

( HPSC **5576** ELIZABETH JESSUP )

## **HIGH PERFORMANCE SCIENTIFIC COMPUTING**

**:: Homework / 8**

**:: Student / Florian Rappi**

**1** problem / **10** points

## Problem 1

### Task:

Write a short program demonstrating the use of MPE's profiling interface, *10 pts*

Use MPE directives to instrument your implementation of *MPI\_Allreduce* from the Week 3 homework. Create custom regions for the following phases of your program, if applicable:

- Memory (malloc calls, if used)
- Communication (send/receive regions)
- Computation (loops summing vectors)

Compile and link using MPE including the library that logs the transmission of all MPI messages. (In your resulting MPE plots, you should see your colored regions, colored regions for MPI calls, and arrows for messages.)

In the week 3 homework, you were asked to test low-to-high and high-to-low bit traversals. You may just choose one here. If you run on a remote system and scp the *.clog\** file to your laptop, you will need to find the visualization toolkit of the same version to be able to open the log file for analysis.

The SLOG SDK releases are at: <ftp://ftp.mcs.anl.gov/pub/mpi/slog2/>

The version on Frost uses an intermediate build. If you run on Frost, you'll need the full MPE source containing the matching SLOG SDK: <ftp://ftp.mcs.anl.gov/pub/mpi/mpe/mpe2-1.0.4.tar.gz>

### Solution:

I did use Frost again in order to simplify the comparison process with my former implementation. I used the [L->H] (Low to High) bit traversal. To include MPE in my project I had to do the following:

<b>include flags</b>	-I /contrib/bgl/mpe2/include
<b>link flags</b>	-L /contrib/bgl/mpe2/lib -llmpe -lmpe

Therefore I compiled it using the following two commands:

```
mpixlc -g -I /contrib/bgl/mpe2/include -c -o allmpe.o all_mpe.c
mpixlc -o all_mpe allmpe.o -g -L /contrib/bgl/mpe2/lib -llmpe -lmpe
```

In order to distinguish between with and without MPE I built in another command line argument, which sets a bit flag to 1 if detected. The MPE commands will only be executed if the flag is at 0. For evaluating the *\*.clog2* files I used Jumpshot-4.

### Question:

Run your *Allreduce()* routine for two data sizes -- 8 doubles, and 1MB worth of doubles -- on 8 or 16 cores. Produce MPE plots showing the behavior of the program for both reduction operations. Can you identify the butterfly structure of the communication in the MPE plots?

Answer:

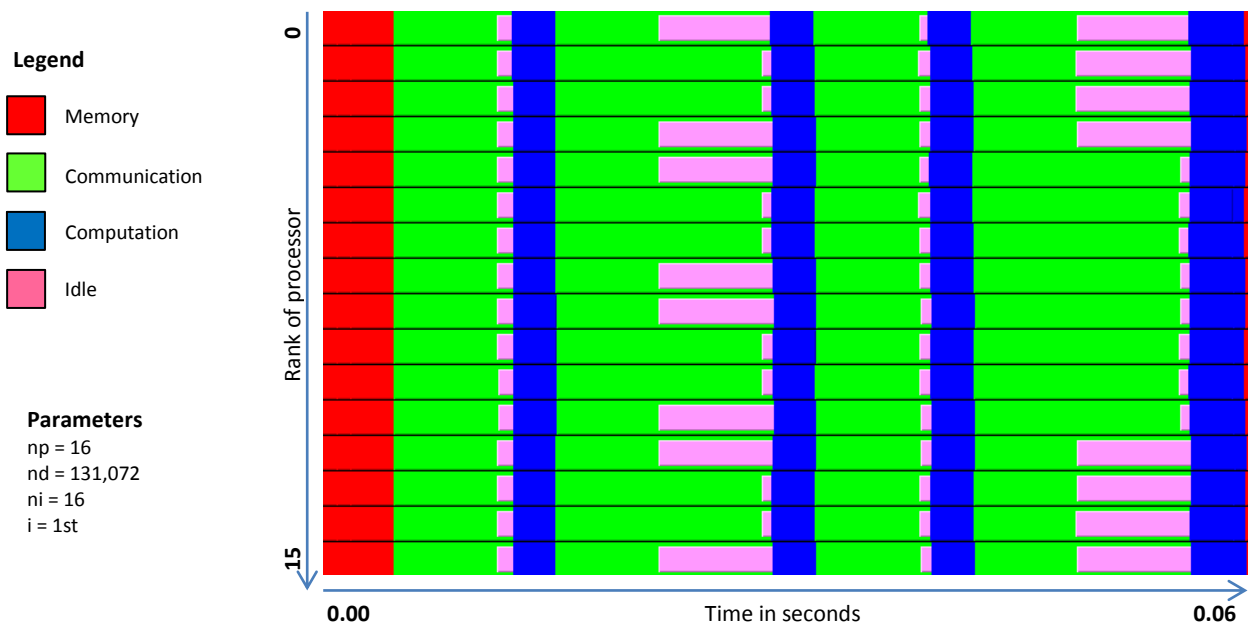


Fig. 1 Measurement with 16 processors and 131072 doubles (1 MB) running one iteration

This plot shows 16 processors performing the *Allreduce()* function (one time). The Butterfly topology is here perfectly visible (and the disadvantage of the Butterfly on the BlueGene). We see that first only next neighbor communication is performed, therefore all arriving at the same time. After that first blue line we see that communication takes longer for 1-3, 4-6 etc. whereas the communication for the others like 0-2, 5-7 etc. is as short as in the previous phase. This is obviously one of the features of BlueGene’s Torus topology, with is not optimized for Butterfly (better said Butterfly is not optimized for the Torus) and therefore produces contention on the wires in this phase. The 3<sup>rd</sup> phase offers again quite short times. Therefore we cannot do any conclusions here only that this phase is better optimized for the Butterfly topology than the previous and the next one. But in the fourth phase we see that obviously the ends are communicating, i.e. those are being slower due to contention on the wires.

*Note 1:* I used the first iteration for this plot. The other plots have a smaller memory time-block, which is nearly as small as a computation time-block.

*Note 2:* With arrow activated for the messages the butterfly structure is even easier to see. However in my opinion the white-arrows are quite hard to see on a printout. I’ve included screenshots in the \*.tar.zip file.

The next plot shows the same (first) iteration with less data, i.e. just 8 doubles. We recognize a lot of black regions, which are obviously something like idle time but only in sense of MPE not knowing what to log. I suppose those things are due to the case selections, i.e. the “if” statements enabling or disabling MPE, and the MPE calls (to get an event ID using *MPE\_Log\_get\_event\_number()*). Since the time scale on this plot is much smaller we can now see those regions whereas in the plot before we could not.

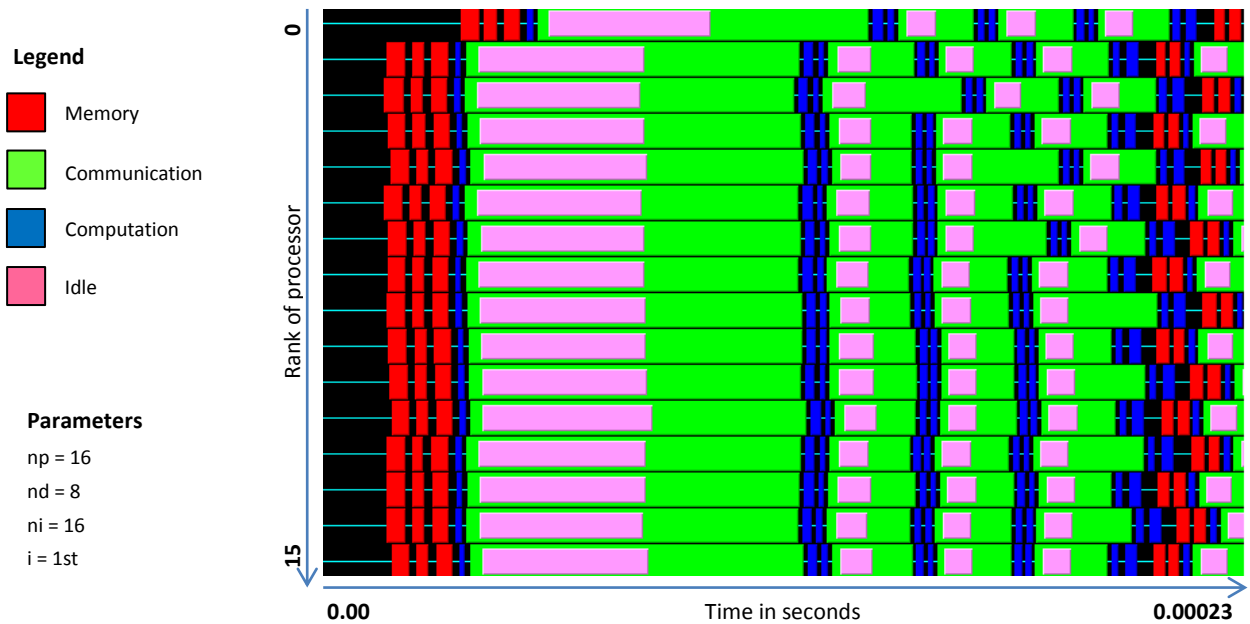


Fig. 2 Measurement with 16 processors and 8 doubles running one iteration

Another interesting thing is the shift of process number 0 (the master process). This shift is due to the *MPE\_Describe\_state()* function calls. Therefore we can clearly see what impact MPE already has on our program. We can also see that this impact is not scaling, since we did not observe such black areas in the plot before (not because they were not there but just because the time scaling was higher and therefore the region too small). Again we can see the butterfly but it is much harder to see when the message size is very small and only latency (which is kind of constant) plays a role. We only observe that all the receive partners from process 0 get stuck behind, just because that process is behind from the beginning. Therefore after the last phase half of the processes are also behind (kind of synchronized with process 0), whereas the other half is still looking quite good. We can also see that we have a lot of idle time in the very first message sending / receiving call. This alone confirms my previous homework sets in that sense that warm up runs are indeed necessary since the other phases have less idle time after the sending statement and also perform the receive much faster.

#### Question:

Run your program with and without MPE, timing the *Allreduce()* execution time. Does the inclusion of MPE profiling reduce performance?

#### Answer:

Indeed it does. I ran the program with the same setup (number of processors: 16, number of doubles: 8, number of iterations: 16) two times. The output is displayed in the next section. The total difference for this setup was  $3 \cdot 10^{-4}$  s or nearly 50% of the execution time without MPE. Overall in the program that included MPE we can say that the MPE instructions made like 30% of the whole program execution time. This is indeed a reduced performance, but since we can assume that the

MPE time does not scale linearly with the program, the bigger a program is the less performance MPE consumes in comparison to the main program.

*Program output:*

```
The final result is 960.000000
PROCESSORS:      16
VECTORLENGTH:    8
ITER 16          TIME 1.002204e-03      MPEP ON
Disabling the clock synchronization...
---
The final result is 960.000000
PROCESSORS:      16
VECTORLENGTH:    8
ITER 16          TIME 7.105557e-04      MPEP OFF
Disabling the clock synchronization...
```

The first output is with MPE active, while the second one used the “*-wompe*” command line argument of the program to disable the execution of MPE profiling.

*Code printout*

```
1 #define ITERATIONS 16
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <string.h>
5 #include "mpi.h"
6 /* global logging event variables */
7 int eventMem_s, eventMem_e; /* for memory */
8 int eventComm_s, eventComm_e; /* for communication */
9 int eventComp_s, eventComp_e; /* for computation */
10 /* prototype of (slightly modified) my_allreduce function */
11 void my_allreduce(double* sndvalue, double* recvalue, int count, unsigned
12 int rank, int processors, int tag, MPI_Comm comm, int wo_mpe);
13
14 int main(int argc, char* argv[]) {
15     int i, j; /* loop counters */
16     int my_rank; /* rank of process */
17     int p; /* number of processes */
18     int tag = 0; /* tag for messages */
19     int count = 1; /* vector size */
20     double *sndvec; /* send vector */
21     double *recvec; /* receive vector */
22     double final = 0.0; /* final result */
23     double starttime; /* timing starttime */
24     double endtime; /* timing endtime */
25     int wo_mpe = 0; /* without mpe flag */
26     /* Command line args parser (shortened for -n & -wompe)*/
27     for(i = 1; i < argc; i++)
28     {
29         // [...] same as in Homework 3 - now only additional:
30         else if(strcmp(argv[i], "-wompe") == 0)
31             wo_mpe = 1;
32     }
33     /* Start up MPI */
34     MPI_Init(&argc, &argv);
35     starttime = MPI_Wtime();
36     if(wo_mpe == 0)
37     {
```

```

38     /* Start up MPE */
39     MPE_Start_log();
40     /* Getting and setting some MPE_Log numbers */
41     eventMem_s = MPE_Log_get_event_number();
42     eventMem_e = MPE_Log_get_event_number();
43     eventComm_s = MPE_Log_get_event_number();
44     eventComm_e = MPE_Log_get_event_number();
45     eventComp_s = MPE_Log_get_event_number();
46     eventComp_e = MPE_Log_get_event_number();
47 }
48
49 /* Find out process rank */
50 MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
51 /* Find out number of processes */
52 MPI_Comm_size(MPI_COMM_WORLD, &p);
53 /* Check if processor count is %2 - else finish */
54 if(p % 2 != 0)
55 {
56     if(my_rank == 0) printf("Limited to np = 2^n.");
57 }
58 else
59 {
60     /* setup the events */
61     if(my_rank == 0 && wo_mpe == 0)
62     {
63         /* define log for M (memory) */
64         MPE_Describe_state(eventMem_s, eventMem_e,
65                             "Memory", "red");
66         /* define log for C (communication) */
67         MPE_Describe_state(eventComm_s, eventComm_e,
68                             "Communication", "green" );
69         /* define log for A (arithmetic) */
70         MPE_Describe_state(eventComp_s, eventComp_e,
71                             "Computation", "blue" );
72     }
73     if(wo_mpe == 0)
74         /* S,LOG(M) */
75         MPE_Log_event(eventMem_s, 0,
76                       "Allocating + filling sendvector");
77     /* create vector for sending data */
78     sndvec = (double*)malloc(count * sizeof(double));
79     /* filling the vector with data - my_rank */
80     for(i = 0; i < count; i++)
81         sndvec[i] = (double)my_rank;
82     if(wo_mpe == 0)
83         /* E,LOG(M) */
84         MPE_Log_event(eventMem_e, 0,
85                       "Sendvector allocated + filled");
86     /* measuring process */
87     for(i = 0; i < ITERATIONS; i++)
88     {
89         final = 0.0;
90         if(wo_mpe == 0)
91             /* S,LOG(M) */
92             MPE_Log_event(eventMem_s, 0,
93                           "Allocating + filling recvector");
94         /* create and set the receive vector */
95         recvec = (double*)malloc(count * sizeof(double));
96         for(j = 0; j < count; j++)
97             recvec[j] = 0.0;
98         if(wo_mpe == 0)

```

```

99         /* E,LOG(M) */
100         MPE_Log_event(eventMem_e, 0,
101                       "Recvvector allocated + filled");
102         /* call the specified function */
103         my_allreduce(sndvec, recvec, count, (unsigned int)my_rank,
104                    p, tag, MPI_COMM_WORLD, wo_mpe);
105         if(wo_mpe == 0)
106             /* S,LOG(A) */
107             MPE_Log_event(eventComp_s, 0,
108                           "Adding everything together");
109         /* gather the data for the final result (chk) */
110         for(j = 0; j < count; j++)
111             final += recvec[j];
112         if(wo_mpe == 0)
113             /* E,LOG(A) */
114             MPE_Log_event(eventComp_e, 0,
115                           "Finished all additions");
116         /* print out final result if all iterations done */
117         if(my_rank == 0 && (i + 1) % ITERATIONS == 0)
118             printf("The final result is %f\n", final);
119
120         if(wo_mpe == 0)
121             /* S,LOG(M) */
122             MPE_Log_event(eventMem_s, 0,
123                           "Clearing recvvector");
124         /* clear memory */
125         free(recvec);
126         if(wo_mpe == 0)
127             /* E,LOG(M) */
128             MPE_Log_event(eventMem_e, 0, "Recvvector cleared");
129     }
130
131     if(wo_mpe == 0)
132         /* S,LOG(M) */
133         MPE_Log_event(eventMem_s, 0, "Clearing sendvector");
134     free(sndvec);
135
136     if(wo_mpe == 0)
137         /* E,LOG(M) */
138         MPE_Log_event(eventMem_e, 0, "Sendvector cleared");
139     if(my_rank == 0)
140     {
141         printf("PROCESSORS:\t%d\n", p);
142         printf("VECTORLENGTH:\t%d\n", count);
143     }
144 }
145
146 if(wo_mpe == 0)
147     /* Shut down MPE */
148     MPE_Stop_log();
149 endtime = MPI_Wtime();
150
151 if(my_rank == 0)
152     printf("ITER %d\tTIME %e\tMPEP %s\n", ITERATIONS,
153           endtime - starttime, wo_mpe == 0 ? "ON" : "OFF");
154
155 /* Shut down MPI */
156 MPI_Finalize();
157
158 return 0;
159 } /* main */

```

```

160
161 void my_allreduce(double* sndvalue, double* recvalue, int count, unsigned
162 int rank, int processors, int tag, MPI_Comm comm, int wo_mpe) {
163     /* is only designed for L->H (up) in no verbose */
164     int i, j; /* Loop counters */
165     unsigned int mask = 1; /* the bit mask */
166     unsigned int dest = 0; /* destination */
167     MPI_Status status; /* status buffer */
168     double *tmpvalue; /* temporary vec */
169     if(wo_mpe == 0)
170         /* S,LOG(M) */
171         MPE_Log_event(eventMem_s, 0, "Allocating and filling Tempvec");
172     /* allocate memory for receive vector */
173     tmpvalue = (double*)malloc(count * sizeof(double));
174     /* get the receive vector set up */
175     for(j = 0; j < count; j++)
176         recvalue[j] += sndvalue[j];
177     if(wo_mpe == 0)
178         /* E,LOG(M) */
179         MPE_Log_event(eventMem_e, 0, "Tempvec allocated and filled");
180     for(i = 1; i < processors; i *= 2)
181     {
182         if(wo_mpe == 0)
183             /* S,LOG(A) */
184             MPE_Log_event(eventComp_s, 0, "Starting Bitshift");
185         /* bit shift to det. the partner */
186         dest = mask ^ rank;
187         if(wo_mpe == 0)
188         {
189             /* E,LOG(A) */
190             MPE_Log_event(eventComp_e, 0, "Bitshift ended");
191             /* S,LOG(C) */
192             MPE_Log_event(eventComm_s, 0, "Starting send/rcv");
193         }
194         /* communication */
195         MPI_Send(recvalue, count, MPI_DOUBLE_PRECISION, dest, tag,
196                 comm);
197         MPI_Recv(tmpvalue, count, MPI_DOUBLE_PRECISION, dest, tag,
198                 comm, &status);
199         if(wo_mpe == 0)
200         {
201             /* E,LOG(C) */
202             MPE_Log_event(eventComm_e, 0, "Send/rcv ended");
203             /* S,LOG(A) */
204             MPE_Log_event(eventComp_s, 0,
205                 "Adding values and bitshifting");
206         }
207         /* do the desired operation - in this case sum up */
208         for(j = 0; j < count; j++)
209             recvalue[j] += tmpvalue[j];
210         /* do the bit shift - for L->H (up) */
211         mask = mask << (unsigned int)1;
212         if(wo_mpe == 0)
213             /* E,LOG(A) */
214             MPE_Log_event(eventComp_e, 0,
215                 "Values added and bitshifted");
216     }
217 } /* my_allreduce */

```