

# Assignment Project Exam Help

<https://eduassistpro.github.io>

University of Leeds  
Add WeChat edu\_assist\_pro

Lecture 2: Introduction to shared memory pa

## Previous lectures

# Assignment Project Exam Help

In the last introductory lecture we saw:

- S.
- <https://eduassistpro.github.io>
- Some general concepts:
  - **Concurrency** (more general than parallelism)
  - **Shared versus distributed memory**
    - Potential performance issues related to shared memory
  - **Flynn's taxonomy.**

## This lecture

# Assignment Project Exam Help

This l  
relev

- <https://eduassistpro.github.io>
- Processes *versus* threads and the threa
- Languages and frameworks suitable for it
- How to set up and run OpenMP.

Add WeChat edu\_assist\_pr

## Multi-core CPUs

# Assignment Project Exam Help

A **core** is a single processing unit that executes instructions.

- Has components that *fetch*, *decode* etc. instructions.

- 

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

As its name suggests, **multi-core** proc  
one such unit.

- **MIMD** in Flynn's taxonomy (sing
- Most common now are **dual core**, **quad core** and **octa core**.
- High-performance chips can have many more, e.g. SW26010 (used in China's Sunway TaihuLight supercomputer) has 260.

## Simultaneous multithreading

Some chips employ **simultaneous multithreading**<sup>1</sup>:

- Two (or more) threads run on the same core.
- If one thread stops execution (e.g. to wait for memory

Appears to require

- Performance improvements only 15%-

When interrogating a framework for the maximum available threads, you may get **more** than the number of cores.

<sup>1</sup>Known as **hyperthreading** on Intel chips.

<sup>2</sup>Rauber and Rünger, *Parallel Programming* 2<sup>nd</sup> ed. (Springer, 2013).

## The processor-memory gap

Memory access rates are increasing far slower than processor performance (taking into account number of cores):

- This is the **processor-memory gap**.

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

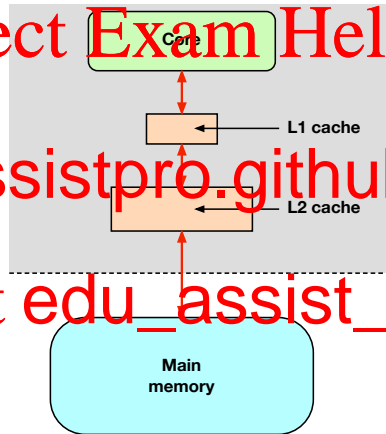
Relative increase in per

Year

Hennessy and Patterson, *Computer Architecture: A Quantitative Approach* (Morgan Kauffman, 2006).

## Single-core memory caches: A reminder

- Small, fast on-chip memory.
- A  
r  
(
- Subsequent accesses return the cache data (a **cache hit** - *fast*) or from main memory (a **cache miss** - *slow*).
- Multiple caches **levels** (e.g. L1, L2, L3) arranged **hierarchically**.



Schematic two-level cache

## Multi-core memory caches

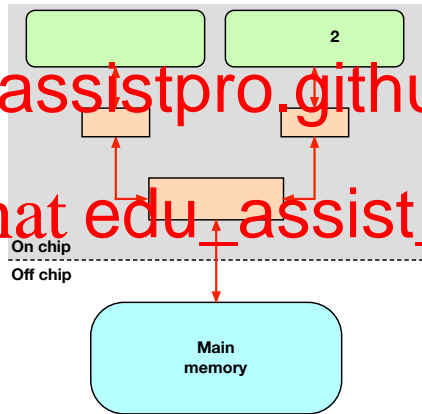
Different manufacturers choose different ways to incorporate caches into multi-core designs.

Often **hierarchy**

- 
- Share higher level caches.

For e.g. quad cores:

- L1 for each core.
- L2 for pairs.
- L3 for all cores.





## Cache coherency

# Assignment Project Exam Help

1 Core 1 reads an address  $x$ , resulting in a line in its L1.

2 Core 2 does the same, resulting in a line in its L1.

3

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

Maintaining consistent memory views for all

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io)

A common way to maintain **cache coherency** is **snooping**:

- The *cache controller* **detects writes** to caches, and updates higher-level caches.

## False sharing

# Assignment Project Exam Help

Maintaining cache coherency incurs a performance loss

- If two cores repeatedly write to the same memory location,

<https://eduassistpro.github.io>

How

locations **on the same cache line**, up

- *i.e.* hardware performance loss

This unnecessary cache coherency is known as **false sharing**

## Potential benefit of cache sharing

# Assignment Project Exam Help

It is also possible for multiple cores to benefit from shared caches:

① Core 1 reads an address  $x$  from main memory.

②

③

④

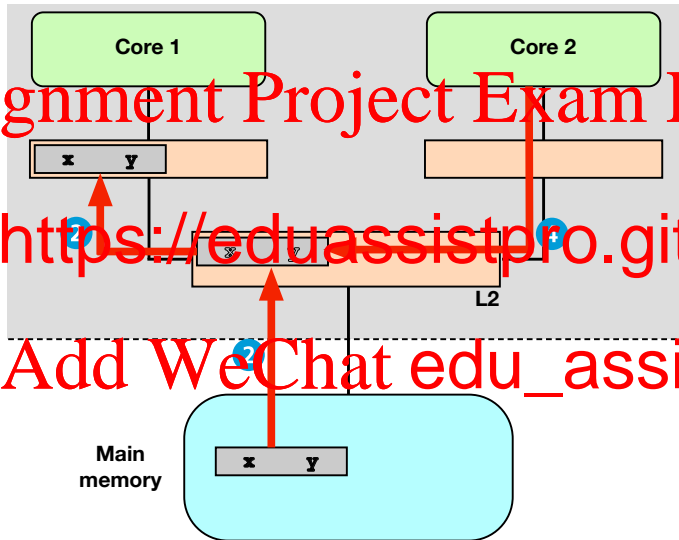
<https://eduassistpro.github.io>

need to

access main memory.

It is therefore possible for fewer accesses to main memory compared to the equivalent serial code.

This can result in parallel speed up **more than the number of cores** (known as superlinear speedup; cf. Lecture 4).



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

## Processes *versus* threads

Control flows can be either **processes** or **threads**:

# Assignment Project Exam Help

### Processes:



- **Expensive** to generate (large heap memory)

### Threads:

- Threads of one process **share** its address space.
- Implicit communication *via* this **shared memory**.
- **Cheap** to generate (no heap memory).

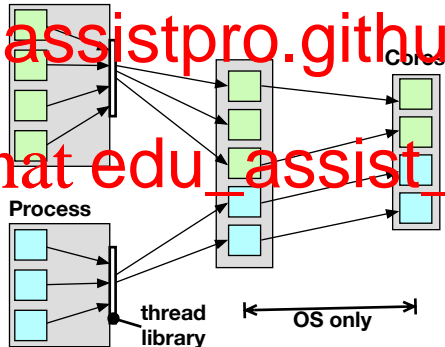
## Kernel *versus* user-level threads

The threads that execute on the core(s) are **kernel threads**.

- Only the OS has direct control over kernel threads.

Programmer generates **user-level threads**.

- Managed by a **thread library**.
- Mapped to kernel threads by the OS **scheduler**.



## Thread programming

**Assignment Project Exam Help**  
The choice of programming framework / API / library / etc. must be suitable for the architecture on which it will run.

- — — — — — lism.

It is pos

- Java supported threads early on, through `Runnable` interface.
- The C library `pthread` implement
- C++11 has language-level concurrency support.
- Python has a threading library, although need to work around its **global interpreter lock** to exploit multi-cores.

## Higher-level threading support

Higher-level options that do not require explicit thread control also exist to reduce development times



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

For C

/ in

*Structured Parallel Programming* (Mo

- Cilk Plus
- TBB (Threading Building Blocks)
- ArBB (Array Building Blocks).

- OpenCL, although primarily used for GPUs.

The first three are not (yet?) widely implemented in compilers.



## OpenMP

# Assignment Project Exam Help

For the SMP component of this module we will use OpenMP.



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

- Currently up to v5.2, although compilers may only support earlier versions

# Add WeChat edu\_assist\_pr

More information available from <http://www.openmp.org>

## Compiling C with OpenMP

For this module we will use gcc (GNU Compiler Collection)<sup>1</sup>.

- Compile with `-fopenmp`
- Must include `omp.h`

All pa

- <https://eduassistpro.github.io>
  - You will each get your own INDIVIDUAL ac
  - Full instructions for logging into your acco
  - A Linux C/S is provided along with the all librari
- this module:
- You can run jobs interactively on 2 cores while debugging (can still have more threads!)
  - You can run batch jobs on up to 16 cores via Slurm

---

<sup>1</sup>Easy to install on Macs with homebrew.

## helloWorld.c

Code on Minerva: helloWorld.c

Assignment Project Exam Help

```
1 #include <stdio.h>
2 #include <omp.h> // Run-time OpenMP library routines.
3
4 int main
5 {
6     // Threads
7     #pragma omp parallel
8     {
9         // Get this thread number, and the maximum.
10        int threadNum = omp_get_thread_num ();
11        int maxThreads = omp_get_max_threads();
12        // Simple message to stdout.
13        printf( "Hello from thread %i of %i!\n", threadNum,
14               maxThreads );
15    }
16    return 0;
17 }
```

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

## Compiling C: Reminder

# Assignment Project Exam Help

gcc -f

Opti

- -fopenmp tells compiler to expect Open
- -Wall turns on all warnings; recommen
- -o helloWorld is the executable nam
- helloWorld.c is the source code.
- Sometimes need e.g. -lm for the maths library.

```
#pragma omp parallel
```

# Assignment Project Exam Help

#pragma directives provide information beyond the language

- All OpenMP pragmas start: `#pragma omp ...`

Here

the ne

from { to }) **in parallel.**

- The code inside this scope is run by
- Outside of this scope there is only

This is why the `printf` statement is repeated multiple times, **even though it only appears once in code.**

We will look at this in more detail next time.

```
#include <omp.h>
```

# Assignment Project Exam Help

Include `omp.h` to use OpenMP runtime library routines:

```
int omp _ _ _
```

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

- May exceed apparent core number with **multithreading** (see earlier).

Add WeChat [edu\\_assist\\_pr](#)

```
int omp_get_thread_num():
```

- Returns the thread number **within the current scope**.
- $0 \leq \text{omp\_get\_thread\_num}() < \text{omp\_max\_thread\_num}()$

## Setting the number of threads

# Assignment Project Exam Help

If you don't want to use the default number of threads:

```
void o _ _ _
```

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

Alternatively, use shell **environment** **va**

- For bash: `export OMP_NUM_THREADS=4`
- Avoids the need to recompile.
- List all environment variables using `env`.
- To see all OMP variables: `env | grep OMP`

## Summary and next lecture

# Assignment Project Exam Help

Today we have started looking at shared memory parallelism (SMP):

- <https://eduassistpro.github.io>
- Various languages, frameworks *etc.* support SMP.
- OpenMP is commonly supported by C/C

Add WeChat edu\_assist\_pr

Next time we will look in more detail at what is actually going on at the thread level, for a more interesting example.