

Lecture 2-1 : OpenMP

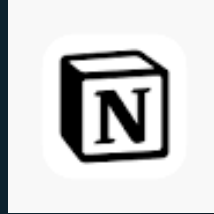
Lecturer : Atseng

Course Information

Join the Discord First !

All the class information call discuss in Discord

<https://discord.gg/wEGVXnT7Mh>



Course Contents

01 Thread intro

02 Pthread

03 OpenMP

04 Hands on time !

01 Core vs. Thread

What is Core ?

- > A core is a physical processing unit inside the cpu.
- > Each core can handle one task or process at a time.
- > More cores mean a CPU can perform multiple tasks simultaneously.
 - a quad-core CPU has 4 cores that work independently.

What is Thread ?

- > A thread is a **virtual component** that represent a task handled by a core.
- > Modern CPUs use ‘Hyper-Threading’(intel), ‘Simultaneous Multi-Threading’(AMD)
- > This increases efficiency and performance.

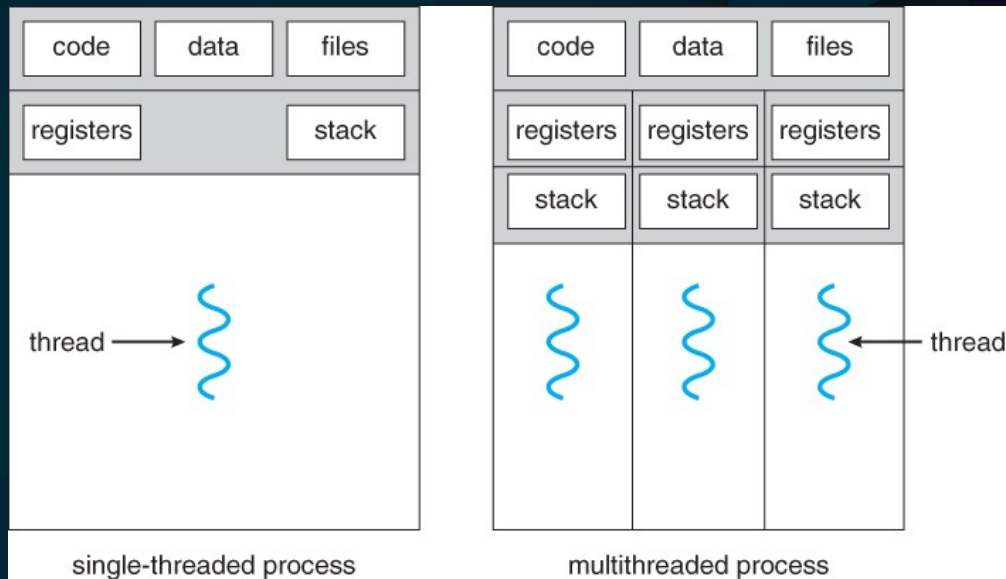
Core vs. Threads

- > **Cores are physical unit; Threads are virtual.**
- > **Threads rely on the cores to execute tasks.**
- > **A CPU with 4 cores and 8 threads can handle up to 8 tasks at the same time**
- > **More threads improve multitasking but don't always double performance.**

Process vs. Threads

> **Process:** Each process has its own
independent memory space

> **Threads :** Threads within the same process
share the same memory space and resources



Shared Memory Programming

Definition

- Multiple processes or threads execute within the same memory space.

Pros :

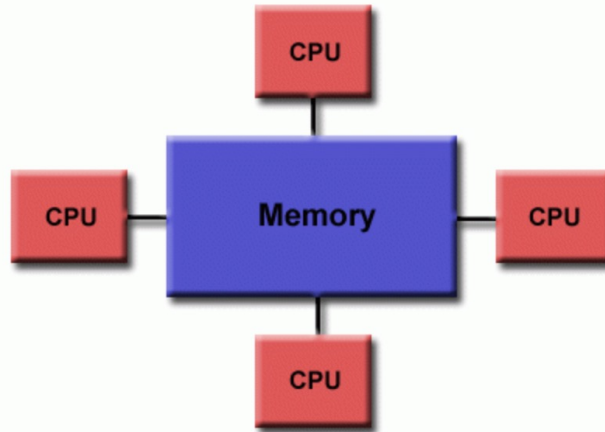
- Fast communication
- Low latency
- Efficient resource sharing

Cons :

- may have issue with Synchronization
- Deadlock
- Cache coherence

Uniform Memory Access (UMA)

- Most commonly represented today by *Symmetric Multiprocessor (SMP)* machines
- Identical processors
- Equal access and access times to memory
- Sometimes called CC-UMA - Cache Coherent UMA. Cache coherent means if one processor updates a location in shared memory, all the other processors know about the update. Cache coherency is accomplished at the hardware level.

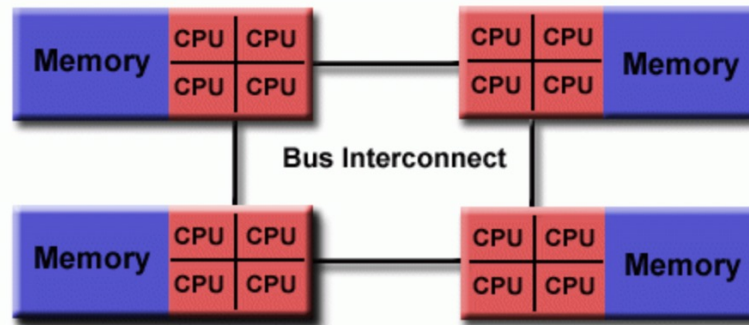


Uniform memory access

<https://hpc.llnl.gov/documentation/tutorials/introduction-parallel-computing-tutorial##SharedMemory>

Non-Uniform Memory Access (NUMA)

- Often made by physically linking two or more SMPs
- One SMP can directly access memory of another SMP
- Not all processors have equal access time to all memories
- Memory access across link is slower
- If cache coherency is maintained, then may also be called CC-NUMA - Cache Coherent NUMA



| Non-uniform memory access

<https://hpc.llnl.gov/documentation/tutorials/introduction-parallel-computing-tutorial##SharedMemory>

02 Pthread

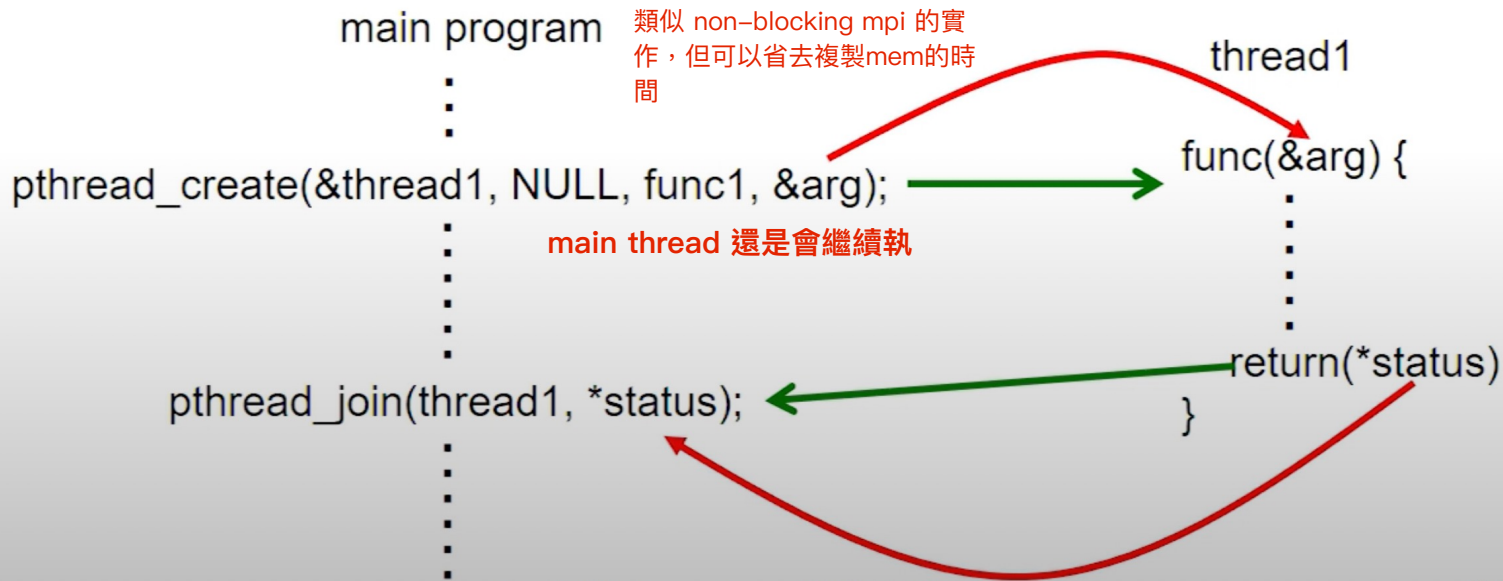
Pthread

> POSIX Thread = Pthread

- **POSIX** (Potable Operating System Interface) standard is specified for portability across Unix-like systems
 - Similar concept as MPI for message passing libraries
- **Pthread** is the implementation of POSIX standard for thread
 - Same relation between MPICH and MPI

<https://www.youtube.com/watch?v=IXXlmEGZ7ZU&list=PLS0SUwIYe8cxqw70UHOE5n4Lm-mXFXbZT&index=19>

Pthread



<https://www.youtube.com/watch?v=IXXImEGZ7ZU&list=PLS0SUwIYe8cxqw70UHOE5n4Lm-mXFXbZT&index=19>

Pthread_create

```
#include <pthread.h>
```

```
int pthread_create(pthread_t *restrict thread,  
                  const pthread_attr_t *restrict attr,  
                  typeof(void *(void *)) *start_routine,  
                  void *restrict arg);
```

https://man7.org/linux/man-pages/man3/pthread_create.3.html

```
pthread_t t1, t2;  
int t1_num = 1;  
int t2_num = 2;  
pthread_create(&t1, NULL, thread_func, &t1_num);
```

Pthread_join

```
#include <pthread.h>

int pthread_join(pthread_t thread, void **retval);
```

https://man7.org/linux/man-pages/man3/pthread_join.3.html

```
36 pthread_join(t1, NULL);
37 pthread_join(t2, NULL);
```

這裡有 blocking

Pthread code

```
19 int main() {
20     pthread_t t1, t2;
21     int t1_num = 1;
22     int t2_num = 2;
23
24     // create first thread
25     if (pthread_create(&t1, NULL, thread_func, &t1_num) != 0) {
26         perror("Failed to create thread 1");
27         return 1;
28     }
29
30     // create second thread
31     if (pthread_create(&t2, NULL, thread_func, &t2_num) != 0) {
32         perror("Failed to create thread 2");
33         return 1;
34     }
35
36     // wait for threads to finish
37     pthread_join(t1, NULL);
38     pthread_join(t2, NULL);
39
40     printf("All threads finished.\n");
41     return 0;
42 }
```

Pthread code

test.c / main()

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <pthread.h>
4
5  // thread function
6  void* thread_func(void* arg) {
7      int thread_num = *(int*)arg;
8      printf("Hello from thread %d\n", thread_num);
9
10     // thread work
11     for (int i = 0; i < 5; i++) {
12         printf("Thread %d is working... (i=%d)\n", thread_num, i);
13     }
14
15     // return
16     pthread_exit((void*)0);
17 }
```

```
19 int main() {
20     pthread_t t1, t2;
21     int t1_num = 1;
22     int t2_num = 2;
23
24     // create first thread
25     if (pthread_create(&t1, NULL, thread_func, &t1_num) != 0) {
26         printf("Failed to create thread 1");
27         return 1;
28     }
29
30     // create second thread
31     if (pthread_create(&t2, NULL, thread_func, &t2_num) != 0) {
32         printf("Failed to create thread 2");
33         return 1;
34     }
35
36     // wait for threads to finish
37     pthread_join(t1, NULL);
38     pthread_join(t2, NULL);
39
40     printf("threads finished.\n");
41     return 0;
42 }
```

Pthread

```
gcc -o thread_program yourfilename.c -lpthread
```

```
./thread_program
```

```
[atseng@TsengdeMacBook-Pro GPUtest % gcc -o thread_program thread.c -lpthread  
[atseng@TsengdeMacBook-Pro GPUtest % ./thread_program  
Hello from thread 1  
Thread 1 is working... (i=0)  
Thread 1 is working... (i=1)  
Thread 1 is working... (i=2)  
Thread 1 is working... (i=3)  
Thread 1 is working... (i=4)  
Hello from thread 2  
Thread 2 is working... (i=0)  
Thread 2 is working... (i=1)  
Thread 2 is working... (i=2)  
Thread 2 is working... (i=3)  
Thread 2 is working... (i=4)  
All threads finished.
```

03 OpenMP

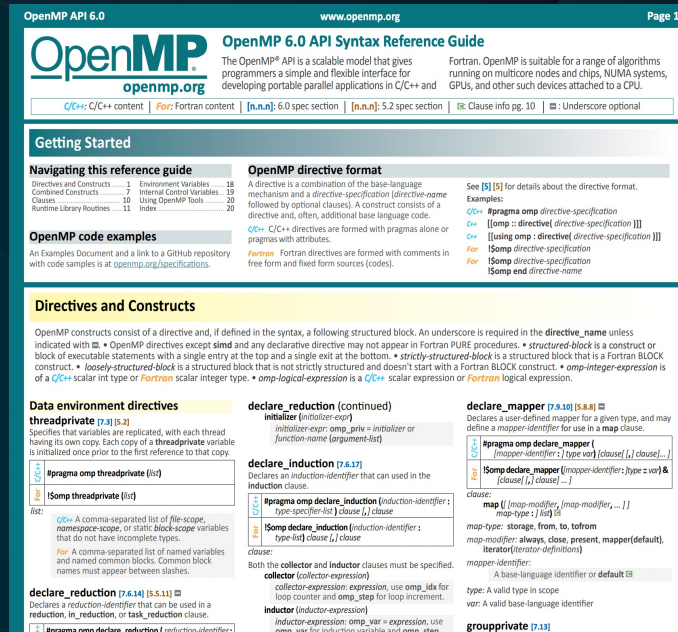
OpenMP

> Easy Parallelization

> Shared memory model

> Compiler Directives

<https://www.openmp.org/resources/refguides/>



<https://www.openmp.org/wp-content/uploads/OpenMP-RefGuide-6.0-OMP60SC24-web.pdf>

Compiler Directives

```
// Parallel region
#pragma omp parallel
{
    // Code here will be executed by multiple threads simultaneously
}

// Work distribution
#pragma omp parallel for
for (int i = 0; i < n; i++) {
    // Loop iterations will be automatically distributed among threads
}
```

Simple example (Vector Addition)

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

int main() {
    int n = 1000000;
    int a[n], b[n], c[n];

    // Initialize arrays
    for (int i = 0; i < n; i++) {
        a[i] = i;
        b[i] = 2 * i;
    }

    // Parallel computation of vector addition
    #pragma omp parallel for
    for (int i = 0; i < n; i++) {
        c[i] = a[i] + b[i];
    }

    printf("c[500000] = %d\n", c[500000]);
    return 0;
}
```

Data sharing and privacy

```
#pragma omp parallel for shared(a, b, c) private(i)
for (int i = 0; i < n; i++) {
    c[i] = a[i] + b[i];
}
```

- `shared` - Specifies variables shared among threads
- `private` - Specifies variables private to each thread

Synchronization Mechanisms

```
// Critical section
#pragma omp critical
{
    // Only one thread can execute this code at a time
    sum += local_sum;
}

// Barrier
#pragma omp barrier
// All threads will wait here until everyone arrives

// Atomic operation
#pragma omp atomic
counter++;
```

確保單一執行緒 (Thread) 對變數的存取是原子性 (Atomic) 的，避免多個執行緒同時修改變數導致不一致性。

僅適用於單一變數的簡單操作（只限於下面那行 code），如 +=, -=, *=, /=, ++, --，不適用於更複雜的程式碼區塊。

比 #pragma omp critical 更輕量，執行效率通常較高，因為 atomic 只保護單一變數操作，而 critical 可能會產生額外的同步開銷。

Task Parallelism

```
#pragma omp parallel
{
    #pragma omp single
    {
        #pragma omp task
        task1();

        #pragma omp task
        task2();
    }
    // Continue after all tasks complete
}
```

Thread Count Control

```
// Set number of threads
omp_set_num_threads(4);

// Get current thread count
int num_threads = omp_get_num_threads();

// Get current thread ID
int thread_id = omp_get_thread_num();
```

04 Hands on time

macOS (install llvm with homebrew)

```
/bin/bash -c "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/HEAD/install.sh)"
```

```
brew install llvm
```

1. Install Homebrew (if not already installed)

```
bash
```

[Copy](#)

```
/bin/bash -c "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/HEAD
```

2. Install LLVM using Homebrew

```
bash
```

[Copy](#)

```
brew install llvm
```

This will install the latest version of LLVM, which includes Clang with full OpenMP support.

macOS (Apple Silicon Macs : m1,m2,m3)

```
echo 'export PATH="/opt/homebrew/opt/llvm/bin:$PATH"' >> ~/.zshrc
```

```
echo 'export PATH="/usr/local/opt/llvm/bin:$PATH"' >> ~/.zshrc
```

```
source ~/.zshrc
```

```
source ~/.zshrc # or ~/.bash_profile if using bash
```

macOS (other Macs)

```
echo 'export PATH="/usr/local/opt/llvm/bin:$PATH"' >> ~/.zshrc
```

```
echo 'export PATH="/usr/local/opt/llvm/bin:$PATH"' >> ~/.zshrc
```

```
source ~/.zshrc
```

```
source ~/.zshrc # or ~/.bash_profile if using bash
```

Using Make file

```
CC = /opt/homebrew/opt/llvm/bin/clang++
CFLAGS = -fopenmp -O2

all: my_program

my_program: test_openmp.cpp
    $(CC) $(CFLAGS) -o my_program test_openmp.cpp

run: my_program
    OMP_DISPLAY_ENV=VERBOSE ./my_program

clean:
    rm -f my_program
```


Sample code

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

int main() {
    std::cout << "OpenMP Version: " << _OPENMP << std::endl;

    #pragma omp parallel
    {
        #pragma omp critical
        {
            std::cout << "Thread " << omp_get_thread_num()
                      << " of " << omp_get_num_threads() << std::endl;
        }
    }

    return 0;
}
```

Sample code2

```
1  #include <stdio.h>
2  #include <omp.h>
3  #include <time.h>
4
5  #define SIZE 100000000
6
7  int main() {
8      double a[SIZE], b[SIZE], c[SIZE];
9      int i;
10     double start_time, end_time;
11
12     // Initialize arrays
13     for (i = 0; i < SIZE; i++) {
14         a[i] = i * 1.0;
15         b[i] = i * 2.0;
16     }
17
18     // Record start time
19     start_time = omp_get_wtime();
20
21     // Parallelize the loop
22     #pragma omp parallel for
23     for (i = 0; i < SIZE; i++) {
24         c[i] = a[i] + b[i];
25     }
26
27     // Record end time
28     end_time = omp_get_wtime();
29
30     printf("Computation completed! Time: %f seconds\n", end_time - start_time);
31     printf("Verification: c[0]=%f, c[%d]=%f\n", c[0], SIZE-1, c[SIZE-1]);
32
33     return 0;
34 }
```



It's your turn!

Course Reference

<https://hpc.llnl.gov/documentation/tutorials/introduction-parallel-computing-tutorial##SharedMemory>

<https://www.youtube.com/watch?v=nE-xN4Bf8XI&list=PLLX-Q6B8xqZ8n8bwjGdzBJ25X2utwnoEG&index=2>

<https://pubs.opengroup.org/onlinepubs/7908799/xsh/pthread.h.html>