

Parallelization with OpenMP and MPI

A Simple Example (Fortran)

Dieter an Mey, Thomas Reichstein

October 26, 2007

1 Introduction

The main aspects of parallelization using MPI (Message Passing Interface) on one hand and OpenMP directives on the other hand shall be shown by means of a toy program calculating π .

Parallelization for computer systems with distributed memory (DM) is done by explicit distribution of work and data on the processors by means of message passing.

Parallelization for computer systems with shared memory (SM) means automatic distribution of loop partitions on multiple processors, or the explicit distribution of work on the processors with compiler directives and runtime function calls (OpenMP).

MPI programs also run on shared memory systems, whereas OpenMP programs do not normally run on distributed memory machines (one exception is Intel's Cluster OpenMP)

The combination of a coarse-grained parallelization with MPI and an underlying fine-grained parallelization of the individual MPI-tasks with OpenMP is an attractive option to use a maximum number of processors efficiently. This method is known as hybrid parallelization.

2 Problem definition, serial program and automatic parallelization

π can be calculated as an integral:

$$\pi = \int_0^1 f(x) dx, \text{ with } f(x) = \frac{4}{(1+x)^2} \quad (2.1)$$

This integral can be numerically approximated through a quadrature method (rectangle method):

$$\pi = \frac{1}{n} \sum_{i=1}^n f(x_i), \text{ with } x_i = \frac{(i - \frac{1}{2})}{n} \text{ for } i = 1, \dots, n \quad (2.2)$$

The following serial program allows to vary the number of nodes n , until entering zero stops the execution.

```
1  !*****
2  ! fpi.f90 - compute pi by integrating f(x) = 4/(1 + x**2)
3  !
4  ! Variables:
5  !
6  ! pi the calculated result
7  ! n number of points of integration.
8  ! x midpoint of each rectangle's interval
9  ! f function to integrate
10 ! sum, pi area of rectangles
11 ! tmp temporary scratch space for global summation
12 ! i do loop index
13 !*****
14 !
15 program main
16 !
17 !.. Implicit Declarations ..
18 implicit none
19 !
20 !.. Local Scalars ..
21 integer :: i,n
22 double precision, parameter :: pi25dt = 3.141592653589793238462643d0
23 double precision :: a,h,pi,sum,x
24 !
25 !.. Intrinsic Functions ..
26 intrinsic ABS, DBLE
27 !
28 !.. Statement Functions ..
29 double precision :: f
30 f(a) = 4.d0 / (1.d0+a*a)
31 !
32 ! ... Executable Statements ...
33 !
34 do
35     ! Input
36     write (6,10000)
37     read (5,10001) n
38     !
39     if (n <= 0) exit ! Finish after input equals zero
40     !
41     h = 1.0d0 / n ! stride
42     ! Calculation of the quadrature formula (summation)
43     sum = 0.0d0
44     do i = 1,n
45         x = h * (DBLE(i)-0.5d0)
46         sum = sum + f(x)
47     end do
48     pi = h * sum
49     ! Output of the solution
50     write (6,10002) pi, ABS(pi-pi25dt)
51 end do
```

```

52  !
53  ! ... Format Declarations ...
54  !
55  10000 format ("Enter the number of intervals: (0 quits)")
56  10001 format (i10)

```

Entering a 10 followed by a 0, the output looks like the following :

```

Enter the number of intervals: (0 quits)
pi is approximately: 3.1415926535981615      Error is: 0.0000000000083684
Enter the number of intervals: (0 quits)

```

The approximated solution is being compared with a solution accurate to the 25th position.

In this simple example the function $f(x)$ being integrated is quite simple and parallelization only pays off when the number of nodes is quite high. Parallelization of an integral with a numerically more expensive function would be a lot more profitable.

The program core, to be parallelized, mainly consists of an inner loop, in which the sum of the values of the function $f(x)$ at the nodes is calculated.

```

1  h = 1.0d0 / n
2  sum = 0.0d0
3  do i = 1,n
4      x = h * (DBLE(i)-0.5d0)
5      sum = sum + f(x)
6  end do
7  pi = h * sum

```

The evaluation of the individual loop iterations have to be distributed to several processors.

In this simple case, the compiler is usually able to automatically parallelize this loop for a shared memory machine. The only problem arises through the recursive use of the variable **sum**, which is read and modified with each loop pass, so that every cycle depends on the previous one. Using the associativity of the summation, parallelization in this case is possible. These rounding errors usually differ from those caused by serial execution. With the the Sun Compiler you therefore have to use both the **-autopar** and the **-reduction** options.

3 Parallelization for distributed memory through Message Passing with MPI

3.1 Preliminary note

Processes with their own address space have to cooperate in order to utilize parallel machines with distributed memory. To ease the communication between separate processes, the MPI Message Passing Library was developed. The sending and receiving of messages is achieved through standardized subroutine calls, so that an MPI program is portable to all machines for which an MPI Library is available. With public domain software packages **mpich2** or **OpenMPI**, every machine that supports the tcp/ip protocol can be used.

MPI programs typically follow the SPMD programming style (Single Program Multiple Data). All involved MPI-processes, execute the same binary program, and after initialization with **MPI_Init**, every process gets the total number of parties involved with the call **MPI_Comm_size** and its own identification by calling **MPI_Comm_rank**. The task with identification zero then usually takes charge (“Master”).

3.2 MPI_Send and MPI_Recv

Here, a so called worker farm is an obvious approach to parallelize this simple example. The master process takes care of the input and output leaving just the evaluation of the inner loop to the other processes. Initially the master-process has to provide all its workers with the necessary data, in this case just the value **n** and towards the end the worker-processes have to send the partial results to the master so that the master can combine them to the overall result.

At the beginning, the master sends the value **n** with **MPI_Send** to all workers. The workers in return have to receive the data with **MPI_Recv**.

The partitioning of the loop indices to all tasks was made in cyclic fashion:

```
do i = myid + 1, n, ntasks
```

Another option would be to divide the loop iterations into chunks:

```
chunksize = ( n + ntasks - 1 ) / ntasks
do i = myid * chunksize + 1, min ( n, ( myid + 1 ) * chunksize ) . . .
```

Here, the master takes part in the computation, which is not necessarily always the case. Finally each worker sends its partial sum **mypi** to the master for collection and adding up the final result.

All tasks leave with **MPI_Finalize** the MPI environment at program end.

```
1 program main
2 ...
3 include 'mpif.h'
4 integer :: i,n,myid,ntasks,ierr,islave
5 integer, dimension(MPI_STATUS_SIZE) :: status
6 integer, parameter :: master=0, msgtag1=11, msgtag2=12
7 double precision, parameter :: pi25dt = 3.141592653589793238462643d0
8 double precision :: a,h,pi,sum,x, mypi
9 ...
10 ! Initialization of the MPI environment
11 call MPI_INIT( ierr )
12 call MPI_COMM_RANK( MPI_COMM_WORLD, myid, ierr )
13 call MPI_COMM_SIZE( MPI_COMM_WORLD, ntasks, ierr )
```

```

14 do
15   if ( myid == 0 ) then ! only the master
16     ! Input
17     write (6,10000)
18     read (5,10001) n
19     ! Distribution of the input data (here just n) to all slaves
20     do islave = 1, ntasks-1
21       call MPI_Send (n,1,MPI_INTEGER,islave, &
22         shtag1,MPI_COMM_WORLD,ierr)
23     end do
24     else ! all slaves
25       ! Receiving the input data
26       call MPI_Recv (n,1,MPI_INTEGER,master, &
27         msgtag1,MPI_COMM_WORLD,status,ierr)
28     end if
29     if (n <= 0) exit ! Finish when input is zero
30     h = 1.0d0 / n ! stride
31     ! parallel calculation of the quadrature formula
32     sum = 0.0d0
33     do i = myid+1, n, ntasks
34       x = h * (DBLE(i)-0.5d0)
35       sum = sum + f(x)
36     end do
37     mypi = h * sum
38     ! Collection of the subtotals
39     if ( myid /= 0 ) then ! Master
40       call MPI_Send (mypi,1,MPI_DOUBLE_PRECISION,master, &
41         msgtag2,MPI_COMM_WORLD,ierr)
42     else ! Slaves
43       pi = mypi
44       do islave = 1, ntasks-1
45         call MPI_Recv (mypi, 1, MPI_DOUBLE_PRECISION,islave,msgtag2, &
46           MPI_COMM_WORLD,status,ierr)
47         pi = pi + mypi
48       end do
49       ! Output of the solution
50       write (6,10002) pi, ABS(pi-pi25dt)
51     endif
52     !
53 end do
54 call MPI_FINALIZE(ierr)
55 ...
56 end program main

```

3.3 MPI_Bcast and MPI_Reduce

The frequent operations “one sends to all“ and ”all send to one“ can be implemented more elegantly through the special MPI calls **MPI_Bcast** and **MPI_Reduce** respectively. These calls are designed in such a way, that to differentiate between sender and receiver no control structures have to be programmed, just the so called root-parameter has to be set.

Attention: When using the reduction function **MPI_Reduce** in conjunction with the parameter **MPI_SUM** its not warranted that you allways receive a numerically identical result, since the MPI library takes advantage of the associativity of the summation. This can result in differend rounding errors.

```

1 program main
2 ...
3 call MPI_INIT( ierr )
4 call MPI_COMM_RANK( MPI_COMM_WORLD, myid, ierr )
5 call MPI_COMM_SIZE( MPI_COMM_WORLD, numprocs, ierr )
6 do
7 ...
8 call MPI_BCAST(n,1,MPI_INTEGER,master,MPI_COMM_WORLD,ierr)
9 if (n <= 0) exit
10 h = 1.0d0 / n
11 sum = 0.0d0
12 do i = myid+1, n, numprocs
13 x = h * (DBLE(i)-0.5d0)

```

```
14 sum = sum + f(x)
15 end do
16 mypi = h * sum
17 ! collect all the partial sums
18 call MPI_REDUCE (mypi,pi,1,MPI_DOUBLE_PRECISION,MPI_SUM,master, &
19 MPI_COMM_WORLD,ierr)
20 if (myid == 0) then
21 write (6,10002) pi, ABS(pi-pi25dt)
22 endif
23 end do
24 call MPI_FINALIZE(rc)
25 ...
26 end program main
```


4 Parallelization for Shared Memory through OpenMP Directives

4.1 Preliminary note

For shared memory programming OpenMP is the de facto standard. The OpenMP API is defined for Fortran, C and C++, it comprises of compiler directives, runtime routines, and environment variables.

At the beginning of the first parallel region of an OpenMP program (these are the program parts between **parallel** and **end parallel** directives), several lightweight processes sharing one address space, so called *Threads*, are started. These threads execute the parallel region redundantly until they reach a so called *Worksharing Construct*, in which the arising work (usually Fortran **DO** , or C/C++ **for** loops) is divided among the *Threads*.

Normally *Threads* can access all data (shared data) likewise.

Attention: In case several *Threads* modify the same shared data, access to it has to be protected in *Critical Regions* (program parts between **critical** and **end critical** directives(Fortran)). *Private Data* in which the individual *Threads* store their temporary data can be used as well. Local data of subprograms, which are called inside of parallel regions, are private too, because they are put on the stack. As a consequence, they do not maintain their contents from one call to the next!

4.2 The parallel and end parallel, the do and end do Directives

Again, the inner loop shall be parallelized, here with OpenMP directives.

```
1 h = 1.0d0 / n
2 sum = 0.0d0
3 do i = 1,n
4     x = h * (DBLE(i)-0.5d0)
5     sum = sum + f(x)
6 end do
7 pi = h * sum
```

The first step is to surround the inner loop with **parallel** and **end parallel** directives respectively. To prevent that all threads execute this loop redundantly, the loop is further enclosed by **do** and **end do** directives in order to distribute the loop iterations to all processors (worksharing).

Since by default all variables are accessible by all threads (shared), the exceptions have to be taken care of. A first candidate for privatization is the loop index **i**. If the loop iterations shall be distributed, the loop index has to be private. This is realized through a **private** clause of the **parallel** directive. As a second candidate for privatization there is the variable **x**, which is used to temporally store the node of the quadrature formula. This happens independently for each loop iteration and the variable contents is not needed after the loop.

With the usage of the summation variable **sum** it gets more complicated. On the one hand, the variable is used by all threads equally to calculate the sum of the quadrature formula, on the other hand it is set to zero prior to the loop and is needed after the loop to calculate the final solution. Would the variable be **shared**, the following problem could arise: a thread reads the value of **sum** from memory and puts it in the cache to add up his newly calculated value of **f(x)**. But before the sum can be written back to memory, another thread may read **sum** from memory to also add up a new function value. This way the contribution of the first thread may be lost.

This situation can be avoided, if only the function values are computed in parallel and stored in an auxiliary array **fx** and the summation is processed by the master thread only.

```

1  program main
2  double precision, allocatable, dimension(:) :: fx
3  ...
4  do
5      ! Input
6      write (6,10000)
7      read (5,10001) n
8      allocate (fx(n),STAT=ierror)
9      if (n <= 0) exit ! Finish after input of zero
10     h = 1.0d0 / n ! stride
11     sum = 0.0d0
12     !$omp parallel private(i,x) shared(h,fx,n)
13     !$omp do
14     do i = 1,n
15         x = h * (DBLE(i)-0.5d0)
16         fx(i) = f(x)
17     end do
18     !$omp end do
19     !$omp end parallel
20     do i = 1,n
21         sum = sum + fx(i)
22     end do
23     pi = h * sum
24     ! Output of the solution
25     write (6,10002) pi, ABS(pi-pi25dt)
26     deallocate (fx)
27 end do
28 ...
29 end program main

```

The array **fx** can confidently be declared **shared** with the corresponding clause of the parallel directive (this also is the default), since the individual threads use different loop indexes **i** and thus access different components of the array **fx**.

4.3 The critical and end critical Directives

The second solution makes use of the possibility to protect code sequences in critical regions, in which several threads modify shared variables. Critical regions are segments of code which can only be executed by a single thread at a time.

```

1  program main
2  ...
3  sum = 0.0d0
4  !$omp parallel
5  !$omp do private(i,x)
6  do i = 1,n
7      x = h * (DBLE(i)-0.5d0)
8      !$omp critical
9      sum = sum + f(x)
10     !$omp end critical
11 end do
12 !$omp end do
13 !$omp end parallel
14 pi = h * sum
15 ...
16 end program main

```

This version however involves quite some overhead, because it introduces a synchronization with every iteration of the inner loop.

The next version introduces an additional private variable, in which the individual threads sum up their contributions. The total sum is then computed in a critical region after the parallel loop.

```

1 program main
2 double precision :: a,h,pi,sum,x,sum_local
3 ...
4 h = 1.0d0 / n
5 sum = 0.0d0
6 !$omp parallel private(i,x,sum_local)
7 sum_local = 0.0d0
8 !$omp do
9 do i = 1,n
10     x = h * (DBLE(i)-0.5d0)
11     sum_local = sum_local + f(x)
12 end do
13 !$omp end do
14 !$omp critical
15 sum = sum + sum_local
16 !$omp end critical
17 !$omp end parallel
18 pi = h * sum
19 ...
20 end program main

```

This solution finally executes with a reasonable speedup.

4.4 The reduction clause

Exactly for this case there exists - analogous to the reduction function in MPI - a **reduction** clause of the **do** directive. Through its usage, the parallel program gets pleasantly short and manageable.

Attention: When the reduction clause is used with the $+$ operator it is not warranted that numerically identical solutions are generated everytime, because the associativity of the summation is utilized. Different rounding errors can occur as a result.

```

1 program main
2 ...
3 h = 1.0d0 / n
4 sum = 0.0d0
5 !$omp parallel private(i,x)
6 !$omp do reduction(+:sum)
7 do i = 1,n
8     x = h * (DBLE(i)-0.5d0)
9     sum = sum + f(x)
10 end do
11 !$omp end do
12 !$omp end parallel
13 pi = h * sum
14 ...
15 end program main

```

This version can be programmed even more concise with just one directive.

```

1 program main
2 ...
3 h = 1.0d0 / n
4 !
5 sum = 0.0d0
6 !$omp parallel do private(i,x) reduction(+:sum)
7 do i = 1,n
8     x = h * (DBLE(i)-0.5d0)
9     sum = sum + f(x)
10 end do
11 pi = h * sum
12 ...
13 end program main

```

4.5 The single and end single, and the barrier Directives

Yet, the usage of OpenMP does not limit itself just on parallelizing (inner)loops. In the following example the entire executable part of the program is enclosed in the parallel region.

```
1 program main
2 ...
3 !$omp parallel private(i,x)
4 do
5     !
6     !$omp single
7     write (6,10000)
8     read (5,10001) n
9     h = 1.0d0 / n
10    sum = 0.0d0
11    !$omp end single
12    !
13    if (n <= 0) exit
14    !$omp do reduction(+:sum)
15    do i = 1,n
16        x = h * (DBLE(i)-0.5d0)
17        sum = sum + f(x)
18    end do
19    !$omp end do
20    !$omp single
21    pi = h * sum
22    write (6,10002) pi, ABS(pi-pi25dt)
23    !$omp end single
24 end do
25 !$omp end parallel
26 ...
27 end program main
```

Therefore read and writes are enclosed in **single** and **end single** directives, causing execution just by a single thread, the first one to reach this point in the program code. The **end single** directive contains an implicit barrier, so that when the if statement is reached, all threads use the current value of the just read variable **n**. The evaluation of the stride **h** and the initialization of the summation variable **sum** is done by this simple thread too. The **barrier** which is included in the **end single** directive before the inner loop is quite important! Without the barrier, it would be possible that an “early” thread already has delivered its contribution the summation, when a “later” thread puts the first summation variable to zero. Thus, the contribution of the “early” thread would be lost.

4.6 Orphaning

The next OpenMP-program version explores the possibility of orphaning. Directives inside of a parallel region do not necessarily have to be included in the same program module. They also can reside in subprograms which are called from inside a parallel region.

```
1 program main
2 ...
3 !$omp parallel
4 do
5     !
6     !$omp single
7     write (6,10000)
8     read (5,10001) n
9     !$omp end single
10    if (n <= 0) exit
11    call calc_pi (n, pi)
12    !$omp single
13    write (6,10002) pi, ABS(pi-pi25dt)
14    !$omp end single
15 end do
16 !$omp end parallel
17 ...
18 end program Main
```

```

19  subroutine calc_pi (n, pi)
20  !
21  integer, intent(in) :: n
22  double precision, intent(out) :: pi
23  !
24  !.. Local Scalars ..
25  double precision, save :: sum, h ! sum and h are shared
26  integer :: i
27  double precision :: a,x
28  !
29  !.. Statement Functions ..
30  double precision :: f
31  f(a) = 4.d0 / (1.d0+a*a)
32  !
33  !$omp single
34  h = 1.0d0 / n
35  sum = 0.0d0
36  !$omp end single
37  !$omp do reduction(+:sum) private(i,x)
38  do i = 1,n
39      x = h * (DBLE(i)-0.5d0)
40      sum = sum + f(x)
41  end do
42  !$omp end do
43  !$omp single
44  pi = h * sum
45  !$omp end single
46  return
47  end subroutine calc_pi

```

Hence the evaluation of the inner loops including the pre- and postprocessing has been sourced out to the subprogram **calc_pi**. The main program now just consists of the input and output parts and the outer loop. Here one has to bear in mind that usually all local variables of such a subprogram are automatically private since they are allocated on the stack. Otherwise multiple threads could not concurrently pass through the same sub program, since they then would destroy each others local variables (thread safety!). In this case however the variables **h** and **sum** are supposed to be used shared! So they have to be explicitly declared static. In Fortran this can be done with **COMMON** blocks, through modules, or with the **SAVE** attribute. In C variables have to be declared **static**.

4.7 The `omp_get_thread_num` and `omp_get_num_threads` functions

The last program version suggests that one is not tied to the parallelization of loops when programming with OpenMP. On the contrary, through the usage of the function calls **omp_get_thread_num**, and **omp_get_num_threads**, which provide the thread-identification and the number of active threads resp., one can develop a program which reminds of the MPI version.

```

1  program main
2  ...
3  integer :: omp_get_thread_num, omp_get_num_threads
4  integer :: myid,numthreads
5  double precision :: sum_local
6  !$omp parallel private(i,x,sum_local,myid)
7  myid = omp_get_thread_num()
8  numthreads = omp_get_num_threads()
9  do
10     !$omp single
11     write (6,10000)
12     read (5,10001) n
13     h = 1.0d0 / n
14     sum = 0.0d0
15     !$omp end single
16     !
17     if (n <= 0) exit
18     sum_local = 0.0d0
19     !$omp barrier
20     do i = myid+1, n, numthreads
21         x = h * (DBLE(i)-0.5d0)

```

```
22     sum_local = sum_local + f(x)
23 end do
24 !$omp critical
25 sum = sum + sum_local
26 !$omp end critical
27 !$omp barrier
28 !$omp single
29 pi = h * sum
30 write (6,10002) pi, ABS(pi-pi25dt)
31 !$omp end single
32 end do
33 ...
34 end program main
```

5 Hybrid Parallelization using MPI and OpenMP

Parallelization with MPI on the top (coarse-grained) layer and with OpenMP on the bottom (fine-grained) layer can easily be combined. Thereby individual MPI-tasks are parallelized with OpenMP, or viewed from another angle, the MPI-library calls take place in the serial regions.

```
1 program main
2 ...
3 integer :: i,n,myid,numtasks,ierr,rc
4 double precision,parameter :: pi25dt=3.141592653589793238462643d0
5 double precision :: a,h,mypi,pi,sum,x
6 ...
7 call MPI_INIT( ierr )
8 call MPI_COMM_RANK( MPI_COMM_WORLD, myid, ierr )
9 call MPI_COMM_SIZE( MPI_COMM_WORLD, numtasks, ierr )
10 do
11     if ( myid .eq. 0 ) then
12         write (6,10000)
13         read (5,10001) n
14     end if
15     call MPI_BCAST(n,1,MPI_INTEGER,0,MPI_COMM_WORLD,ierr)
16     if (n <= 0) exit
17     h = 1.0d0 / n
18     sum = 0.0d0
19     !$omp parallel do reduction(+:sum) private(i,x)
20     do i = myid+1, n, numtasks
21         x = h * (DBLE(i)-0.5d0)
22         sum = sum + f(x)
23     end do
24     mypi = h * sum
25     call MPI_REDUCE(mypi,pi,1,MPI_DOUBLE_PRECISION,MPI_SUM,0, &
26     MPI_COMM_WORLD,ierr)
27     if (myid .eq. 0) then
28         write (6,10002) pi, ABS(pi-pi25dt)
29     endif
30     !
31 end do
32 call MPI_FINALIZE(ierr)
33 ...
34 end program main
```

In this simple example a single OpenMP directive has to be introduced into the MPI version, to demonstrate a valid hybrid program.

