Advanced MPI Erwin Laure Director PDC-HPC, Guest Professor KTH



1









7

Profiling Interface

- MPI allows to log certain events to a log file that can be analyzed post-mortem
- Part of the MPI MultiProcessing Environment
 - Prefix MPE
 - Tracing Library This traces all MPI calls. Each MPI call is preceded by a line that contains the rank in MPI_COMM_WORLD of the calling process, and followed by another line indicating that the call has completed..
 - Animation Library This is a simple form of real-time program animation and requires X window routines.
 - Logging Library This is the most useful and widely used profiling libraries in MPE. They form the basis to generate log files from user MPI programs. There are currently 3 different log file formats allowed in MPE.



















Main Cartesian Commands

- MPI_CART_CREATE: creates a new communicator using a Cartesian topology
- MPI_CART_COORDS: returns the corresponding Cartesian coordinates of a (linear) rank in a Cartesian communicator.
- MPI_CART_RANK: returns the corresponding process rank of the Cartesian coordinates of a Cartesian communicator.
- MPI_CART_SUB: creates new communicators for subgrids of up to (N-1) dimensions from an N-dimensional Cartesian grid.
- MPI_CART_SHIFT: finds the resulting source and destination ranks, given a shift direction and amount. 17

Example

```
#include "mpi.h"
MPI_Comm old_comm, new_comm;
int ndims, reorder, periods[2], dim_size[2];
old_comm = MPI_COMM_WORLD;
ndims = 2;
                  /* 2-D matrix/grid */
                   /* rows */
dim_size[0] = 3;
dim_size[1] = 2;
                   /* columns */
periods[0] = 1;
                   /* row periodic (each column forms a
                       ring) */
periods[1] = 0;
                    /* columns nonperiodic */
reorder = 1;
                    /* allows processes reordered for
                       efficiency */
MPI_Cart_create(old_comm, ndims, dim_size,
               periods, reorder, &new_comm);
                                                        19
```

Example Cont'd			
	-1,0 (4)	-1,1 (5)	
0,-1(-1)	0,0 (0)	0,1 (1)	0,2(-1)
1,-1(-1)	1,0 (2)	1,1 (3)	1,2 (-1)
2,-1(-1)	2,0 (4)	2,1 (5)	2,2 (-1)
	3,0 (0)	3,1 (1)	
periods(0)=.true.;periods(1)=.false.			

















MPI File Structure

- MPI defines how multiple processes access and modify data in a shared file.
- Necessary to think about how data is partitioned within this file
 - Similar to how derived datatypes define data partitions within memory
- MPI-IO works with simple datatypes and derived datatypes
 - Derived datatypes are preferred because of performance benefits
- A view defines the current set of data, visible and accessible, from an open file.
 - Each process has its own view of the shared file that defines what data it can access.
 - A view can be changed by the user during program execution.

29



















39

MPI_Win_create exposes local memory to RMA operation by other processes in a communicator Collective operation Creates window object MPI_Win_free deallocates window object MPI_Put moves data from local memory to remote memory MPI_Get retrieves data from remote memory into local memory MPI_Accumulate updates remote memory using local values Data movement operations are non-blocking Subsequent synchronization on window object needed to ensure operation is complete







43

And finally ...

• The top MPI Errors according to

Advanced MPI: I/O and One-Sided Communication, presented at SC2005, by William Gropp, Rusty Lusk, Rob Ross, and Rajeev Thakur

http://www.mcs.anl.gov/research/projects/mpi/tutorial/ advmpi/sc2005-advmpi.pdf)



