

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat University of Leeds edu_assist_pr

Lecture 8: Introduction to distributed me

Previous lectures

Assignment Project Exam Help

In the last six lectures we looked at **shared memory parallelism** (SMP) relevant to e.g. multi-core CPUs:



<https://eduassistpro.github.io>



Without proper **synchronisation**, r
non-deterministic



Dependencies can lead to **data race**



Can reach **deadlock** if threads wait for synchronisation events that never occur.

This lecture

Assignment Project Exam Help

This lecture is the first of six on **distributed memory parallelism**, and we will see that some (but not all) of these issues remain relevant

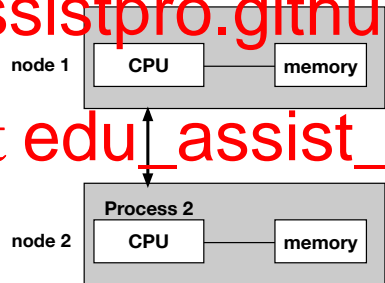
- <https://eduassistpro.github.io>
- **tion.**
- No **data races**.
- Performance considerations remain, in primary parallel overhead is **com**
- Improper synchronisation can still lead to **non-determinism** and **deadlock**.

Distributed memory systems

Assignment Project Exam Help

Multiple processes (rather than threads) that communicate via an interconnection network or 'interconnect'.

- Each process has its own heap memory
- If a process needs data currently held on another node's memory, must **communicate** over the network.



Current fastest supercomputer¹

Fujitsu Fugaku, RIKEN, Kobe, Japan

- ARM-based A64FX CPU.
- 4 as
- T
- No GPUs.
- Draws nearly 30MW of power.
- Benchmarked ≈ 442 PFLOPS.
- $1 \text{ PFLOPS} = 10^{15} \text{ FLOPS}$.
- $1 \text{ FLOPS} = 1$ floating point operation per second.

¹As of Nov. 2021; top500.org.

Clusters as distributed systems

Assignment Project Exam Help

Supercomputers share features with other distributed systems such as data centres:



<https://eduassistpro.github.io>

- May have high energy demand and cooling r

Here focus on High Performance Computing

- Individual cluster nodes use the same **operating system**.
- Cannot usually be **addressed individually**.
- Requires a special **job scheduler**.

The interconnection network or 'interconnect'

For the local area networks within HPC clusters, communication between nodes is carried over high performance **interconnects**.



1.



<https://eduassistpro.github.io>

These numbers are improving with time but

CPU performance

Add WeChat edu_assist_pro

The need to reduce communication overhead
more important in the foreseeable future.

¹As of Nov. 2021; see top500.org.

Network topology

If data sent via intermediary nodes, latency is increased

- Each node must parse data packet and decide where to send.

Then

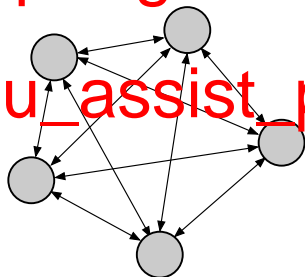
Net



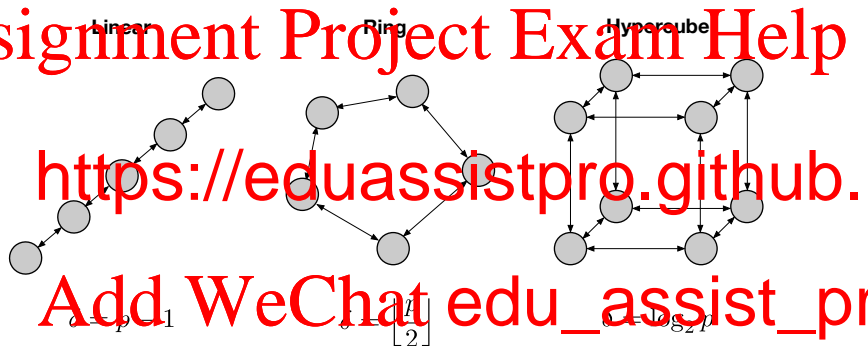
- E = connections (edges).

Want G with smallest **diameter** δ
(largest path length between nodes).

A **complete graph** (right) has $\delta = 1$,
but is impractical (*too many connections for each machine*).



Example topologies for p nodes



Hypercube topology preferred due to its short path lengths¹.

¹Rauber and Rünger, *Parallel programming for multicore and cluster systems* (Springer, 2013).

Processes *versus* threads

Recall from Lecture 2 that **processes** communicate with other processes using e.g. sockets.



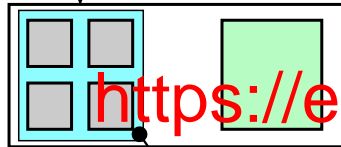
<https://eduassistpro.github.io>
For m nodes, with one thread per core.

- Avoids communication **within** a node
- Combination of OpenMP and MPI is quite efficient

For simplicity, we consider one **single-threaded process per core**, and therefore **multiple processes per node**.

Example for quad core nodes

One 4-thread process per node

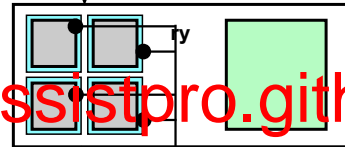


single process

memory

node 2

4 one-thread processes per node



node 1

node 2

Books

Assignment Project Exam Help

Willinson and Allen [Lecture 1] covers distributed memory parallelism (MPI), and a little OpenMP, but no GPU.



- <https://eduassistpro.github.io>

A more practical book for MPI coding is:

- **Parallel Programming with MPI**
(Morgan-Kaufman).

- Old (1997), only covers distributed memory systems and MPI.
- Many code examples and snippets.

Distributed HPC programming

Assignment Project Exam Help

For distributed HPC, there is essentially only one option¹ **MPI**

- Stands for **Message Passing Interface**.

-

ing').

-

-

- Fully supports C, C++ and FORTRAN.

- Most online examples are in one of these languages.

- Unofficial bindings for Java, MATLAB, P

¹Has superseded PVM = Parallel Virtual Machine (1989). Others such as Spark, Chapel etc. not (yet?) widely used in HPC.

Implementations

The MPI standard only defines the interface, it's still down to a vendor to provide an **implementation**.



There

- **MPICH**: www.mpich.org

- **OpenMPI**: www.open-mpi.org

- Don't confuse OpenMPI with OpenMP ...!

There are also commercial implementations:

- e.g. Intel MPI, Spectrum MPI (IBM).

Installing MPI

The system `cloud-hpc1.leeds.ac.uk` has OpenMPI¹ installed:

```
module load mpi/openmpi-x86_64
```

For pe
(cf. *lin*)

- Mac users might like to try home

On Windows machines, Microsoft MPI

- Based on MPICH.

¹Note the linux command “`module avail`” shows what modules are installed.

²<https://docs.microsoft.com/en-us/message-passing-interface/microsoft-mpi>

Building an MPI program

Assignment Project Exam Help

Need to use a special compiler for MPI programs:

- Standard installation includes `mpicc`, `mpic++`, `mpifort`.

-

- <https://eduassistpro.github.io>

For example, to compile a file `helloW`

Add WeChat edu_assist_pro

- `mpicc -Wall -o helloWpr hello`
- Will generate the executable `hel`
- All warnings on (`'-Wall'`).
- Add e.g. `-lm` for the maths library.

Executing an MPI program

Also need a special **launcher** to execute an MPI program¹.

For multiple processes all on the same local machine:

```
mpiexec -n 2 ./helloWorld
```



<https://eduassistpro.github.io>



mpirun is the same/very similar to

Best to develop/debug code on a single machine (e.g. of cloud-hpc1.leeds.ac.uk), then run in batch mode for e.g. timing runs.

¹Executing as usual ('./helloWorld') will launch *one* process, i.e. serial.

²With OpenMPI, can override with the argument `-oversubscribe`.

Launching via the batch queue

The system `cloud-hpc1.leeds.ac.uk` has been set up to allow access to two 8-core nodes via `slurm`.

enMP:

- <https://eduassistpro.github.io>

```
#!/bin/bash
```

```
#Request a single node, and 8 cores (adjust as n
```

```
#SBATCH -N1 -n8
```

```
module add mpi/openmpi3-x86_64
```

```
mpiexec -n 8 ./helloWorld
```

A 'Hello World' example

```
1 #include "stdio.h"
2 #include "mpi.h"
3 // Need to include mpi.h
4
5 int main
6 {
7     int n;
8
9     MPI_Init( &argc, &argv );
10    MPI_Comm_size( MPI_COMM_WORLD, &numprocs );
11    MPI_Comm_rank( MPI_COMM_WORLD, &rank );
12
13    printf( "Process %d of %d.\n", rank, numprocs );
14
15    MPI_Finalize();
16    return EXIT_SUCCESS;
17 }
```

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Initialising and finalising

Assignment Project Exam Help

The first MPI call **must** be `MPI_Init()`:

- Pass command line arguments `argc` and `argv`.

-

- <https://eduassistpro.github.io>

The final MPI call **must** be `MPI_Finalize()`:

- Note the JS spelling *finalize* not

Add WeChat edu_assist_pro

Any MPI calls before `MPI_Init()` or after `MPI_Finalize()` will result in a runtime error.

Number of processes and rank

Assignment Project Exam Help

`MPI_Comm_size(MPI_COMM_WORLD, &numprocs)`



<https://eduassistpro.github.io>

`MPI_Comm_rank(MPI_COMM_WORLD, &rank)`



Sets rank to the process number, known as the rank.



Ranges from 0 to numprocs-1 inclusive.



Similar to `omp_get_thread_num()`.

Communicators

Assignment Project Exam Help

For our purposes, whenever you see an MPI call with the argument **communicator**, just use `MPI_COMM_WORLD`:

-
- <https://eduassistpro.github.io>

In general, communicators allow processes to be

- e.g. when developing a parallel library, if processes to accidentally communicate processes.
- An advanced feature we won't consider.

Summary and next lecture

Assignment Project Exam Help

Today we have started looking at **distributed memory parallelism**:

- <https://eduassistpro.github.io>
 - For HPC, use **MPI = Message Passing Interface**.
 - Seen how to build and execute a 'Hello World'
- Add WeChat edu_assist_pr

Next time we will see how MPI supports communication between processes, and use this to solve real problems.