

# MPI Lab

from I. Moulitsas

October 2018

## 1 Working with the Crescent Cluster

In order to use MPI and run on an HPC cluster, you will be working with Crescent - Cranfield University-wide small-scale research HPC facility. Crescent has a total of 77 SGI Chassis compute nodes, which are a mixture of other node types. The largest number of nodes are SGI Chassis'. They have two Intel E5-2660 (Sandy Bridge) CPUs giving 16 CPU cores and 64GB of shared memory, with 8 of these nodes having 128GB of memory. Taken together these give a total of 1488 available cores. Crescent has a single GPU node. This GPU node has two Intel E5-2698 v3 (Haswell) CPUs giving 32 CPU cores, 256GB of shared memory, and 4 Tesla K80 GPU cards. You will have access to Crescent for the duration of your MSc CSTE studies.

When developing/testing MPI programs some users tend to run small MPI programs interactively on the login node. This is not good practice and should be avoided. In order to run computations which measure execution time, it is only advisable that you use *PBS* - a batch system which manages jobs on the Crescent and ensures that each of your MPI processes gets allocated to a separate CPU.

### 1.1 Accessing Crescent

A substantial amount of information regarding research computing at Cranfield University can be found at <https://intranet.cranfield.ac.uk/it/Pages/HighPerformance.aspx>

In order to login to Crescent from the command line:

```
ssh -X <s123456>@crescent.central.cranfield.ac.uk
```

Where <s123456> is your Cranfield University username. If you are logging in from off-campus you will first need to **ssh** into **hpcgate.cranfield.ac.uk**.

If you are accessing Crescent from a Windows machine you can use putty, or some other similar program.

### 1.2 Some useful commands

#### Transferring files

To copy a file, for example **test.c**, from directory **X** on Crescent to current directory:

```
scp <s123456>@crescent.central.cranfield.ac.uk:CRESCENT-DRIVE/X/test.c ./
```

To copy all files from directory **X** on Crescent to current directory:

```
scp <s123456>@crescent.central.cranfield.ac.uk:CRESCENT-DRIVE/X/* ./
```

To copy a file, for example **test.c**, from current directory to directory **X** on Crescent:

```
scp ./test.f <s123456>@crescent.central.cranfield.ac.uk:CRESCENT-DRIVE/X/
```

To copy all files from current directory to directory **X** on Crescent:

```
scp ./* <s123456>@crescent.central.cranfield.ac.uk:CRESCENT-DRIVE/X/
```

If you are accessing Crescent from a Windows machine you can use winscp3, or some other similar program.

## Standard Commands

Some Unix commands you might find useful:

- `ls` - list the files in the current directory
- `ls -laF` - list the files in the current directory including details such as modification date and file size
- `cd A` - change directory to A
- `cd ..` - change directory - go one level up
- `mkdir` - create a directory
- `cp <pathA>/file.A <pathB>` - copy file.A from directory given by pathA to directory given by pathB
- `rm file.A` - remove file file.A
- `rmdir A` - remove directory A
- `cat file.A` - print the contents of file.A on the screen
- `less file.A` - open file.A in a viewer (q to exit the viewer)
- `history` - print last commands executed
- `pwd` - print working directory

In general you have access to further information and the arguments of a Unix command by:

- `man <command_name>`

## Editing files on Crescent

You can either run the editor on Crescent remotely, i.e. after logging in via ssh:

```
nedit &
```

Or edit files locally and copy to Crescent using `scp`.

There are more advanced editors for Unix based environments, such as `vim` or `emacs`. For these you will need to get acquainted to the command interface of such editors, it is recommended you [google](#) for a cheatsheet.

## 2 MPI

### 2.1 Compiling MPI programs

In order to use MPI one first needs to set their environment variables accordingly.

To use an MPI implementation use one of the following modules

```
module load OpenMPI
module load MPICH
module load intel
```

In order to compile an MPI C or FORTRAN program, run:

```
mpicc <options> <file name>
```

or

```
mpif90 <options> <file name>
```

mpicc and mpif90 are in fact wrappers for the C or Fortran compiler, so all those options are applicable. For more details

```
man <compiler>
```

### 2.2 Running MPI programs

To run an mpi program, say a.out, the command is:

```
mpirun -np <number of processes to start> ./a.out
```

As discussed, the scheduling system should be used in order to run a program. To submit an MPI program to the scheduler you will need to use:

```
qsub <script>
```

command where <script> is the job submission script. A sample script can also be found at /apps/examples/pbs/mpi.sub. The script you submit should correspond to the following template:

```
#!/bin/bash
##
## MPI submission script for PBS on CRESCENT
## -----
##
## Follow the 6 steps below to configure your job
##
## STEP 1:
##
## Enter a job name after the -N on the line below:
##
#PBS -N mpi_example
##
## STEP 2:
##
## Select the number of cpus/cores required by modifying the #PBS -l select lin\
e below
##
```

```

## Normally you select cpus in chunks of 16 cpus
## The Maximum value for ncpus is 16 and mpirprocs MUST be the same value as ncp\
us.
##
## If more than 16 cpus are required then select multiple chunks of 16
## e.g. 16 CPUs: select=1:ncpus=16:mpiprocs=16
##      32 CPUs: select=2:ncpus=16:mpiprocs=16
##      48 CPUs: select=3:ncpus=16:mpiprocs=16
##      ..etc..
##
#PBS -l select=1:ncpus=16:mpiprocs=16
##
## STEP 3:
##
## Select the correct queue by modifying the #PBS -q line below
##
## half_hour      - 30 minutes
## one_hour       - 1 hour
## half_day       - 12 hours
## one_day        - 24 hours
## two_day        - 48 hours
## five_day       - 120 hours
## ten_day        - 240 hours (by special arrangement)
##
#PBS -q half_hour
##
## STEP 4:
##
## Replace the hpc@cranfield.ac.uk email address
## with your Cranfield email address on the #PBS -M line below:
## Your email address is NOT your username
##
#PBS -m abe
#PBS -M hpc@cranfield.ac.uk
##
## =====
## DO NOT CHANGE THE LINES BETWEEN HERE
## =====
#PBS -j oe
#PBS -W sandbox=PRIVATE
#PBS -k n
ln -s $PWD $PBS_O_WORKDIR/$PBS_JOBID
## Change to working directory
cd $PBS_O_WORKDIR
## Calculate number of CPUs
cpus='cat $PBS_NODEFILE | wc -l'
## set some MPI tuning parameters to use the correct transport
export I_MPI_FABRICS=shm:dapl
export I_MPI_DAPL_UD=enable

```

```

export I_MPI_PLATFORM=bdw
export I_MPI_ADJUST_ALLREDUCE=5
## =====
## AND HERE
## =====
##
## STEP 5:
##
## Load the default application environment
## For a specific version add the version number, e.g.
## module load intel/2016b
##
module load intel
##
## STEP 6:
##
## Run MPI code
##
## The main parameter to modify is your mpi program name
## - change YOUR_EXECUTABLE to your own filename
##

mpirun -machinefile $PBS_NODEFILE -np ${cpus} YOUR_EXECUTABLE

## Tidy up the log directory
## DO NOT CHANGE THE LINE BELOW
## =====
rm $PBS_O_WORKDIR/$PBS_JOBID
#

```

Where you specify your values for the job name, output/error file names, number of processors and executable name. Once the job has been submitted you can use `qstat` command to show the status of the queue and your job's position. As noted, a sample script can also be found at `/apps/examples/pbs/mpi.sub`

### 3 Tasks

In these lab sessions you will write four MPI programs:

- Simple Hello World Send-Receive test
- Parallel Dot Product computation
- Parallel Pi computation
- Parallel Ping-Pong network timing

You may find the following online MPI references useful:

```

http://www.netlib.org/utk/papers/mpi-book/mpi-book.html
http://www.llnl.gov/computing/tutorials/mpi/

```

### 3.1 Hello World

- Write an MPI program which:
  - Initialises MPI
  - Obtains the rank of the process
  - For each process prints “Hello World” message and the rank of the process
  - Do not forget to check the return value of every MPI operation to detect errors
- Run this program interactively on 2 processors.
- Run the program through the scheduler on 16 processors

### 3.2 Send-Receive Test

- Modify your Hello World program so that:
  - Process 0 Sends each process a message with a real array with n variables, where n is the user input
  - Each of the other processes receives the array and prints “Hello World” message, rank of the process and the first element of the array received
  - Process 0 measures the **overall send** time using MPI\_WTIME function
  - Do not forget to check the return value of every MPI operation to detect errors
- Run this program interactively on 2 processors.
- Run the program through the scheduler on 16 processors

### 3.3 Dot Product Computation

Based on the following serial implementation of the dot product computation, write an MPI program to implement the same tasks. Vectors will be distributed among the processes and use MPI\_Allreduce when necessary.

```
c serial_dot.f -- compute a dot product
c                      serially -- on a single processor.
c
c Input:
c   n: vector size
c   x, y: vectors
c
c Output:
c   the dot product of x and y.
c
c   PROGRAM serial_dot
c   implicit none
c   integer MAX_ORDER, i, n
c   parameter (MAX_ORDER = 20000)
c   real dot, x(MAX_ORDER), y(MAX_ORDER)
c   double precision t1, t2
c
c   real serial_dot_product
```

```

c
    call timing(t1)

    print *, 'Enter the vector size'
    read *, n
    call read_vector('the first vector', x, n)
    call read_vector('the second vector', y, n)

    do i=1, 1000000
        dot = serial_dot_product(x, y, n)
    enddo
    write (*,*) 'The dot product is ', dot

    call timing(t2)
    write (*,*) 'user+system time=',t2-t1

    stop
end

c
c *****
c      subroutine read_vector(prompt, v, n)
c      implicit none
c      integer i, n
c      character *20  prompt
c      real  v(n)

c
c      print *, 'Enter ', prompt, 'data (return after each entry): '
c      do i = 1, n
c          read (*,*) v(i)
c      enddo

c      return
c      end

c
c *****
c      real function serial_dot_product(x, y, n)
c      implicit none
c      integer i, n
c      real  x(n), y(n)
c      real  sum

c
c      sum = 0.0
c      do i = 1, n
c          sum = sum + x(i)*y(i)
c      enddo
c      serial_dot_product = sum

c      return
c      end

c
c *****
c      subroutine timing (tnow)
c      double precision tnow
c      real*4 d(2)

```

```

    rtnow=etime(d)
    tnow=d(1)+d(2)

c    write (*,*) 'rtnow=', rtnow

    return
end

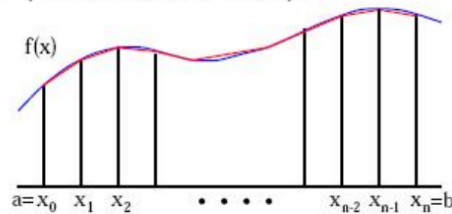
```

### 3.4 Computation of PI

- Well-known formula:

$$\int_0^1 \frac{4}{1+x^2} dx = \pi.$$

- Numerical integration (Trapezoidal rule):



$$\int_a^b f(x) dx \approx h \left[ \frac{1}{2} f(x_0) + f(x_1) + \cdots + f(x_{n-1}) + \frac{1}{2} f(x_n) \right].$$

$$x_i = a + ih, \quad h = (b - a)/n, \quad n = \# \text{ of subintervals.}$$

Figure 1: Computation of  $\pi$

Parallelisation can be achieved based on the split of the integration interval into sub-intervals with each processor computing its part of the integral sum and returning it to the master processor which gathers the result. The sequential code is given by (or you can use the code which you have written during the introduction week) :

```

program compute_pi
    integer n, i
    double precision w, x, sum, pi, f, a
! function
    f(a) = 4.d0 / (1.d0 + a*a)
    print *, ' Enter number of intervals: '
    read *, n
! interval size
    w = 1.0d0/n
    sum = 0.0d0
    do i = 1, n
        x = w * (i - 0.5d0)

```



```

        sum = sum + f(x)
    end do
    pi = w * sum
    print *, 'computed pi = ', pi
    stop
end program compute_pi

```

- Time the execution of the sequential code using `CPU_TIME` FORTRAN subroutine (see Intel's FORTRAN reference on the Blackboard).
- Parallelise the algorithm based on the split of the integration interval into  $n$  subintervals where  $n$  is the number of processes
- Write the following two parallel versions of this program:
  - Version based on the reduction operation `MPI_REDUCE`.
  - Version based on point-to-point communications `MPI_SEND`, `MPI_RECV`.
  - Each version should calculate the execution using `MPI_WTIME` function.
- Run these programs interactively on 2 processors.
- Run these programs through the scheduler on 16 processors

### 3.5 Ping-Pong

- Generate a ping-pong code which exchanges 1000 messages of double arrays of size  $n$ , where  $n$  is user input, between two processes and measures the total communication time.
- Run the program locally on the Crescent login node
- Run the program through the scheduler
- Run the program using explicit hosts specification:
 

```
mpirun -np <number of processors> <executable>
```
- Based on the comparison of the communication time difference between different nodes and on the same node comment on the network speed versus the memory speed.