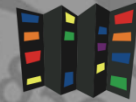
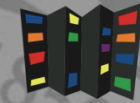


Introduction to Parallel Programming

Martin Čuma
Center for High Performance Computing
University of Utah
m.cuma@utah.edu



- Types of parallel computers.
- Parallel programming options.
- How to write parallel applications.
- How to compile.
- How to debug/profile.
- Summary, future expansion.



Single processor:

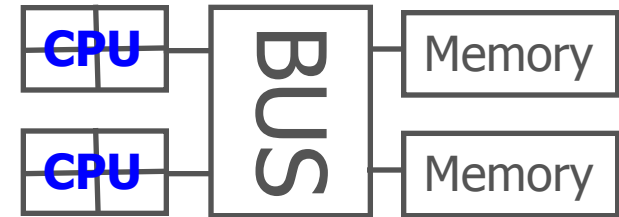
- SISD – single instruction single data.

Multiple processors:

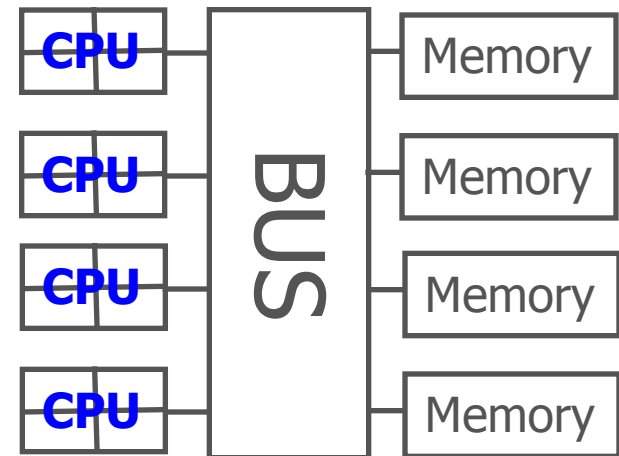
- SIMD - single instruction multiple data.
- MIMD – multiple instruction multiple data.
 - Shared Memory
 - Distributed Memory

- All processors have access to local memory
- Simpler programming
- Concurrent memory access
- More specialized hardware
- CHPC :
Linux clusters, 2, 4, 8, 12, 16 core nodes
GPU nodes

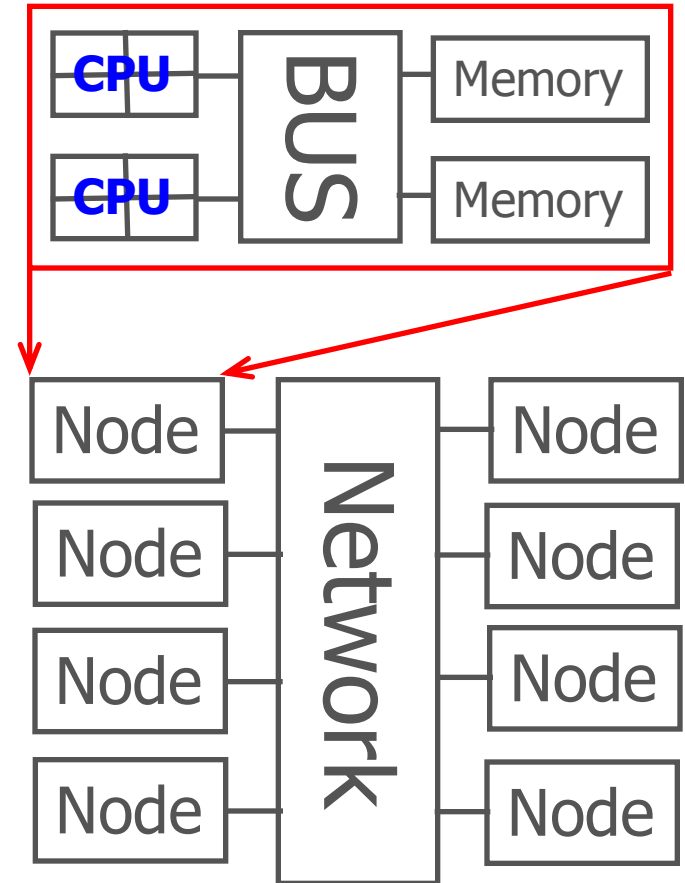
Dual quad-core node



Many-core node (e.g. SGI)



- Process has access only to its local memory
- Data between processes must be communicated
- More complex programming
- Cheap commodity hardware
- CHPC: Linux clusters



8 node cluster (64 cores)

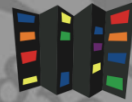


Shared Memory

- Threads – POSIX Pthreads, **OpenMP** (CPU, MIC), OpenACC, CUDA (GPU)
 - Thread – own execution sequence but shares memory space with the original process
- Message passing – processes
 - Process – entity that executes a program – has its own memory space, execution sequence

Distributed Memory

- Message passing libraries
 - Vendor specific – non portable
 - General – **MPI**, PVM



- Compiler directives to parallelize
 - Fortran – source code comments

```
!$omp parallel/!$omp end parallel
```
 - C/C++ - #pragmas

```
#pragma omp parallel
```
- Small set of subroutines
- Degree of parallelism specification
 - OMP_NUM_THREADS or
omp_set_num_threads(INTEGER n)

- Communication library
- Language bindings:
 - C/C++ - `int MPI_Init(int argv, char* argc[])`
 - Fortran - `MPI_Init(INTEGER ierr)`
- Quite complex (100+ subroutines)
but only small number used frequently
- User defined parallel distribution

- Complex to code
- Slow data communication

- Ported to many architectures
- Many tune-up options for parallel execution

- Easy to code
- Fast data exchange

- Memory access (thread safety)
- Limited usability
- Limited user's influence on parallel execution



- saxpy – vector addition: $\bar{z} = a\bar{x} + \bar{y}$
- simple loop, no cross-dependence, easy to parallelize

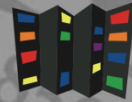
```
subroutine saxpy_serial(z, a, x, y, n)
integer i, n
real z(n), a, x(n), y(n)

do i=1, n
    z(i) = a*x(i) + y(i)
enddo
return
```

```
subroutine saxpy_parallel_omp(z, a, x, y, n)
integer i, n
real z(n), a, x(n), y(n)

!$omp parallel do
do i=1, n
    z(i) = a*x(i) + y(i)
enddo
return
```

```
setenv OMP_NUM_THREADS 4
```



```
subroutine saxpy_parallel_mpi(z, a, x, y, n)
integer i, n, ierr, my_rank, nodes, i_st, i_end
real z(n), a, x(n), y(n)
```

```
call MPI_Init(ierr)
call MPI_Comm_rank(MPI_COMM_WORLD,my_rank,ierr)
call MPI_Comm_size(MPI_COMM_WORLD,nodes,ierr)
i_st = n/nodes*my_rank+1
i_end = n/nodes*(my_rank+1)
```

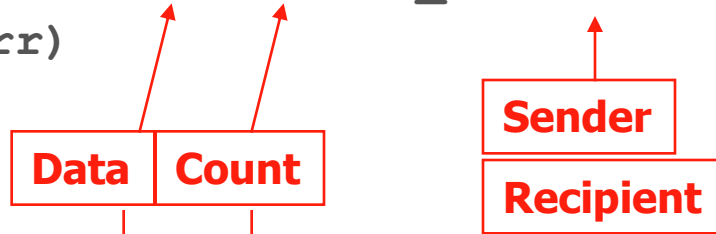
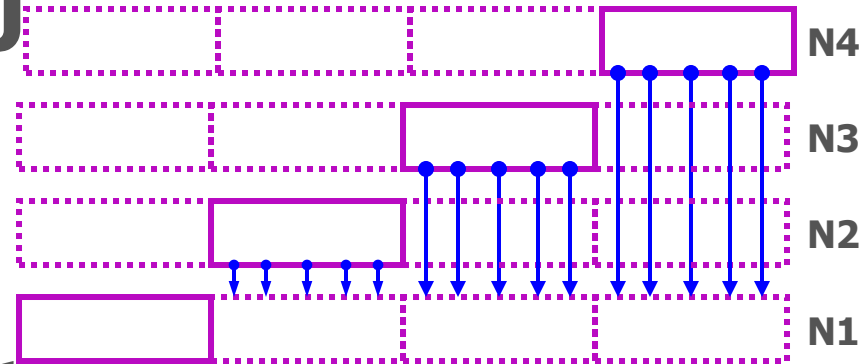
```
do i=i_st, i_end
  z(i) = a*x(i) + y(i)
enddo
call MPI_Finalize(ierr)
return
```

z(i) operation on 4 processes

z(1 ... n/4)	z(n/4+1 ... 2*n/4)	z(2*n/4+1 ... 3*n/4)	z(3*n/4+1 ... n)
-----------------	-----------------------	-------------------------	---------------------

• Result on the first CPU

```
include "mpif.h"
integer status(MPI_STATUS_SIZE)
if (my_rank .eq. 0 ) then
  do j = 1, nodes-1
    do i= n/nodes*j+1, n/nodes*(j+1)
      call MPI_Recv(z(i),1,MPI_REAL,j,0,MPI_COMM_WORLD,
&      status,ierr)
    enddo
  enddo
else
  do i=i_st, i_end
    call MPI_Send(z(i),1,MPI_REAL,0,0,MPI_COMM_WORLD,ierr)
  enddo
endif
```



- Collective communication

```
real zi(n)
```

```
j = 1
```

```
do i=i_st, i_end
  zi(j) = a*x(i) + y(i)
  j = j + 1
enddo
```

```
enddo
```

```
call MPI_Gather(zi, n/nodes, MPI_REAL, z, n/nodes, MPI_REAL,
& 0, MPI_COMM_WORLD, ierr)
```

Send data

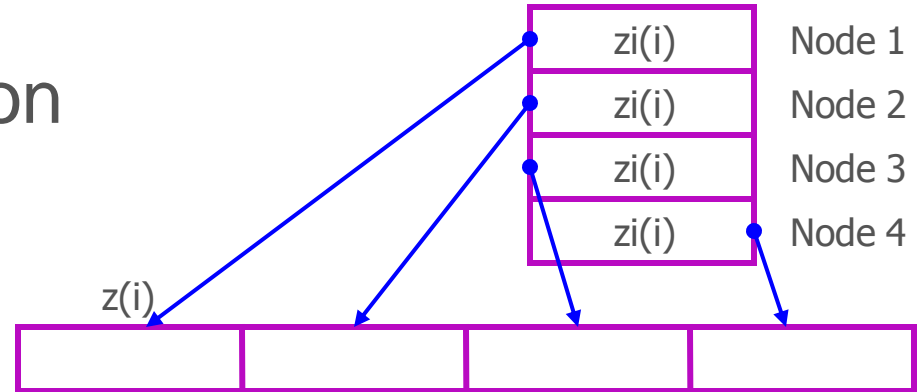
Receive data

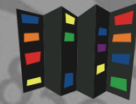
Root process

- Result on all nodes

```
call MPI_AllGather(zi, n/nodes, MPI_REAL, z, n/nodes,
& MPI_REAL, MPI_COMM_WORLD, ierr)
```

No root process





- First log into one of the clusters

```
ssh sanddunearch.chpc.utah.edu – Ethernet, InfiniBand
```

```
ssh updraft.chpc.utah.edu – Ethernet, InfiniBand
```

```
ssh ember.chpc.utah.edu – Ethernet, InfiniBand
```

```
ssh kingspeak.chpc.utah.edu – Ethernet, InfiniBand
```

- Then submit a job to get compute nodes

```
qsub -I -l nodes=2:ppn=4,walltime=1:00:00
```

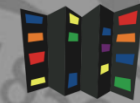
```
qsub script.pbs
```

- Useful scheduler commands

```
qsub – submit a job
```

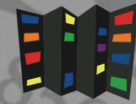
```
qdel – delete a job
```

```
showq – show job queue
```



- No clear text passwords use ssh and scp
- You may not share your account under any circumstances
- Don't leave your terminal unattended while logged into your account
- Do not introduce classified or sensitive work onto CHPC systems
- Use a good password and protect it

- Do not try to break passwords, tamper with files etc.
- Do not distribute or copy privileged data or software
- Report suspicions to CHPC (security@chpc.utah.edu)
- Please see <http://www.chpc.utah.edu/docs/policies/security.html> for more details



Arches

- Supported by most compilers, `-openmp`, `-fopenmp` or `-mp` switch

e.g. `pgf77 -mp source.f -o program.exe`

- Dual-processor and quad-processor (Sanddunearch) nodes
- Further references:

Compilers man page – `man ifort`

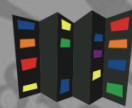
Compilers websites

<http://www.intel.com/software/products/compilers>

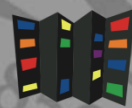
<http://gcc.gnu.org>

<http://www.pgroup.com/doc/>

<http://www.pathscale.com/node/70>



- Three common network interfaces
 - Ethernet, Myrinet, InfiniBand
- Different MPI implementations
 - MPICH2 - Ethernet, shmem
 - OpenMPI – Ethernet, InfiniBand
 - MVAPICH2 - InfiniBand
 - Intel MPI, Scali MPI - commercial



- **Clusters** – MPICH2, OpenMPI, MVAPICH2

```
/MPI-path/bin/mpixx source.x -o program.exe
```

- MPI-path = **location of the distribution**

```
/uufs/chpc.utah.edu/sys/pkg/mpich2/std MPICH2 TCP-IP
```

```
/uufs/$UUFSCELL.arches/sys/pkg/openmpi/std OpenMPI
```

InfiniBand

```
/uufs/$UUFSCELL.arches/sys/pkg/mvapich2/std
```

MVAPICH2 InfiniBand

= must specify full path to `mpixx (/MPI-path/bin)` or add this path to `PATH` environment variable

- **MPICH Interactive batch (incl. Updraft)**

```
qsub -I -l nodes=2:ppn=2,walltime=1:00:00
```

```
... wait for prompt ...
```

```
/MPI-path/bin/mpirun -np 4 -machinefile
```

```
$PBS_NODEFILE program.exe - Updraft, Ember, Kingspeak
```

```
/MPI-path/bin/mpirun -launcher rsh -np 4
```

```
-machinefile $PBS_NODEFILE program.exe -
```

```
sanddunearch
```

- **MPICH Batch**

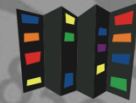
```
#PBS -l nodes=2:ppn=2,walltime=1:00:00
```

- **OpenMP Batch**

```
#PBS -l nodes=1:ppn=2,walltime=1:00:00
```

```
setenv OMP_NUM_THREADS 2
```

```
program.exe
```



- Useful for finding bugs in programs
- Several free
 - `gdb` – GNU, text based, limited parallel
 - `ddd` – graphical frontend for `gdb`
- Commercial that come with compilers
 - `pgdbg` – PGI, graphical, parallel but not intuitive
 - `pathdb`, `idb` – Pathscale, Intel, text based
- Specialized commercial
 - `totalview` – graphical, parallel, CHPC has a license
 - `ddt` - Distributed Debugging Tool
- How to use:
 - http://www.chpc.utah.edu/docs/manuals/software/par_devel.html

- Parallel debugging more complex due to interaction between processes
- Totalview is the debugger of choice at CHPC
 - Expensive but academia get discount
 - How to run it:
 - compile with `-g` flag
 - automatic attachment to OpenMP threads
 - extra flag (`-tv`) to `mpirun/mpiexec`

- **Details:**

`http://www.chpc.utah.edu/docs/manuals/software/totalview.html`

- **Further information**

`http://www.totalviewtech.com/alltvdocs.htm`

Etnus TotalView 5.0

File Edit View Tools Window Help

Attached Unattached Groups Log

1001795	T	mpirun (in __select)
1003847	B4	mpirun<trapp_r>.0 (in main)
1/1003847	B4	in main
1004955	B4	mpirun<trapp_r>.1 (in main)
1/1004955	B4	in main
1004633	B4	mpirun<trapp_r>.2 (in main)
1/1004633	B4	in main
1000345	B4	mpirun<trapp_r>.3 (in main)
1/1000345	B4	in main

Process view

mpirun<trapp_r>.0

File Edit View Group Process Thread Action Point Tools Window Help

Group Control Go Halt Next Step Out Run To Next Step P- P+ T- T+

Process 1003847: mpirun<trapp_r>.0 (At Breakpoint 4)

Thread 1003847.1: mpirun<trapp_r>.0 (At Breakpoint 4)

Stack Trace

C	main,	FP=7fff2f20	Function "main":
	__start,	FP=7fff2f30	argc: 0x00000001 (1)
			argv: 0x7fff2f34 -> 0x7fff30
			Local variables:
			p: 0x00000000 (0)
			tag: 0x00000000 (0)
			my_rank: 0x7fff2fc0 (2147430336)
			(Compound Object)
			status: 0x00000001 (1)
			n: 0x00000000 (0)
			i: 0x00000000 (0)
			local_n: 0x0fb70250 (263651920)
			...

Function main in trapp.c

```

33 MPI_Send(&b, 1, MPI_FLOAT, i, tag, MPI_COMM_WORLD);
34 tag = 2;
35 MPI_Send(&n, 1, MPI_INT, i, tag, MPI_COMM_WORLD);
36
37 }
38 }
39 else
40 {
41 tag = 0;
42 MPI_Recv(&a, 1, MPI_FLOAT, 0, tag, MPI_COMM_WORLD, &status);
43 tag = 1;
44 MPI_Recv(&b, 1, MPI_FLOAT, 0, tag, MPI_COMM_WORLD, &status);
45 tag = 2;
46 MPI_Recv(&n, 1, MPI_INT, 0, tag, MPI_COMM_WORLD, &status);
47 }
48
49 STOP h = (b-a)/n;
50 local_n = n/p;
51
52 local_a = a + my_rank*h*local_n;
53 local_b = local_a + h*local_n;
54
55 printf("%d %f %f %d\n", my_rank, local_a, local_b, local_n);
    
```

Source code view

Thread (1)

1/1003847	B4	in main	STOP	4	line 16	at main+0x24	in "trapp.c"
			BARR	7	line 27	at main+0x80	in "trapp.c"
			STOP	8	line 49	at main+0x1e8	in "trapp.c"

Action Points

main

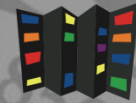
File Edit View Tools Window Help

my_rank (Laminated)

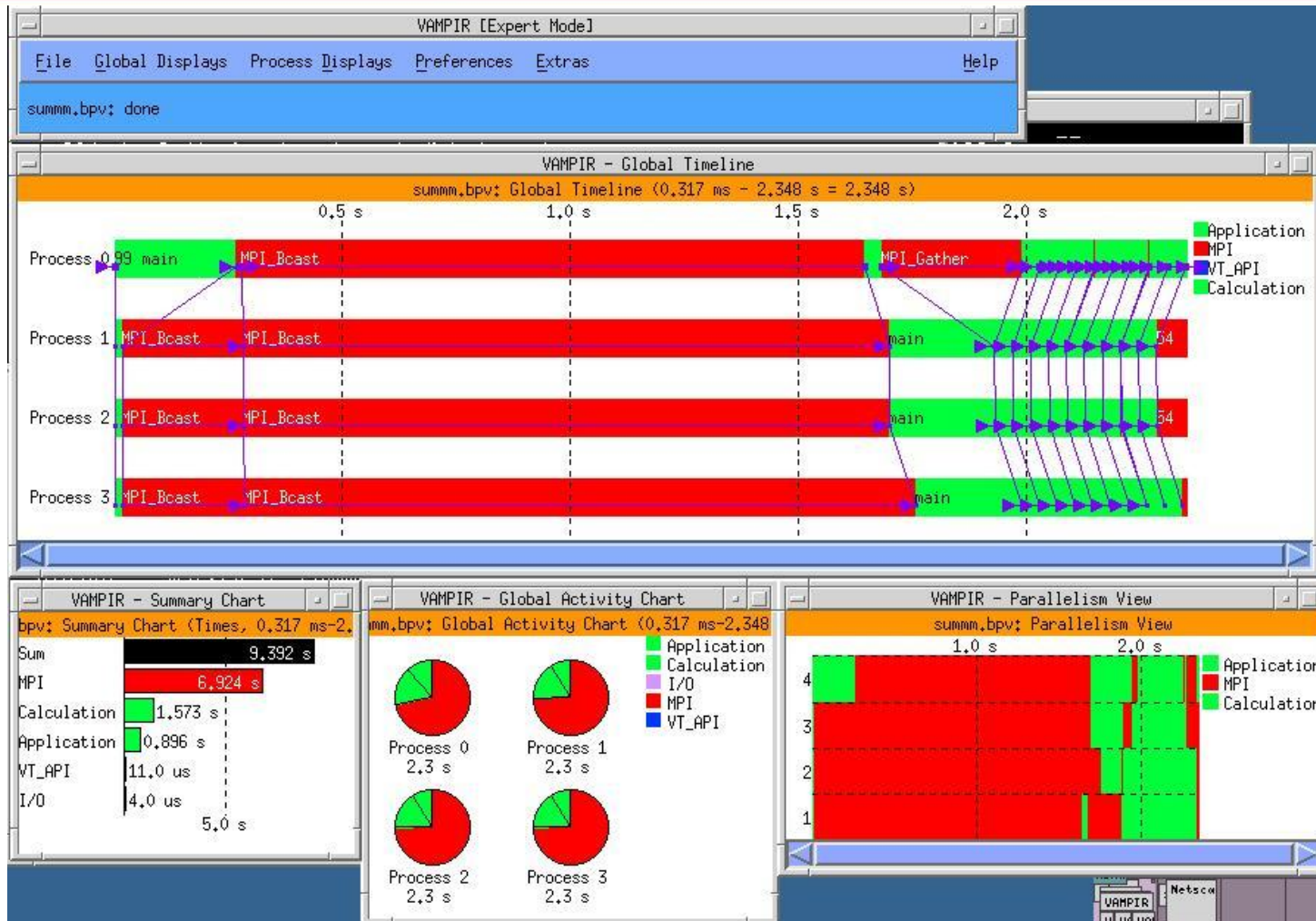
(at 0x7fff2ea4) Type: int
Filter:

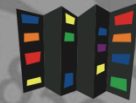
Process	Value
mpirun<trapp_r>.0	0x00000000 (0)
mpirun<trapp_r>.1	0x00000001 (1)
mpirun<trapp_r>.2	0x00000002 (2)
mpirun<trapp_r>.3	0x00000003 (3)

Data inspection



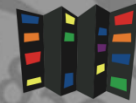
- Measure performance of the code
- Serial profiling
 - discover inefficient programming
 - computer architecture slowdowns
 - compiler optimizations evaluation
 - gprof, pgprof, pathopt2
- Parallel profiling
 - target is inefficient communication
 - Intel Trace Collector and Analyzer





- Shared vs. Distributed memory
- OpenMP
 - limited on Arches
 - Simple parallelization
- MPI
 - Arches
 - Must use communication

http://www.chpc.utah.edu/docs/presentations/intro_par



- OpenMP

<http://www.openmp.org/>

Chandra, et. al. - Parallel Programming in OpenMP

Chapman, Jost, van der Pas – Using OpenMP

- MPI

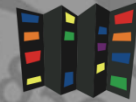
<http://www-unix.mcs.anl.gov/mpi/>

Pacheco - Parallel Programming with MPI

Gropp, Lusk, Skjellum - Using MPI 1, 2

- MPI and OpenMP

Pacheco – An Introduction to Parallel Programming



- Introduction to MPI
- Introduction to OpenMP
- Debugging with Totalview
- Profiling with TAU/Vampir
- Intermediate MPI and MPI-IO
- Mathematical Libraries at the CHPC