

Assignment Project Exam Help

<https://eduassistpro.github.io>

University of Leeds
Add WeChat edu_assist_pr

Lecture 10: Parallel data reorg

Previous lectures

Assignment Project Exam Help

In the last lecture we saw how to perform **explicit point-to-point communication** in a distributed memory system:



<https://eduassistpro.github.io>



- Sending process calls `MPI_Send`

- Receiving process calls `MPI_Recv`



Both **blocking** calls that do not return until safely modified.



Can result in **deadlock**, e.g. cyclic communication pattern.

This lecture

Assignment Project Exam Help

In this lecture we are going to look at one of the key considerations for the **performance** of distributed memory systems: **Data reorg**

- <https://eduassistpro.github.io>
- For distributed systems, data reorganisa significant **parallel overhead**.
- Improved performance using **collective routines**.
- Will go through a worked example of a simple **distributed counting algorithm**.

Data reorganisation

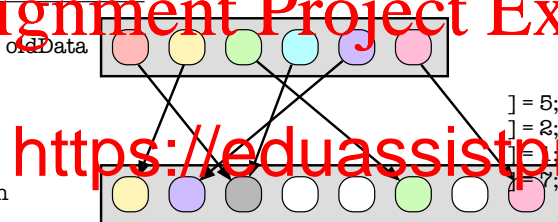
Assignment Project Exam Help

Many algorithms require some form of **large-scale data reorganisation**

- <https://eduassistpro.github.io> or, stack, queue etc.) or a **database**.
- **Numerical algorithms**, i.e. reordering matrix
- **Compression** (e.g. bzip, gzip etc)
- ...

Generalised scatter and gather

General scatter:¹



In serial code:

```
1 for( i=0; i<N; i++ )  
2   newData[ index[i] ] = oldData[i];
```

General **gather** is similar, but indices give **read** locations.

¹McCool *et al.*, *Structured parallel programming* (Morgan-Kaufman, 2012).

Shared *versus* distributed

Data reorganisation in shared memory systems can lead to a **data race or race condition**:



- <https://eduassistpro.github.io> with associated performance penalty.

Although data races are not relevant to distributed systems, data reorganisation is very important:

The primary overhead in distributed systems is **communication**, which is a form of data reorganisation

Communication performance

Although most of the overheads in Lecture 4 apply to distributed systems, one typically dominates. The communication time.

If the scaling is poor, the communication time dominates the calculation.

<https://eduassistpro.github.io>

For communication to *not* adversely affect the ratio

$$\frac{t_{\text{comm}}}{t_{\text{comp}}}$$

to be as small as possible.

¹Recall from Lecture 4 that t_p is the parallel execution time.

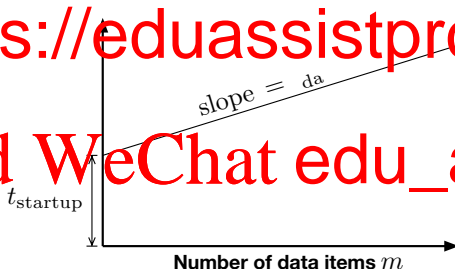
Analysis of t_{comm}

For a single message of size m , a good approximation to t_{comm} is¹:

$$t_{\text{comm}} = t_{\text{startup}} + mt_{\text{data}}$$

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



¹Wilkinson and Allen, *Parallel programming* 2nd ed. (Pearson, 2005).

Measurement of t_{comm} from SoC lab machines (in Leeds)

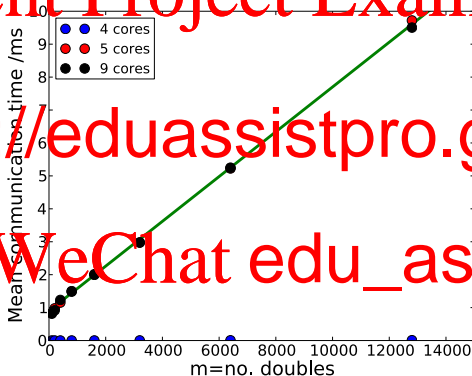

Code on Minerva: `measure_tComm.c`

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

$\approx 0.9\text{ms}$

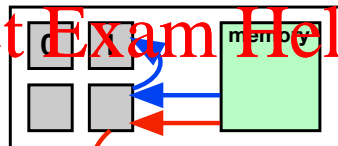


For faster interconnects both t_{startup} and t_{data} about 10 times smaller, but **communication remains the primary overhead.**

Intra- versus inter-node communication

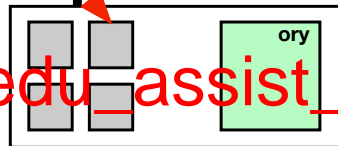
Process 3 sends data to process 1, it is copied within the same machine's memory (blue arrows).

- F



If process 3 now sends the same data to process 5, it is sent over the network (red arrows).

- Slow



¹Could be removed by using one *multi-threaded* process per node [Lecture 8].

Strategies to reduce communication times

Assignment Project Exam Help

This comm suggests we should merge messages when possible

- For **two** messages of size m and n :

<https://eduassistpro.github.io>

- For **one** message of size $m + n$:

Add WeChat edu_assist_pro

$$t_{\text{comm}} = t_{\text{startup}} +$$

So we have **saved** t_{startup} in total communication time.

We will see another strategy in Lecture 12.

Collective communication

Assignment Project Exam Help
An alternative and often easier way to reduce communication times is to use **collective communication**:



<https://eduassistpro.github.io>

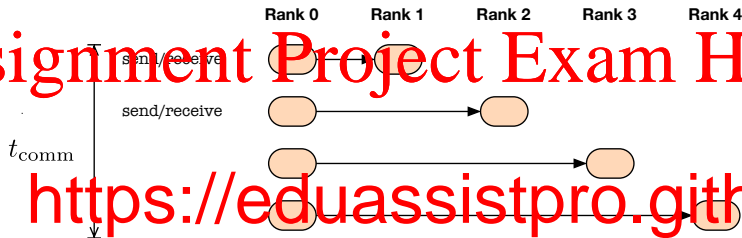
Distributed programming APIs include

common communication

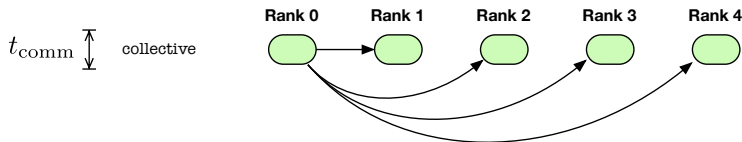
Add WeChat edu_assist_pro

- Can **drastically reduce** the communication overhead.
- Implementation varies, but typically **overlap** communications to reduce t_{comm} .

Point-to-point communication:



Collective communication (ideal case):



Common forms of collective communication

Distribution	Type	Meaning
One-to-all	Broadcast	Same data from one process to all
Many-to-one	Gather	Data from many processes to one
Many-to-many	Reduction	Data from many processes to one, then broadcast back to all

Other variants (*i.e.* *many-to-many* such as *multi-broadcast*) exist but are less commonly used and not considered here.

We will consider *reduction* next lecture.

Collective communication in MPI

Code on Minerva: `distributedCount.c`

Assignment Project Exam Help

To demonstrate collective communication in MPI, we will use a simple worked example: A **distributed count** algorithm.

1

2

<https://eduassistpro.github.io>

3

Each rank (including rank 0) counts how many data are below some threshold.

4

All ranks > 0 send their counts to rank 0, which calculates the total.

Note we assume only rank 0 knows the total data size.

- e.g. if rank 0 had loaded the data from a file.

Step 1: Broadcasting: MPI_Bcast()

Assignment Project Exam Help

Sending the variable `localSize` to all processes can be performed using point-to-point communication.

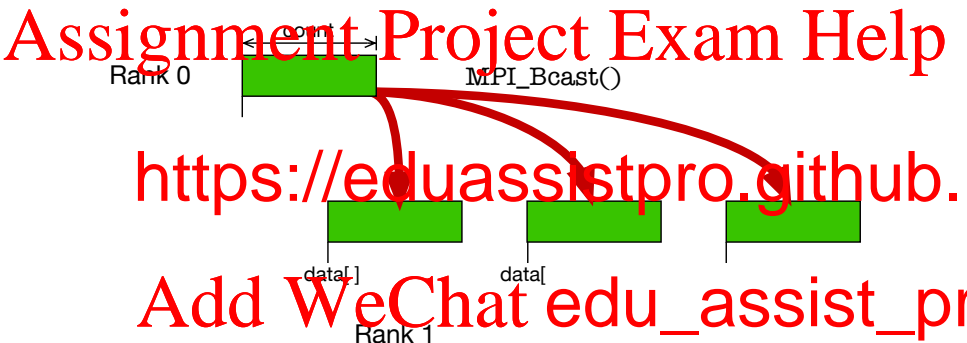
```
1 if( rank
2     fo
3
4 else
5     MPI_Recv(&localSize,1,MPI_INT,0,...);
```

The same thing can be achieved using

```
1 MPI_Bcast(&localSize,1,MPI_INT,0,MPI_COMM_WORLD
```

- First 3 arguments same as MPI_Send()/MPI_Recv().
- Fourth argument is the rank on which `localSize` is defined.

Broadcasting: Schematic



Note that using '`&variable`' for the data argument 'fools' MPI into thinking the variable count is an array of size 1.

Common pitfall - careful!

When using collective communication, it is important to realise that **all** processes are involved.



This

<https://eduassistpro.github.io>

```
1 if( rank==0 ) MPI_Bcast(...);
```

- `MPI_Bcast()` does not return until called
- Ranks > 0 do **not** call `MPI_Bcast`
- Rank 0 will wait forever - **deadlock**.

The name *broadcast* is misleading as it suggests only **sending** is involved, whereas in fact it also includes the **receiving**.

Step 2: Scattering: MPI_Scatter()

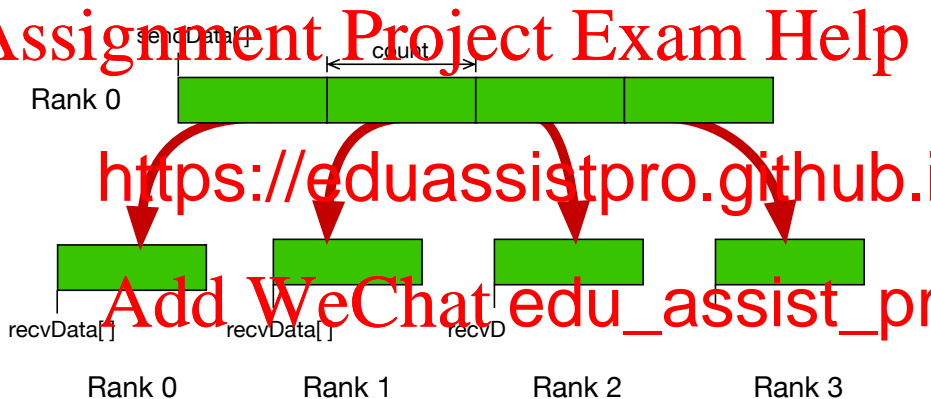
Need to break up an array into equal sized chunks and send one chunk to each process *[cf. vector addition last lecture]*:

```
1 if( rank==0 )  
2     fo  
3  
4 else  
5     MP
```

This can be replaced with a single call:

```
1 MPI_Scatter(  
2     globalData,localSize,MPI_INT, // Sent from  
3     localData ,localSize,MPI_INT, // Received to  
4     0, MPI_COMM_WORLD             // Source rank 0  
5 );
```

Scattering: Schematic



Note also copies to `recvData[]` on rank 0.

Step 4: Gathering: MPI_Gather()

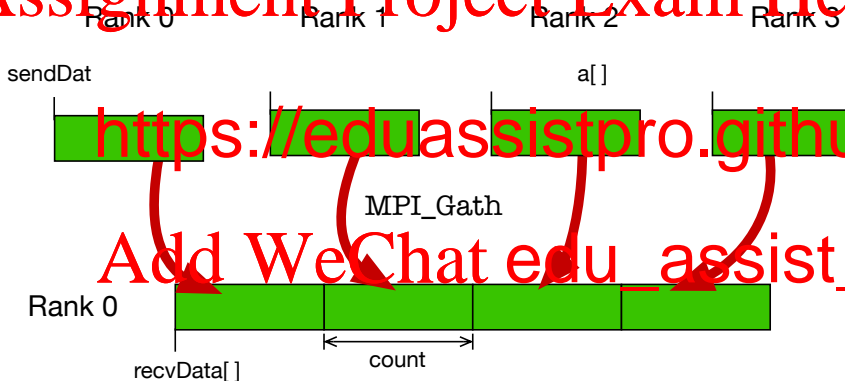
Gathering is the opposite of scattering:

```
1 MPI_Gather(  
2     &count, 1, MPI_INT, // Sent from  
3     pa  
4     0, MPI // D  
5 );
```

This gives an array of `Procs`,
from which the total can be counted. As with scatter

- Data is ordered by rank.
- There is no tag.
- The data size is the **local** size, both times.
- Can **in principle** use different data sizes or types for sending and receiving, but not recommended.

Gathering: Schematic



Summary and next lecture

Assignment Project Exam Help

Today we have looked at **data reorganisation and collective communication**

- <https://eduassistpro.github.io>
- **Broadcasting** (e.g. `MPI_Bcast`)
- **Scattering** (e.g. `MPI_Scatter`)
- **Gathering** (e.g. `MPI_Gather`).

In fact, the last stage of our example involved data reorganisation **and calculation**.

- This **reduction** is the subject of the next lecture.