Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

University of Leeds

Lecture 11: Reducti

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Previous lectures
Today's lecture

## Previous lectures

In the last lecture we looked at **data reorganisation** and **colle**

- Collective communication invol one-to-many, many-to-one or many-to

- Reduce the communication time *t* point-to-point communications.

- In MPI: `MPI_Bcast()`, `MPI_Scatter()`, `MPI_Gather()`.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Previous lectures
Today's lecture

# Today's lecture

Today we will look at a common combination of data reor

- https://eduassistpro.github.i
-
- Support for many parallel APIs, including
- Often optimised using a **binary tre**
- Binary trees also useful for collective comm

Add WeChat edu_assist_pr

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
MapReduce
Library support for reduction

## Reminder: Serial reduction

- Start with a large data set.
  - Apply **binary operations** to *reduce* to a smaller set.

Example 1: Finding the sum

```
1  sum = 0;
2  for( i=0; i<N; i++ )
3    sum += a[i];
```

Example 2: Finding the maximum element

```
1  max = a[0];
2  for( i=1; i<N; i++ )
3    if( a[i]>max ) max = a[i];
```

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

**Reduction in serial**
Reduction in parallel
MapReduce
Library support for reduction

Note each operation is performed **sequentially**
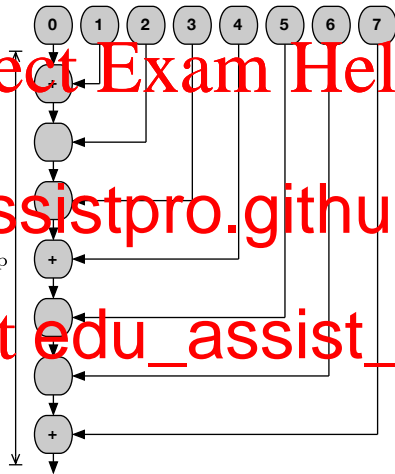
Total computation time $t_{\mathrm{comp}}$
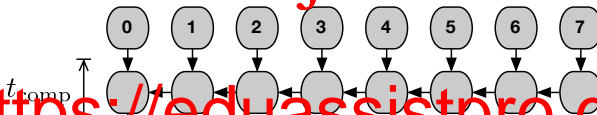
array s

- *i.*
  **complexity** is $\mathcal{O}(n)$.

If these were **processing units**, most would be **idle** throughout most of the calculation.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
MapReduce
Library support for reduction

## Parallel reduction

Ideally, we would want to perform all calculations **simultaneously**:



This *would* have a time complexity of $t = (1)$, but is not possible to achieve in practice.

For now, note that:

> **Any** parallel reduction **must** change the sequence of calculations

Some concrete examples will be given later in this lecture.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
MapReduce
Library support for reduction

## Recap: Commutativity and associativity

Let $\otimes$ denote a general binary operator: $c = a \otimes b$

As parallel reduction alters the sequence in which calculations are perfo

https://eduassistpro.github.i⊗d

If $\otimes$ is only **approximately associative**
reduction will be (slightly) **different fro**

Some parallel reduction algorithms also requir
**commutative**:

> An operator $\otimes$ is **commutative** if $a \otimes b = b \otimes a$

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
MapReduce
Library support for reduction

## Commutativity and associativity (examples)

Assignment Project Exam Help

| Combination | Examples |
| --- | --- |
| Ass | |
| | https://eduassistpro.github.i |
| Ass | |
| Commutative; **not** associative | **Fin** t |
| | add [1] |
| | sig |
| | Add WeChat edu_assist_pr |
| **Neither** commutative **nor** associative | Subtraction, division |

[1] Only *approximately* associative. See Worksheet 2 Question 6.
[2] *e.g.* `fn(a,b)=(a+b<1?a+b:1)` with a=0.8, b=0.5 and c=−0.3.

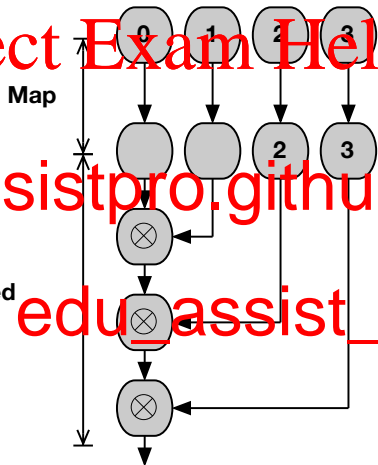| Overview | Reduction in serial |
| Parallel reduction: Overview and library support | Reduction in parallel |
| Parallel reduction: Implementation | **MapReduce** |
| Summary and next lecture | Library support for reduction |

# MapReduce

An important application of reduct

**Map**



**Map**

**Red**

- Fusion of a **map** followed by a **reduction**.
- Can avoid the need for **synchronisation** after the map.

---

[1]McCool *et al.*, *Structured parallel programming* (Morgan-Kaufmann, 2012).

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
**MapReduce**
Library support for reduction

## Distributed systems example

Suppose a database is **distributed** over nodes *in a cluster*.

- Each node has access to part of the full database.

Suppose https://eduassistpro.github.i

- Each node searches its local database *('map')*.
- Local results are combined to give the requir *('reduce')*.

Add WeChat edu_assist_pr

This **MapReduce** was developed by Google and was one of the reasons for their early success.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
MapReduce
Library support for reduction

## Example: Vector dot product

Consider the **vector dot** product (*aka* inner or scalar product)[1]:

$$\mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^{n} a_i b_i$$

In serial:

```
1  float dot=0.0;
2  for( i=0; i<n; i++ )
3    dot += a[i] * b[i];
```

Note this is a **map** (*the multiplication*) followed by a **reduction** (*the summation*).

---

[1]Recall maths indexing starts from 1 but computer indexing starts from 0.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
MapReduce
Library support for reduction

## Reduction in OpenMP
Code on Minerva: `dotProduct_OpenMP.c`

In OpenMP (i.e. shared memory systems), reduction is supported
by the **reduction clause**:

```
1  float do
2  #prag
3  for( i=0
4    do  +
```

- Specify the **binary operation** ('+'
  ('dot').
- Compiler and runtime will implement an **efficient** reduction
  **for the given architecture**.
- Details of the implementation **opaque** to the user.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
MapReduce
Library support for reduction

## Reduction in MPI
Code on Minerva: `dotProduct_MPI.c`

For MPI, distribute the full arrays on rank 0 to local arrays on each process using `MPI_Scatter()`[1]:

```
1  MPI_S
2  MPI_S
```

Each process then calculates its own local dot prod

```
1  float local_dot=0.0;
2  for( i=0; i<numPerProc; i++ )
3    local_dot += local_a[i]*local_b[i];
```

---

[1]This step is the same as for vector addition; cf. Lecture 9.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
MapReduce
Library support for reduction

# `MPI_Reduce()`

The MPI standard supports reduction through `MPI_Reduce()`:

```
1  float dot;
2  MPI_R
```

- <span style="color:red">https://eduassistpro.github.i</span>

- Applied to `local_dot` on all process

- Reduced to `dot` on rank 0 (the $6^{t}$

- Other operations are supported, e
  `MPI_MIN`, logical and binary boolean operators.

- Implementation opaque to the user, but *should* be optimised
  for the system on which it is installed.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Efficient parallel reduction
Binary tree reduction
Need for synchronisation
Binary trees in collective communication

## Efficient parallel reduction

How OpenMP and MPI implement reduction is not specified by their respective standards.

- 

Usual https://eduassistpro.github.i
to consider possible implementation details to help understand performance and identify potential issues.

Parallel reduction starts after each of ads, processes) have completed their **local reduction**.

- That is, calculated the **partial sums** of all the data each processing unit is 'responsible' for.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Efficient parallel reduction
Binary tree reduction
Need for synchronisation
Binary trees in collective communication

# Binary trees

Assignment Project Exam Help

The most common method to implement parallel reduction is with
a **binary tree**:

- https://eduassistpro.github.i

- Perform calculations **in parallel** at each level.

- Reduction time is then $\mathcal{O}(\log_2(p)$ an
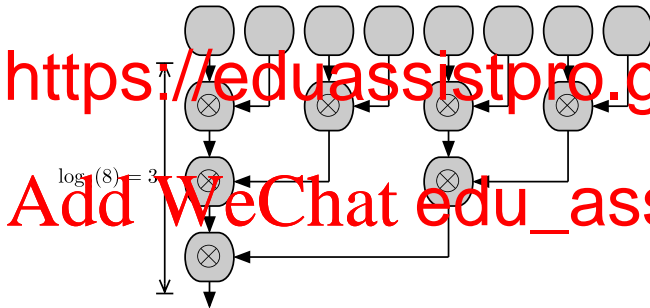$\mathcal{O}(p)$ for large $p$

Add WeChat edu_assist_pr

---

[1] If $p$ is not a power of 2, round up.

Overview
Parallel reduction: Overview and library support
**Parallel reduction: Implementation**
Summary and next lecture

Efficient parallel reduction
**Binary tree reduction**
Need for synchronisation
Binary trees in collective communication

# Binary tree: Example 1

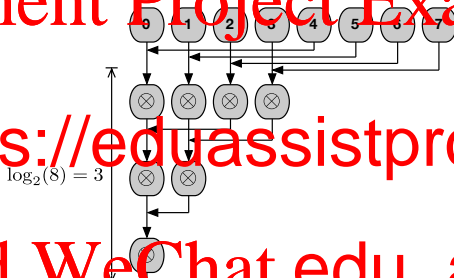Not all binary trees are valid for all binary operators $\otimes$.

For instance, this version requires that $\otimes$ be **associative**:



$\log_2(8) = 3$

The **indexing**, *i.e.* which processing units are performing the operations at each level, can be performed using bitwise arithmetic.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Efficient parallel reduction
Binary tree reduction
Need for synchronisation
Binary trees in collective communication

# Binary tree: Example 2

For this example, $\otimes$ must also be **commutative**:



$\log_2(8) = 3$

Indexing is easier than the previous example:

- In the first level, units 0 to $p/2$ perform the operations.
- In the next level, units 0 to $p/4$ perform the operations.
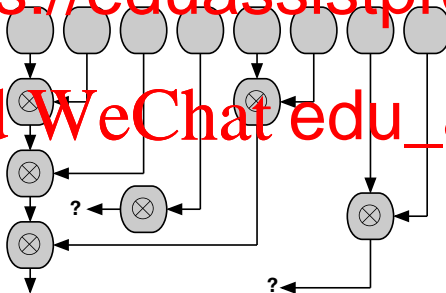- . . .

Overview
Parallel reduction: Overview and library support
**Parallel reduction: Implementation**
Summary and next lecture

Efficient parallel reduction
Binary tree reduction
**Need for synchronisation**
Binary trees in collective communication

## Synchronisation between levels

Note we must ensure each level's calculation have been **completed** before continuing to the **next** level.

This e
at bes

Overview
Parallel reduction: Overview and library support
**Parallel reduction: Implementation**
Summary and next lecture

Efficient parallel reduction
Binary tree reduction
**Need for synchronisation**
Binary trees in collective communication

## Barriers

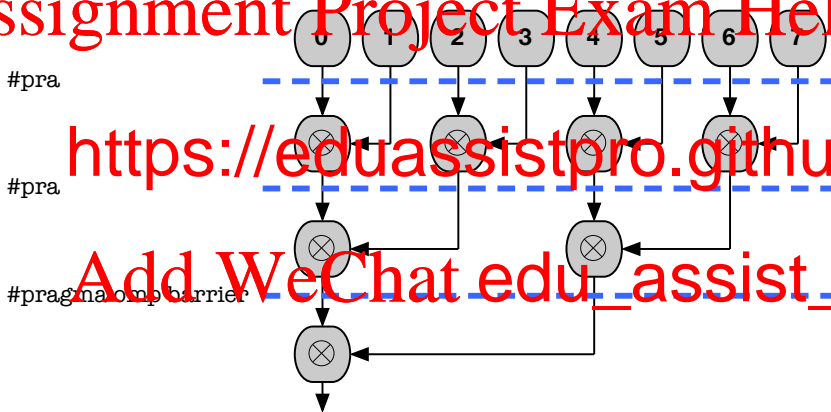Most parallel APIs provide a means to synchronise all processing units

For instance, in OpenMP (in a parallel region):

```
1  #pragma omp barrier
```

- No processing unit (*i.e.* thread) will proceed past the barrier command until **all** units have reached it.

Overview
Parallel reduction: Overview and library support
**Parallel reduction: Implementation**
Summary and next lecture

Efficient parallel reduction
Binary tree reduction
**Need for synchronisation**
Binary trees in collective communication

# Barrier synchronisation in a binary tree

Overview
Parallel reduction: Overview and library support
**Parallel reduction: Implementation**
Summary and next lecture

Efficient parallel reduction
Binary tree reduction
**Need for synchronisation**
Binary trees in collective communication

# Synchronisation in MPI

MPI also provides a barrier operation:

```
1  MPI_B
```

How https://eduassistpro.github.i

can be

- `MPI_Send()`, `MPI_Recv()` will not re
  been sent or received
- Provides the necessary synchronisation
  processes.

Overview
Parallel reduction: Overview and library support
**Parallel reduction: Implementation**
Summary and next lecture

Efficient parallel reduction
Binary tree reduction
Need for synchronisation
**Binary trees in collective communication**

## Binary trees in collective communication

Assignment Project Exam Help

Note that `MPI_Reduce()` is a **collective communication**:

- 

https://eduassistpro.github.i

The binary tree pattern is typically used for all collective communication.

- Communication time $t_{comm} = \mathcal{O}($ Add WeChat edu_assist_pr

- Faster than the $\mathcal{O}(p)$ for a loop of send-a

- 'Inverted' in the case of `MPI_Bcast()` and `MPI_Scatter()`.

Overview
Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Summary and next lecture

# Summary and next lecture

Today we have looked at **parallel reduction**:

-
-
- In MPI, the necessary synchronisation pro
  **blocking communication**.

Next time we will look at **non-blocking**, or **asynchronous**, communication.