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

# XJCO3221 Parallel Computation

Peter Jimack

University of Leeds

Lecture 11: Reduction

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 **collective communication**:

- **Communication** is usually the most significant overhead for distributed systems.
- **Collective communication** involves multiple processes in a one-to-many, many-to-one or many-to-many pattern.
- Reduce the communication time $t_{comm}$ compared to a loop of 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 reorganisation and computation: **Reduction**.

- **Reduces** a data set to one of a smaller size.

- Important for both shared and distributed memory systems.

- Support for many parallel APIs, including OpenMP and MPI.

- Often optimised using a **binary tree**.

- Binary trees also useful for collective communication.

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: Summing the elements of an array

```
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
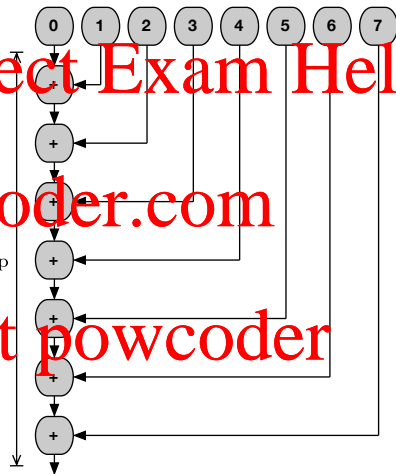MapReduce
Library support for reduction

Note each operation is performed **sequentially**.

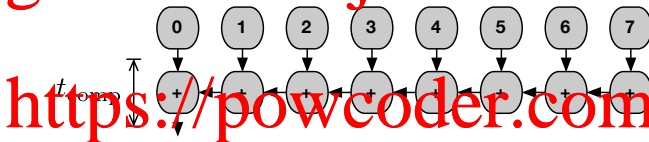Total computation time $t_{\mathrm{comp}}$ is proportional to the array size $n$.

- *i.e.* the **time 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 $t_{\mathrm{comp}} = \mathcal{O}(1)$, 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 performed, $\otimes$ must be **associative**:

> An operator $\otimes$ is **associative** if $a \otimes (b \otimes c) = (a \otimes b) \otimes c$

If $\otimes$ is only **approximately associative**, the result of parallel reduction will be (slightly) **different from serial reduction**.

Some parallel reduction algorithms also require $\otimes$ to be **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)

| Combination | Examples |
|---|---|
| Associative **and** commutative | max, min, Boolean AND, OR, XOR, **exact** addition and multiplication |
| Associative; **not** commutative | Matrix multiplication |
| Commutative; **not** associative | **Finite precision** floating point addition and multiplication[1], signed saturated addition[2] |
| **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
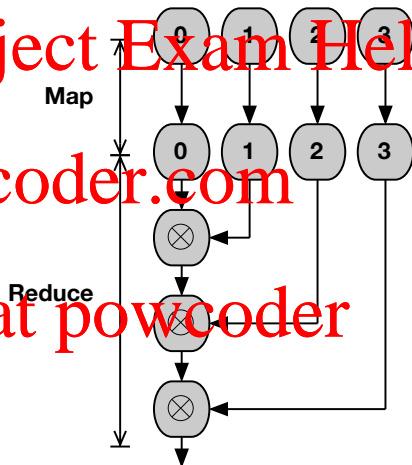Parallel reduction: Overview and library support
Parallel reduction: Implementation
Summary and next lecture

Reduction in serial
Reduction in parallel
MapReduce
Library support for reduction

# MapReduce



An important application of reduction is as part of the **MapReduce** pattern[1]:

- 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 a user initiates a search. We could use **MapReduce**:

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

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):

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

In serial[1]:

```
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  dot=0.0;
2  #pragma omp parallel for reduction(+:dot)
3  for( i=0; i<N; i++ )
4    dot += a[i] * b[i];
```

- Specify the **binary operation** ('+') and the **target variable** ('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_Scatter(a,numPerProc,MPI_FLOAT,local_a,numPerProc,
       MPI_FLOAT,0,MPI_COMM_WORLD);
2  MPI_Scatter(b,numPerProc,MPI_FLOAT,local_b,numPerProc,
       MPI_FLOAT,0,MPI_COMM_WORLD);
```

Each process then calculates its own local dot product:

```
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_Reduce(&local_dot,&dot,1,MPI_FLOAT,MPI_SUM,0,
        MPI_COMM_WORLD);
```

- Binary operator specified (`MPI_SUM`).

- Applied to local_dot on all processes.

- Reduced to dot on rank 0 (the $6^{\text{th}}$ argument).

- Other operations are supported, e.g. MPI_PROD, MPI_MAX,
  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.

- Allows **optimisation** for specific hardware architectures.

Usually best to use the support as provided, but sometimes useful to consider possible implementation details to help understand performance and identify potential issues.

Parallel reduction starts after each of $p$ **processing units** (threads, 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

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

- One 'leaf' node for each processing unit.
- For $p$ processing units, there are $\log_2(p)$ levels[1]
- Perform calculations **in parallel** at each level.
- Reduction time is then $\mathcal{O}(\log_2(p))$, which is **much** faster than $\mathcal{O}(p)$ for large $p$.
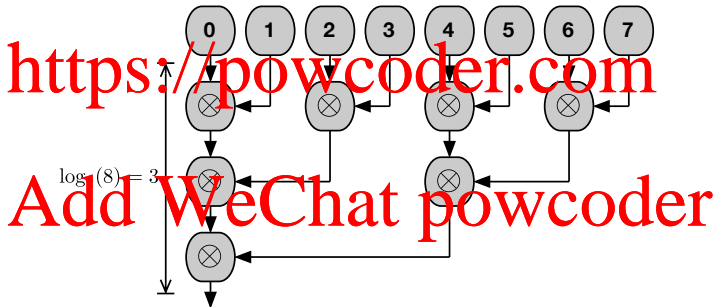
---

[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

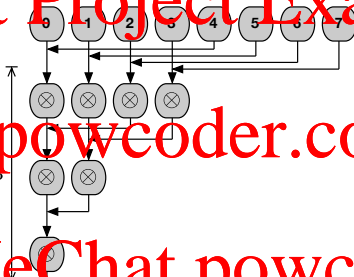Not all binary trees are valid for all binary operators $\otimes$.

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



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.
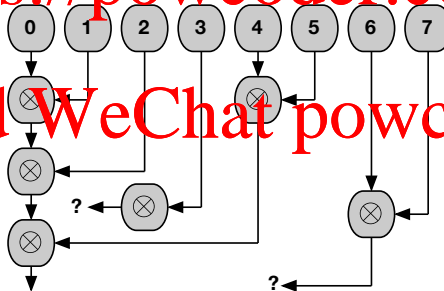- . . .

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 calculations have been **completed** before continuing to the **next** level.

This example, where units 3, 6 and 7 are delayed, would result in at best an incorrect calculation, and at worst **deadlock**:

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 at a specific point in code.

- Often called **barriers**
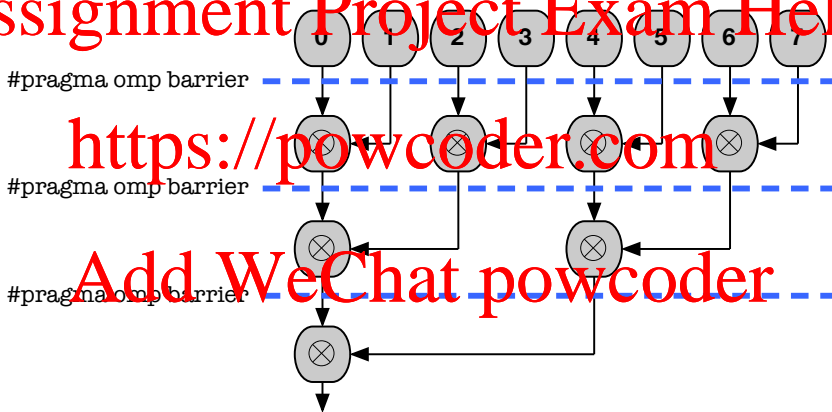
For instance, in OpenMP (in a parallel region):

```
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_Barrier( MPI_COMM_WORLD );
```

However, there is usually no need as the necessary synchronisation can be achieved using **blocking communication**.

- MPI_Send(), MPI_Recv() will not return until message has been sent or received.
- Provides the necessary synchronisation between pairs of 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

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

- Must be called by **all** ranks.

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

- Communication time $t_{comm} = \mathcal{O}(\log_2(p))$.
- Faster than the $\mathcal{O}(p)$ for a loop of send-and-receives.
- 'Inverted' in the case of `MPI_Bcast()` and `MPI_Scatter()`.

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

Lecture

Summary and next lecture

# Summary and next lecture

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

- Supported by most libraries, including OpenMPI and MPI.
- Typically implemented as a **binary tree**.
- Famous example was Google's **MapReduce**.
- In MPI, the necessary synchronisation provided by using **blocking communication**.

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