

# OpenMP Theory

OpenMP is a standardized API that enables you to add parallelism to your C, C++, or Fortran programs using simple compiler directives. It's designed for shared-memory systems (such as multi-core desktops or servers), allowing you to write code that can be compiled and run both in serial and in parallel without changing the original algorithm. Here's a detailed look at the theory behind OpenMP and practical guidelines for using it in C++.

## What Is OpenMP?

- **Definition:**

OpenMP (Open Multi-Processing) is an API that uses compiler directives (pragmas), runtime library routines, and environment variables to control parallel execution. It was first introduced in 1997 and is now supported by most major compilers (GCC, Clang, Intel, MSVC, etc.)

- **Parallel Model:**

OpenMP follows the *fork-join* model:

- **Fork:** The master thread (running serially) encounters a parallel region and “forks” a team of threads.
- **Join:** After executing the parallel code, threads synchronize and “join” back to the master thread.

- **Shared Memory:**

In OpenMP, threads share the same memory space. This simplifies programming because data does not need to be explicitly passed between threads, although it requires careful management to avoid race conditions.

## Core Concepts and Constructs

### 1. Parallel Regions

A parallel region is the block of code that runs concurrently on multiple threads. In C++, you mark such a region with:

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

int main() {
    #pragma omp parallel
    {
        int id = omp_get_thread_num();
        std::cout << "Hello from thread " << id << "\n";
    }
    return 0;
}
```

- **omp\_get\_thread\_num():** Returns the unique ID of the current thread.
- **omp\_get\_num\_threads():** Returns the total number of threads executing the region.

## 2. Work-Sharing Constructs

To distribute work among threads, OpenMP provides constructs that split loops or independent code sections:

- **omp for:** Divides iterations of a loop among the available threads.

```
#pragma omp parallel for
for (int i = 0; i < 100; i++) {
    // Each thread processes a chunk of iterations.
}
```

- **omp sections:** Used to specify distinct blocks of code that can run in parallel.

```
#pragma omp parallel sections
{
    #pragma omp section
    { /* Code block 1 */ }
    #pragma omp section
    { /* Code block 2 */ }
}
```

- **omp single/master:** These directives ensure that a particular block is executed by only one thread (either any one or specifically the master thread).

## 3. Data Environment Management

Since all threads share memory, you must control which variables are shared or private:

- **shared:** Variables are accessible by all threads (default for global variables).
- **private:** Each thread gets its own instance of the variable (uninitialized by default).
- **firstprivate:** Similar to private, but each thread's instance is initialized with the value from before entering the parallel region.
- **lastprivate:** After a loop, the variable's value from the last iteration (in sequential order) is copied back to the original variable.

Example:

```
int sum = 0;
#pragma omp parallel for reduction(+:sum)
for (int i = 0; i < 100; i++) {
    sum += i;
}
std::cout << "Sum is " << sum << "\n";
```

Here, the `reduction` clause makes each thread maintain a private copy of `sum` and then combines them using addition.

## 4. Synchronization and Mutual Exclusion

Parallel programs often require synchronization to avoid race conditions:

- **critical:** Ensures that a code block is executed by only one thread at a time.

```
#pragma omp critical
{
    // Critical section: only one thread at a time can execute this.
}
```

- **atomic:** Provides a lighter-weight option to perform specific memory updates atomically.
- **barrier:** Forces all threads to wait until each has reached the barrier point.

## 5. Scheduling

When you parallelize loops, you can control how iterations are distributed using the `schedule` clause:

- **static:** Divides iterations evenly among threads.
- **dynamic:** Threads request new chunks of iterations as they finish previous ones, useful when iterations have varying workloads.
- **guided:** Similar to dynamic but with decreasing chunk sizes to reduce overhead as the loop progresses.

Example:

```
#pragma omp parallel for schedule(dynamic, 10)
for (int i = 0; i < 1000; i++) {
    // Process iteration i.
}
```

# How to Use OpenMP in C++

## 1. Enabling OpenMP

To use OpenMP, include the header and compile with the OpenMP flag:

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

- **GCC/Clang:** Compile with `-fopenmp`  
Example: `g++ -fopenmp -O2 my_program.cpp -o my_program`
- **MSVC:** Enable OpenMP support in project settings or compile with `/openmp`.

## 2. Practical Steps

- **Identify Parallel Regions:**  
Look for loops or independent sections of code that can be executed in parallel.
- **Apply Directives:**  
Add the appropriate `#pragma omp` directives to create parallel regions, work-sharing constructs, and data-sharing clauses.

- **Manage Data Sharing:**  
Carefully decide which variables are shared (common to all threads) and which should be private to avoid race conditions.
- **Synchronize When Needed:**  
Use synchronization constructs like `critical`, `atomic`, and `barrier` where multiple threads might update the same data concurrently.
- **Tune Thread Count:**  
You can set the number of threads using the runtime function `omp_set_num_threads( )` or by setting the `OMP_NUM_THREADS` environment variable.
- **Test and Debug:**  
Parallel programming can introduce subtle bugs like race conditions or deadlocks. Use OpenMP's runtime functions (like `omp_get_thread_num( )`) to help debug and verify that your code is running correctly in parallel.

## Example: Combining It All

Below is a simple example that demonstrates a parallel loop with a reduction in C++:

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

int main() {
    const int N = 100;
    int sum = 0;

    // Parallelize the loop; each thread gets its private copy of sum, then
    // they are combined.
    #pragma omp parallel for reduction(+:sum)
    for (int i = 0; i < N; i++) {
        sum += i;
    }

    std::cout << "Sum of numbers 0 to " << N-1 << " is " << sum << "\n";
    return 0;
}
```

Compile with:

```
g++ -fopenmp -O2 example.cpp -o example
```

## Summary

- **Theory:** OpenMP is based on a fork-join model where the master thread spawns multiple threads for parallel execution. It simplifies parallel programming by providing high-level constructs for dividing work, managing data, and synchronizing threads.
- **Usage in C++:** You simply add `#pragma omp` directives to your code, manage data sharing with clauses like `shared`, `private`, and `reduction`, and compile your program with the appropriate flag (e.g., `-fopenmp`). This allows you to transform sequential loops into parallel loops with minimal changes to your code.