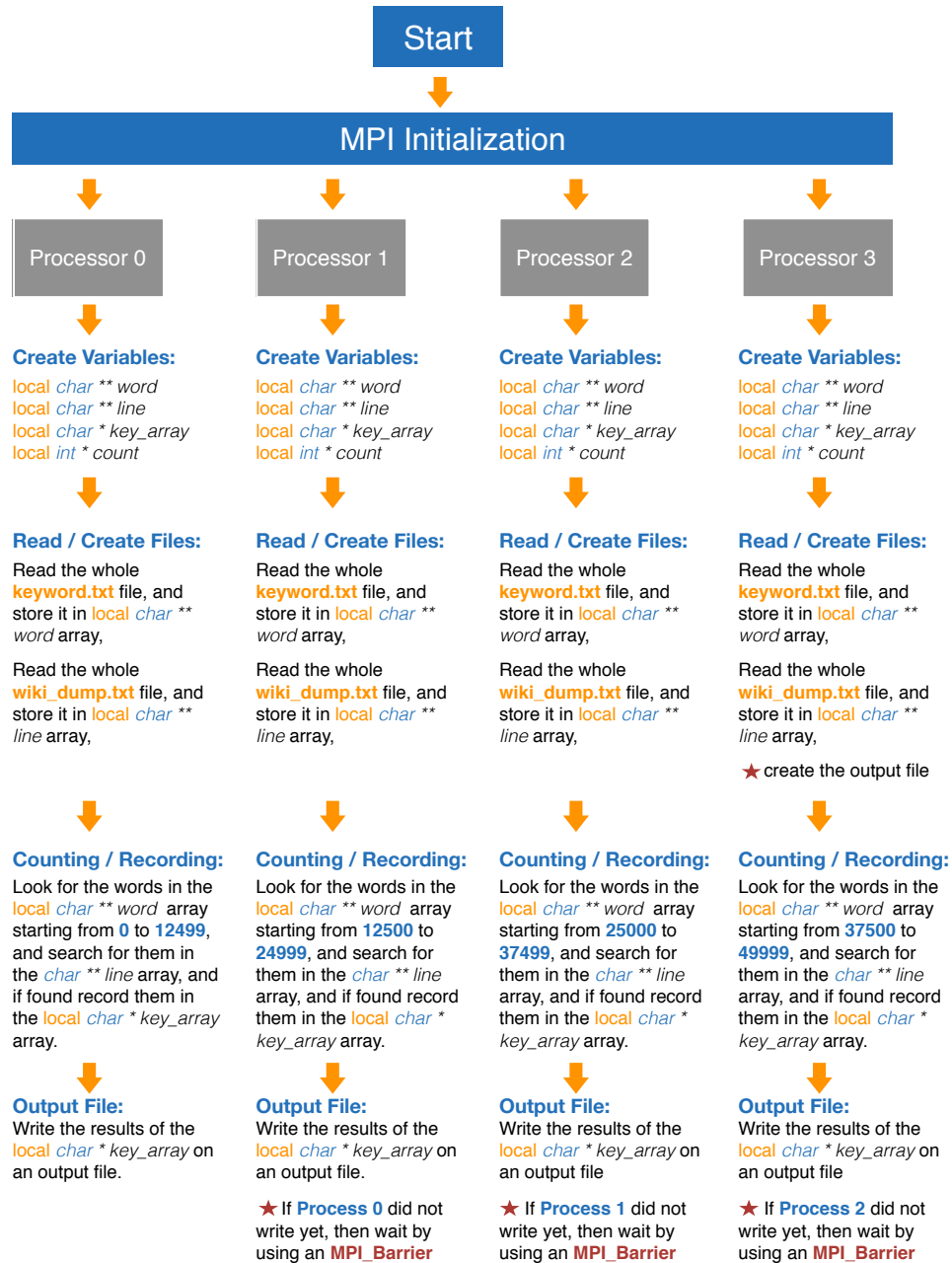


Figure 2: Parallel implementation of the serial code shown in figure 1

4 Results

For the current case the MPI implementation was divided in to two type, one type on will implement the cores on a *single machine*, and the other will divide the cores across *multiple machines*. The table below shows how many cores were run for each case.

	Number of Cores						
Type	1	2	4	8	16	32	64
<i>Single Machine</i>	X	X	X	X	X		
<i>Multiple Machines</i>			X	X	X	X	X

Table 2: Number of cores implemented for each parallel execution model

In addition to running on single or multiple machine, the MPI code was tested for different ranges of the the variable that controls the number of lines starting from 200,000 and going up to 1,000,000 with an increment of 200,000. The results of the plot are shown in the figure 3.

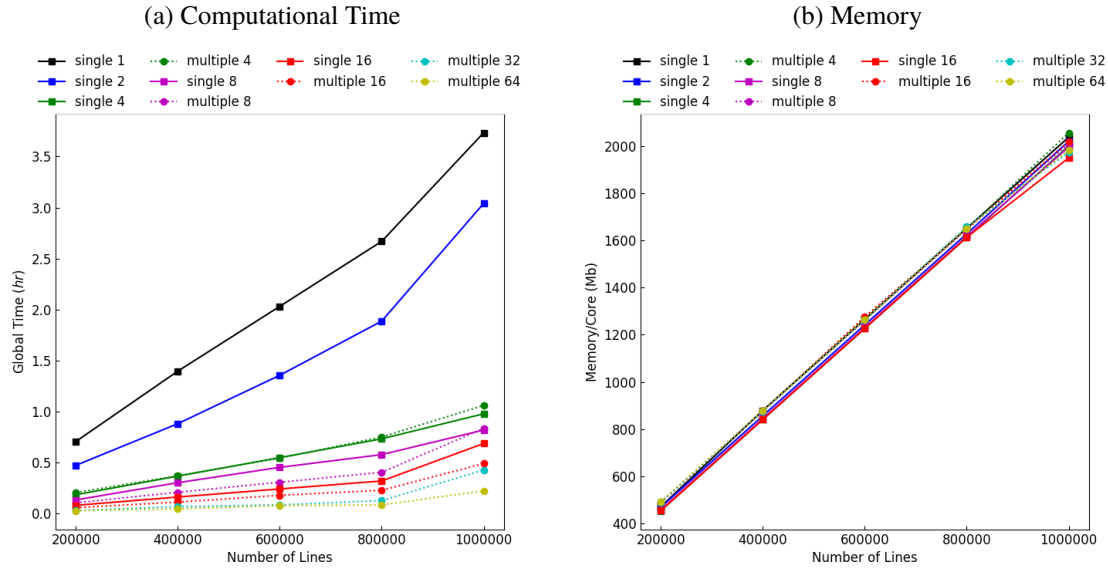
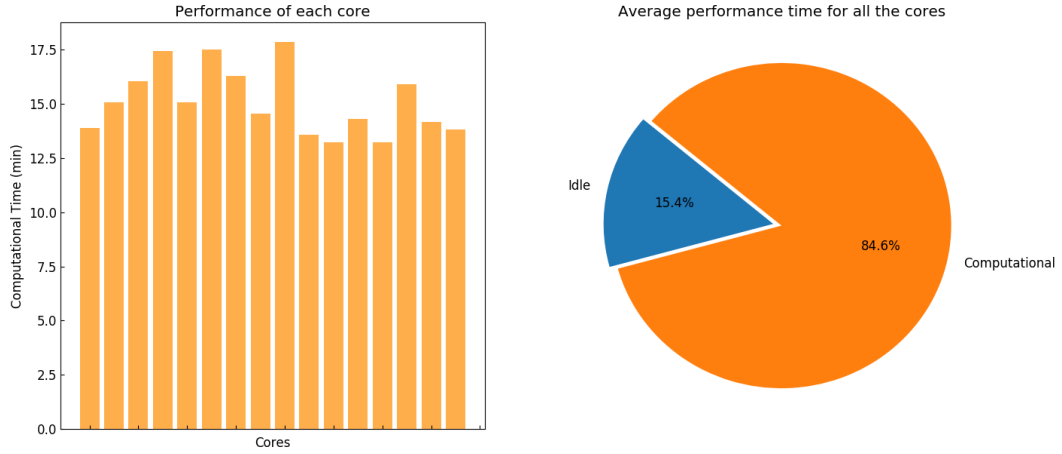


Figure 3: (a) Computational time versus number of lines, (b) Memory used per core vs the number of lines

From figure 3a it is clear that as the number of lines increases the computational time of the code increases, what is interesting is that the run on a single or on a multiple machine with the same number of cores had the same amount of computational time. The reason behind this behavior is because the MPI code is implemented in a way where cores do not rely on each other, and the communication time between the different machines was reduced to almost null, enabling both machines to behave in a similar manner. A better example to explain this behavior is shown in figure 4, where the performance of a 16 core on a single machine is displayed versus the performance of a multiple machine for a 1 million line. The results show that on average the cores on single or multiple machine operate at the same speed. However there are some cores on the multiple machine that took double the time, this can be attributed to the fact that the cache on that machine were being used by another user. But on general the mpi code on single and multiple machine operate at the same time. The downside of this approach is that all the cores will take the same amount of memory, even though they will not use all of it as is shown in figure 3b. A better way was to store only the memory of that will be needed but that will be a hassle to code

(a) *Signle Machine*



(b) *Multiple Machines*

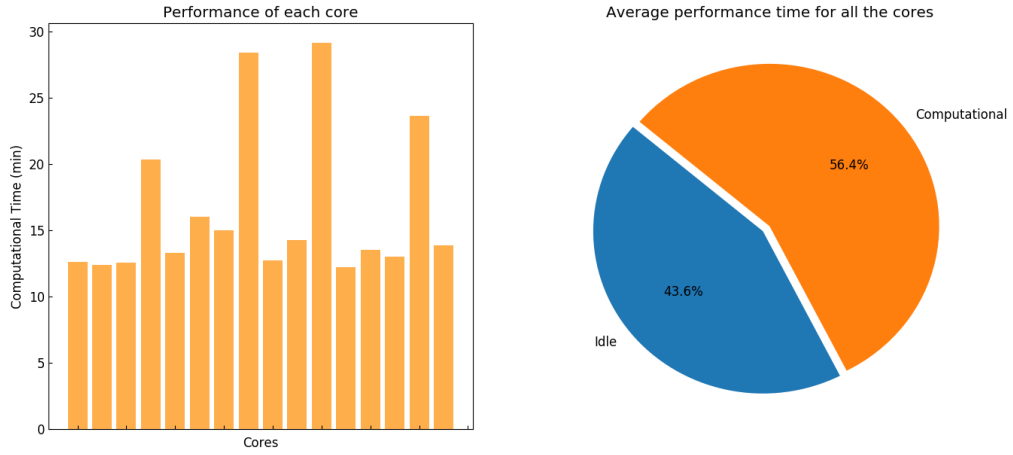


Figure 4: Computational time for each core and how much on average did the cores work or stay idel for number of lines equal 1,000,000 on a (a) single machine, (b) multiple machines

Figures 5a and 5b show the performance of the code versus the number of cores for the case with lines equal to 1,000,000. The figure clearly shows that as the number of cores increases the run will take less time starting with around 3.75 hrs for 1 core, and ending with around 8 minutes for 64 cores.

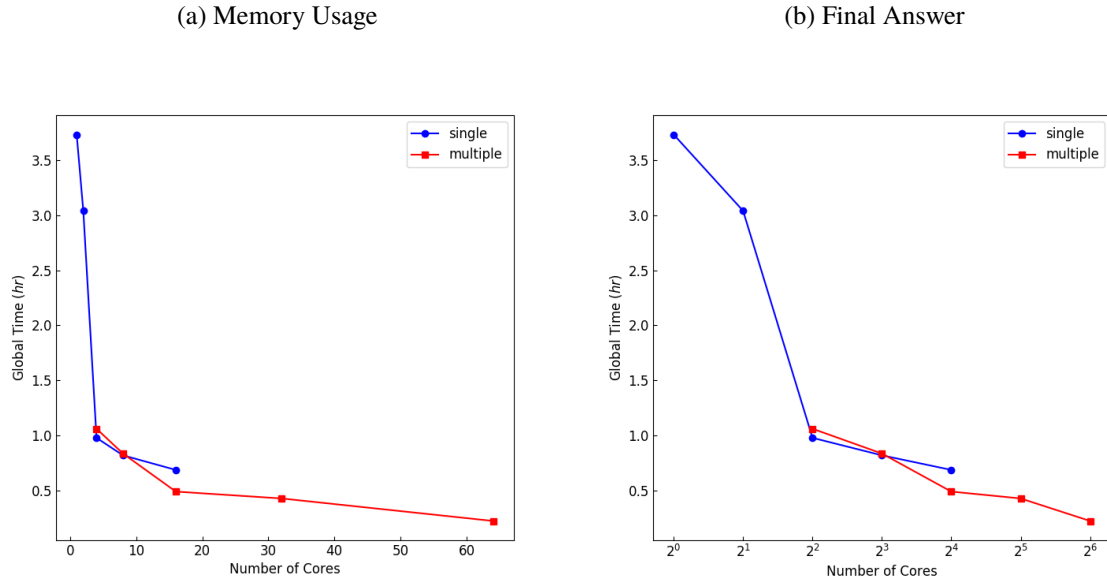


Figure 5: (a) Performance of the code versus the number of cores for the case with lines equal to 1,000,000, (b) the same but \log_2 for the x-axis

5 References:

1. https://support.beocat.ksu.edu/BeocatDocs/index.php/Main_page
2. https://support.beocat.ksu.edu/BeocatDocs/index.php/Compute_Nodes

6 Appendice:

Listing 1: 'Find_keys_mpi.c'

```
1  #include <stdlib.h>
2  #include <stdio.h>
3  #include <string.h>
4  #include <time.h>
5  #include <mpi.h>
6  #include <sys/resource.h>
7
8  # define version 1
9
10 int maxwords = 50000;
11 int maxlines;
12 int nwords;
13 int nlines;
14 int err, *count, nthreads = 1;
15 double tstart, ttotal;
16 FILE *fd;
17 char **word, **line;
18 char *key_array;
19 long key_array_size;
20 int filenumber;
21 char newword[15], hostname[256], filename[256];
22
23 /* myclock: (Calculates the time)
24 -----*/
25
26 double myclock() {
27     static time_t t_start = 0;  // Save and subtract off each time
28
29     struct timespec ts;
30     clock_gettime(CLOCK_REALTIME, &ts);
31     if( t_start == 0 ) t_start = ts.tv_sec;
32
33     return (double) (ts.tv_sec - t_start) + ts.tv_nsec * 1.0e-9;
34 }
35
36 /* init_list: (Initaite word list and lines)
37 -----*/
38 void init_list(void *rank)
39 {
40     // Malloc space for the word list and lines
41     int i;
42     int myID = *((int*) rank);
43     word = (char **) malloc( maxwords * sizeof( char * ) );
44     count = (int *) malloc( maxwords * sizeof( int ) );
45     for( i = 0; i < maxwords; i++ ) {
46         word[i] = malloc( 10 );
47         count[i] = 0;
48     }
49
50     line = (char **) malloc( maxlines * sizeof( char * ) );
51     for( i = 0; i < maxlines; i++ ) {
52         line[i] = malloc( 2001 );
```

```

53     }
54
55     key_array_size = 1000*maxlines;
56     key_array = malloc(sizeof(char)*key_array_size);
57     key_array[0] = '\0';
58 }
59
60
61 /* read_dict_words: (Read Dictionary words)
62 -----*/
63 void read_dict_words()
64 {
65     // Read in the dictionary words
66     fd = fopen("625/keywords.txt", "r" );
67     nwords = -1;
68     do {
69         err = fscanf( fd, "%[^\n]\n", word[++nwords] );
70     } while( err != EOF && nwords < maxwords );
71     fclose( fd );
72
73     //printf( "Read in %d words\n", nwords);
74 }
75
76
77 /* read_lines: (Read wiki lines)
78 -----*/
79 void read_lines()
80 {
81     // Read in the lines from the data file
82     double nchars = 0;
83     fd = fopen( "625/wiki_dump.txt", "r" );
84     nlines = -1;
85     do {
86         err = fscanf( fd, "%[^\n]\n", line[++nlines] );
87         if( line[nlines] != NULL ) nchars += (double) strlen( line[nlines] );
88     } while( err != EOF && nlines < maxlines);
89     fclose( fd );
90 }
91
92 /* reset_file: (create an output file)
93 -----*/
94 void *reset_file()
95 {
96     snprintf(filename,250,"wiki-mpi-%d.out", filename);
97     fd = fopen(filename, "w" );
98     fclose( fd );
99 }
100
101 /* count_words: (count the words)
102 -----*/
103 void *count_words(void *rank)
104 {
105     int i,k;
106     int myID = *((int*) rank);
107     int startPos = ((long) myID) * (nwords / nthreads);

```



```

108     int endPos = startPos + (nwords / nthreads);
109
110     for( i = startPos; i < endPos; i++ ) {
111         for( k = 0; k < nlines; k++ ) {
112             if( strstr( line[k], word[i] ) != NULL )
113             {
114                 if (count[i] == 0)
115                 {
116                     snprintf(newword,8,"%s:",word[i]);
117                     strcat(key_array,newword);
118                 }
119                 count[i]++;
120                 snprintf(newword,8,"%d,",k);
121                 strcat(key_array,newword);
122             }
123         }
124         //if (i % 1000 == 0) printf("Loop %d of %d done\n", i, nwords);
125         if (count[i] != 0)
126         {
127             snprintf(newword,5,"\n");
128             strcat(key_array,newword);
129         }
130     }
131 }
132
133
134 /* dump_words: (write the output on file)
135 -----*/
136 void *dump_words()
137 {
138     fd = fopen(filename, "a" );
139     int results =fputs(key_array,fd);
140     fclose( fd );
141 }
142
143
144 /*-----
145                               Main
146 -----*/
147
148 int main(int argc, char* argv[])
149 {
150
151     int i, rc;
152     int numtasks, rank;
153     double tstart_init, tend_int, tstart_count, tend_count, tend_reduce;
154     struct rusage ru;
155
156     MPI_Status Status;
157     maxlines = 100000; // Default Value
158     filenumber = rand() %100000;
159     if (argc >= 2){
160         maxlines = atol(argv[1]);
161         filenumber = atol(argv[2]); // Name of the output file
162     }

```

```

163
164 rc = MPI_Init(&argc,&argv);
165 if (rc != MPI_SUCCESS){
166     printf ("Error starting MPI program. Terminating.\n");
167     MPI_Abort(MPI_COMM_WORLD, rc);
168 }
169
170 MPI_Comm_size(MPI_COMM_WORLD,&numtasks); // Number of cores
171 MPI_Comm_rank(MPI_COMM_WORLD,&rank);     // rank of each core
172 nthreads = numtasks;
173
174 gethostname(hostname,255);
175 // The last core will reset the output file
176 if(rank == nthreads-1) reset_file();
177
178 tstart = myclock(); // Global Clock
179 // Initialization
180 tstart_init = MPI_Wtime(); // Private Clock for each core
181 init_list(&rank);
182 read_dict_words();
183 read_lines();
184 tend_int = MPI_Wtime();
185
186 getrusage(RUSAGE_SELF, &ru);
187 long MEMORY_USAGE = ru.ru_maxrss; // Memory usage in Kb
188 // Counting
189 tstart_count = MPI_Wtime();
190 count_words(&rank);
191 tend_count = MPI_Wtime();
192
193 // Outputting on File
194 tend_reduce = MPI_Wtime();
195 for(i = 0; i < nthreads;++i)
196 {
197     if(rank == i) dump_words();
198     MPI_Barrier(MPI_COMM_WORLD);
199 }
200 printf("initialization time %lf, counting time %lf, writing time %lf,
201 size %d, rank %d, hostname %s, memory %ld Kb\n", tend_int - tstart_init,
202 tend_count - tstart_count, tend_reduce - tend_count, numtasks,
203 rank, hostname, MEMORY_USAGE);
204
205 fflush(stdout);
206
207 if ( rank == 0 ) {
208     tttotal = myclock() - tstart;
209     printf("version %d, cores %d, total time %lf seconds, words %d,
210 lines %d\n", version, nthreads, tttotal, nwords, maxlines);
211 }
212 MPI_Finalize();
213 return 0;
214 }

```

Listing 2: Bash script example for multiple

```
1 #!/bin/bash
2 $$ -l mem=1G
3 $$ -l h_rt=24:00:00
4 $$ -l killable
5 $$ -cwd
6 $$ -q \*\@elves
7 $$ -pe mpi-2 64
8
9 for i in 1000000 800000 600000 400000 200000
10 do
11 mpirun -np 64 /homes/mcheikh/CIS_625/hw3/MPI_V2.out $i 264
12 hostname
13 echo -e "----Done---\n"
14 done
15 ##/homes/mcheikh/CIS_625/hw2/a.out
```

Listing 3: Bash script example for single

```
1 #!/bin/bash
2 $$ -l mem=1G
3 $$ -l h_rt=24:00:00
4 $$ -l killable
5 $$ -cwd
6 $$ -q \*\@elves
7 $$ -pe signle 16
8
9 for i in 1000000 800000 600000 400000 200000
10 do
11 mpirun -np 16 /homes/mcheikh/CIS_625/hw3/MPI_V2.out $i 168
12 hostname
13 echo -e "----Done---\n"
14 done
15 ##/homes/mcheikh/CIS_625/hw2/a.out
```
