

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat University of Leeds edu_assist_pr

Lecture 4: Theory of parallel perf

Previous lecture

Assignment Project Exam Help

In the last lecture we started to look at solving problems in parallel:



<https://eduassistpro.github.io>

`#pragma omp parallel for`



Mandelbrot set which has a **ne**



Both **data parallel** problems (‘in
loop are independent.



Still difficult to achieve good performance for the Mandelbrot set.

This lecture

Assignment Project Exam Help

Now we will look at some general considerations for **parallel perf**

- <https://eduassistpro.github.io>
- Classic **models** for predicting parallel highlighting potential pitfalls.
- How these relate to **scaling**, i.e. the number of processors and the problem size.

Notation

Assignment Project Exam Help

For this lecture we will use the following notation:

S

<https://eduassistpro.github.io>

ce

t_s

t_p

Serial execution time

Parallel execution
time

'Optimal'

f

Serial fraction

Amdahl, Gustafson-Barsis

What we are trying to achieve

Assignment Project Exam Help

Assume a problem can be solved by a **serial** algorithm in a time t_s .

- We assume this is **optimal**, i.e. cannot be improved (in serial).

In pra

- <https://eduassistpro.github.io>
- May be known, but take too long to implement.

Usually consider the serial algorithm which is 'eq
parallel one.

- For instance, if developing a parallel bubblesort, would probably compare to **serial bubblesort** (rather than quicksort, mergesort, heapsort etc.).

Parallel acceleration

Assignment Project Exam Help

One way to improve on the serial execution time t_s is to implement a **parallel** solution on **parallel** hardware.



US



<https://eduassistpro.github.io>

Denote the (not necessarily optimal) parallel ex



Add WeChat edu_assist_pro



Measured in same units as t_s .

On 'as similar as possible' hardware.



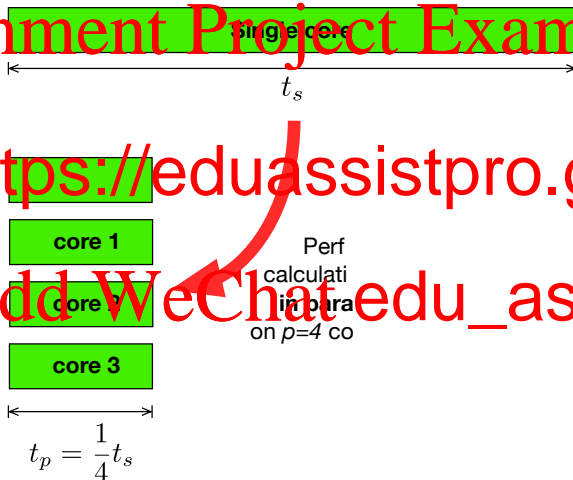
Sometimes known as the **wall clock** time, as it is what 'a clock on the wall' would measure.

Simultaneous calculations (ideally)

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](#)



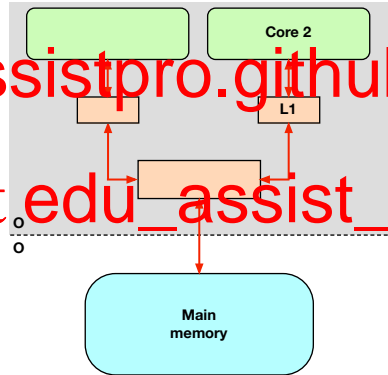
Multi-core memory cache

Assignment Project Exam Help

Recall from Lecture 2 that using multiple cores can make better use of **memory cache**.

Fewer

- C
c
required by another core.
- Depending on how data is arranged in memory and accessed, a parallel code may result in **fewer cache misses** overall than the equivalent serial algorithm.



Challenges to parallel performance

Assignment Project Exam Help

These potential benefits must be offset by the many challenges to achieving good parallel performance.

In Lecture 3: <https://eduassistpro.github.io>

- Hardware performance loss in maintaining **coherence** when two cores repeatedly write to the same **memory location** even though they never read the other core's data

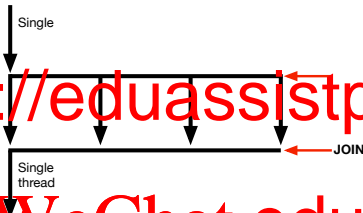
Add WeChat [edu_assist_pro](#)

Over the coming lectures we will see two important, general challenges: **synchronisation** and **load balancing**.

Synchronisation

Assignment Project Exam Help

In the **fork-join** construct from last lecture, multiple threads complete before the main thread continues.



This **joining** requires resources.

- Main thread may repeatedly **probe** worker thread status.
- Alternatively, workers may **signal** their completion to main.
- An example of **synchronisation**.

Load balancing

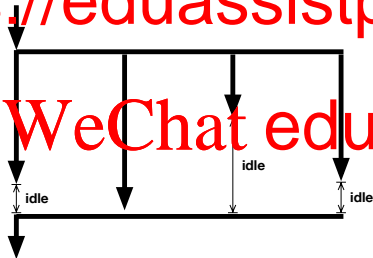
A related issue is **load balancing**:

- What if the worker threads did not all finish at the same time?



<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



This happens in the Mandelbrot set since each thread performs different numbers of calculations [*cf. last lecture; Lecture 13*].

Parallel overheads

Assignment Project Exam Help

Even if these challenges could be overcome, there are inevitable **overheads**. For example:

- <https://eduassistpro.github.io> threads
- **Communication** between threads/processes not present in the serial equivalent.
- **Computation** not present in serial, problem size between threads.

The impact may be small or large depending on parallel algorithm and hardware architecture.

Metrics for parallel performance

Assignment Project Exam Help

There are various measurements of parallel vs. serial performance.

The m

<https://eduassistpro.github.io>

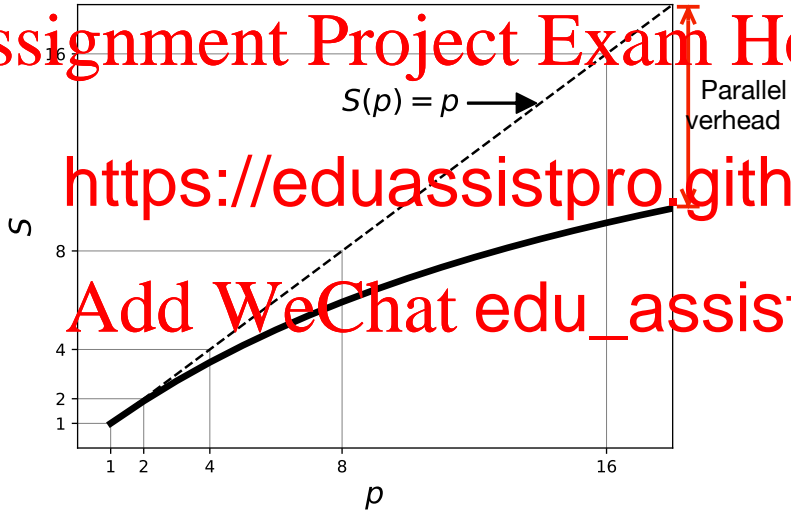
- If the parallel version was p times f

Add WeChat edu_assist_pro

$$t_p = \frac{1}{p} t_s \implies S = \frac{1}{\frac{1}{p} t_s} = p$$

- Rarely realised in practice due to **parallel overheads**.

Speedup example



Superlinear speedup

Assignment Project Exam Help

$S(p) > p$ is possible.

- Usually due to memory

C

This is

<https://eduassistpro.github.io>

Example (right): Benchmark computation of fluid dynamics algorithm.

Add WeChat edu_assist_pro

However, this is rare - most commonly see $S(p) < p$.

$S(p)=p$

From *Parallel Programming in OpenMP*, Chandra et al.
(Academic, 2001).

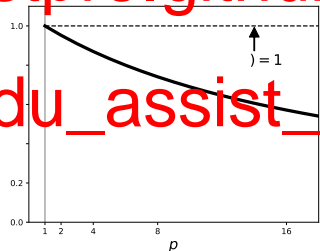
Efficiency

Another common parallel performance metric is the efficiency E :

$$\frac{t_s}{S}$$

<https://eduassistpro.github.io>

- corresponds to $E = 1$.
- Often expressed as a percentage:
 $E = 1 = 100\%$.
- Typically $E < 1$ due to parallel overheads.
- Superlinear speedup gives $E > 1$.



Models for parallel performance

Assignment Project Exam Help

Desirable to theoretically predict t_p for parallel algorithms.

- Select the 'best' without development and testing.
-

Chal

- Need to include e.g. memory cache
- Involve many unknown parameters required
- Would need re-calibration for new hardware

However, even **simple** models can predict **trends**.

- **Parallel scaling**, which refers to the variation with p .

Amdahl's law

Suppose a fraction f of t_s cannot be parallelised

$$t_s = ft_s + (1 - f)t_s$$

$$\frac{1}{t_p} = \frac{f}{ft_s} + \frac{(1-f)}{(1-f)t_s}$$

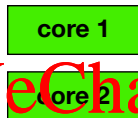
This is Amdahl's law¹ (1967).

For large p it predicts $S \leq \frac{1}{f}$ **regardless of p .**

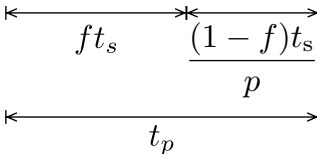
- e.g. $f = 0.2$, maximum speedup of 5, **even for $p=\infty$!**

¹Amdahl, *AFIPS Conference Proceedings* **30**, 483 (1967).

Schematic for Amdahl's law ($p = 3$)



divided between



Gustafson-Barsis law

Assignment Project Exam Help

However, Amdahl assumed t_s — and hence n — was fixed.

- Suppose instead n increases with p such that t_p is fixed.

<https://eduassistpro.github.io>

$$\implies S \leq f + p(1 - f)$$

Now $S \leq (1 - f)p$ for large p — no upper bound

Add WeChat [edu_assist_pro](#)

This is the **Gustafson-Barsis law**, or just **Gustafson's law**¹.

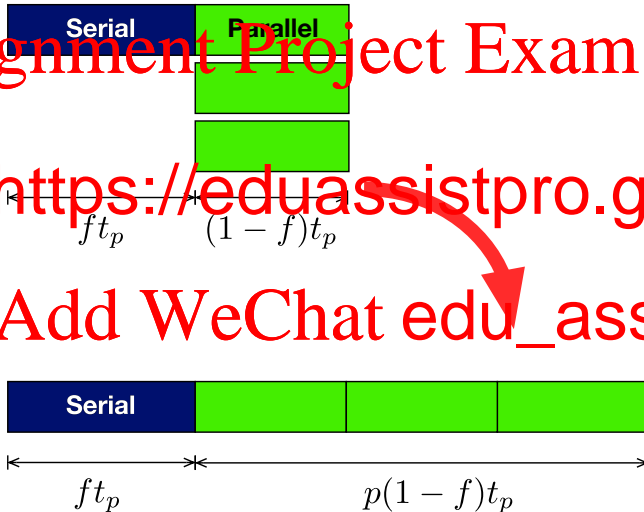
¹Gustafson, *Comm. ACM* **31**, 532 (1988).

Schematic for Gustafson-Barsis law ($p=3$)

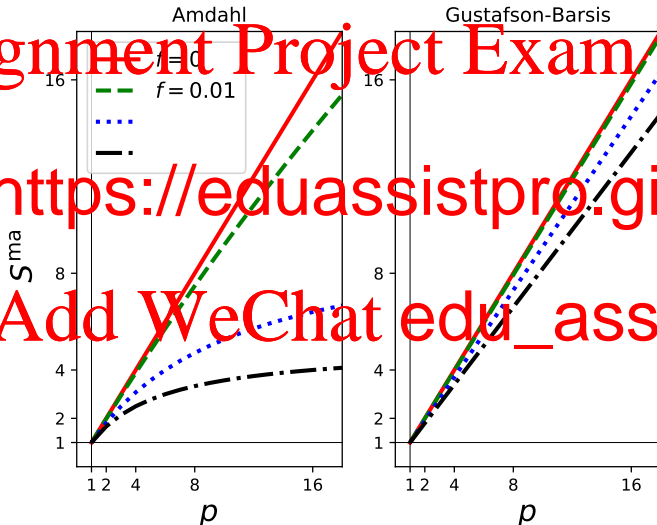
Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Amdahl *versus* Gustafson-Barsis



Weak *versus* strong scaling

The differences are encapsulated in **weak** *versus* **strong** scaling:

Strong scaling: Increasing p with n fixed.

-
-
-

<https://eduassistpro.github.io>

Weak scaling: Increasing n with p .

- **Gustafson-Barsis law.**
- Have freedom to vary n .
- e.g. higher resolution meshes for scientific/engineering applications; more/larger layers in neural networks.

Summary and next lecture

Assignment Project Exam Help

Today we have looked at parallel performance:

- Two common metrics: **speedup** and **efficiency**.



<https://eduassistpro.github.io>



Gustafson-Barsis law.

- Correspond to **strong** and **weak**

Add WeChat edu_assist_pr

Next time we will look more closely at **data dependencies** in parallel loops.