



UNIVERSIDAD CATÓLICA SAN PABLO

Ejercicios del Chapter 4: C++11 MULTITHREADING

Computer Science — Parallel and Distributed Computing

Harold Alejandro Villanueva Borda

1. Ejercicio 2

Considere dos matrices cuadradas $A, B \in \mathbb{R}^{n \times n}$ cada una de ancho $n = 1024$. A continuación queremos calcular su suma $C = A + B$ mediante la suma indexada de entradas $C_{ij} = A_{ij} + B_{ij}$ para todo $i, j \in \{0, \dots, n - 1\}$. Implemente un programa paralelo que utilice p subprocesos para calcular las entradas de C . Experimente con diferentes patrones de distribución de subprocesos estáticos.

1. Configuración Inicial:

- Definimos el tamaño de las matrices $n = 1024$.
- Inicializamos las matrices A y B con valores aleatorios o predefinidos.
- Definimos la matriz C que almacenará el resultado.

2. Distribución de Trabajo:

- Dividimos el trabajo entre p hilos. Cada hilo procesará un bloque de filas de las matrices A y B y calculará las filas correspondientes de C .

3. Creación y Ejecución de Hilos:

- Creamos y lanzamos los hilos para realizar la suma de las matrices en paralelo.

4. Unión de Hilos:

- Esperamos a que todos los hilos completen su trabajo y se unan al hilo principal.

5. Código en C++11:

```
1 #include <iostream>
2 #include <vector>
3 #include <thread>
4 #include <random> // Para std::random_device y std::mt19937
5
6 const int n = 1024; // Tamaño de las matrices
7 const int p = 8;    // Número de hilos
8
9 // Función que cada hilo ejecutará para sumar las matrices
10 void sum_matrices(int thread_id, const std::vector<std::vector<int>>& A, const std::vector<
    std::vector<int>>& B, std::vector<std::vector<int>>& C) {
11     int rows_per_thread = n / p;
12     int start_row = thread_id * rows_per_thread;
```

```

13     int end_row = (thread_id + 1) * rows_per_thread;
14
15     for (int i = start_row; i < end_row; ++i) {
16         for (int j = 0; j < n; ++j) {
17             C[i][j] = A[i][j] + B[i][j];
18         }
19     }
20 }
21
22 int main() {
23     // Inicializaci n de matrices A y B con valores aleatorios
24     std::vector<std::vector<int>>> A(n, std::vector<int>(n));
25     std::vector<std::vector<int>>> B(n, std::vector<int>(n));
26     std::vector<std::vector<int>>> C(n, std::vector<int>(n));
27
28     std::random_device rd; // Generador de n meros aleatorios
29     std::mt19937 gen(rd()); // Motor de generaci n de n meros aleatorios
30     std::uniform_int_distribution<> dis(0, 99); // Distribuci n uniform
31
32     for (int i = 0; i < n; ++i) {
33         for (int j = 0; j < n; ++j) {
34             A[i][j] = dis(gen);
35             B[i][j] = dis(gen);
36         }
37     }
38
39     // Vector para almacenar los hilos
40     std::vector<std::thread> threads;
41
42     // Creaci n y lanzamiento de hilos
43     for (int i = 0; i < p; ++i) {
44         threads.emplace_back(sum_matrices, i, std::ref(A), std::ref(B), std::ref(C));
45     }
46
47     // Uni n de hilos
48     for (int i = 0; i < p; ++i) {
49         threads[i].join();
50     }
51
52     // mostramos una parte de la matriz C
53     std::cout << "C[0][0] = " << C[0][0] << std::endl;
54     std::cout << "C[n-1][n-1] = " << C[n-1][n-1] << std::endl;
55
56     return 0;
57 }

```

Implementation

Explicaci3n del C3digo

- **Inicializaci3n de Matrices:** Las matrices A y B se inicializan con valores aleatorios usando `std::rand()`. La matriz C se inicializa para almacenar los resultados de la suma.
- **Funci3n de Suma:** La funci3n `sum_matrices` toma un `thread_id` para identificar qu3 parte de las filas debe procesar. Cada hilo calcula su porci3n correspondiente de la suma.

- **Lanzamiento y Unión de Hilos:** Se crean p hilos y se les asigna la tarea de sumar las matrices. Luego, el hilo principal espera a que todos los hilos terminen su trabajo con `join()`.

Consideraciones sobre la Distribución de Hilos

- **Distribución Estática:** En esta implementación, cada hilo procesa un bloque de filas de tamaño fijo. Esto se conoce como distribución estática y es fácil de implementar y eficiente si las filas tienen un tiempo de procesamiento similar.
- **Balanceo de Carga:** Si algunas filas fueran más costosas de procesar que otras, una distribución dinámica (donde los hilos solicitan nuevas filas para procesar después de terminar las actuales) podría ser más eficiente. Esto evitaría que algunos hilos terminen mucho antes que otros.

2. Ejercicio 4

Considera el siguiente fragmento de código:

```
#include <iostream>
#include <thread>
int main () {
    auto child = [] () -> void {
        std::cout << "child" << std::endl;
    };
    std::thread thread(child);
    thread.detach();

    std::cout << "parent" << std::endl;
}
```

En la mayoría de los casos, la salida del programa es exclusivamente “parent” ya que llegamos al final de la función principal antes de la declaración de impresión en el cierre del hijo. Sin embargo, el uso de `thread.join()`; en lugar de `thread.detach()`; resulta en la salida “child” seguida de “parent”. Ahora supón que tuviéramos que implementar nuestro propio método join usando variables de condición. Procederíamos de la siguiente manera:

- Introducimos una variable booleana global `done = false`. Además, necesitamos un mutex `m` y una variable de condición `c`.
- Después de la declaración de impresión en child, establecemos `done = true` en un alcance bloqueado y luego notificamos la variable de condición `c`.
- Nuestro método join personalizado realiza una espera condicional en un alcance bloqueado siempre que `done == false`.

Implementa el método join personalizado descrito. Prueba diferentes implementaciones: experimenta con el bloqueo explícito de mutexes y bloqueos de alcance.

- Explicación del comportamiento:

- En la mayoría de los casos, la salida del programa es exclusivamente “parent” ya que se llega al final de la función `main` antes de que se ejecute la declaración de impresión en el cierre `child`.

- El uso de `thread.join()`; en lugar de `thread.detach()`; resulta en la salida `child` seguida de `parent`.

2. Implementación del Método `join` Personalizado:

a. Introducción de variables globales:

```
#include <iostream>
#include <thread>
#include <mutex>
#include <condition_variable>

bool done = false;
std::mutex m;
std::condition_variable c;
```

- `done`: Variable booleana global que indica si el hilo hijo ha terminado su ejecución.
- `m`: Mutex global para sincronizar el acceso a `done`.
- `c`: Variable de condición global para notificar cuando el hilo hijo ha terminado.

b. Modificación de la función del hilo hijo:

```
auto child = [] () -> void {
    std::cout << "child" << std::endl;
    std::unique_lock<std::mutex> lock(m);
    done = true;
    c.notify_one();
};
```

- Después de imprimir `child`, se adquiere un bloqueo único (`std::unique_lock`) en el mutex `m`.
- Se establece `done` a `true` dentro del alcance del bloqueo.
- Se notifica a la variable de condición `c` que el hilo ha terminado.

c. Implementación del método `join` personalizado:

```
void custom_join() {
    std::unique_lock<std::mutex> lock(m);
    c.wait(lock, []{return done;});
}
```

- `custom_join` adquiere un bloqueo único en el mutex `m`.
- La función `c.wait` espera a que `done` sea `true`, liberando el bloqueo temporalmente y volviéndolo a adquirir cuando se despierta.

d. Modificación de la función `main`:

```

int main() {
    auto child = [] () -> void {
        std::cout << "child" << std::endl;
        std::unique_lock<std::mutex> lock(m);
        done = true;
        c.notify_one();
    };
    std::thread thread(child);
    custom_join();
    std::cout << "parent" << std::endl;
}

```

- Se inicia el hilo `child` y luego se llama a `custom_join` para esperar a que el hilo `child` termine antes de continuar.
- Después de que `custom_join` retorna, se imprime "parent".

3. Prueba de Diferentes Implementaciones:

Con bloqueo explícito:

```

void custom_join() {
    std::unique_lock<std::mutex> lock(m);
    while (!done) {
        c.wait(lock);
    }
}

```

- Aquí, `custom_join` usa un ciclo `while` para esperar a que `done` sea `true`.

4. Con bloqueos de alcance:

```

void custom_join() {
    std::lock_guard<std::mutex> lock(m);
    c.wait(m, []{return done;});
}

```

- Esta versión no es válida ya que `std::condition_variable::wait` requiere un `std::unique_lock`.

5. Código Completo con Implementación del join Personalizado:

```

1 #include <iostream>
2 #include <thread>
3 #include <mutex>
4 #include <condition_variable>
5
6 bool done = false;
7 std::mutex m;
8 std::condition_variable c;

```

```

9
10 //bloqueo join personalizado
11 void custom_join1() {
12     std::unique_lock<std::mutex> lock(m);
13     c.wait(lock, []{ return done; });
14 }
15
16 //bloqueo explicito
17 void custom_join2() {
18     std::unique_lock<std::mutex> lock(m);
19     while (!done) {
20         c.wait(lock);
21     }
22 }
23
24 /*
25 //bloqueo de alcance, forma invalida
26 void custom_join3() {
27     std::lock_guard<std::mutex> lock(m);
28     c.wait(m, []{ return done; });
29 }
30 /**/
31
32 int main () {
33     auto child = [] () -> void {
34         std::cout << "child" << std::endl;
35
36         std::unique_lock<std::mutex> lock(m);
37         done = true;
38         c.notify_one();
39     };
40
41     std::thread thread(child);
42     custom_join1();
43
44     std::cout << "parent" << std::endl;
45 }

```

Implementation

El método `join` personalizado garantiza que el hilo principal espere hasta que el hilo hijo haya terminado su ejecución, imprimiendo siempre `child` seguido de `parent`.

3. Ejercicio 5

Revisa el ejercicio mencionado anteriormente. Ahora implementa un método `join` personalizado utilizando one-shot synchronization como se describe en la Sección 4.5.

Paso 1: Implementar un método join personalizado utilizando la sincronización de una sola vez

Paso 2: Creación de la variable global y los componentes necesarios: Primero, necesitamos una variable global `done`, un `std::promise` y un `std::future` para manejar la sincronización de una sola vez.

```
#include <iostream>
#include <thread>
#include <future>
#include <mutex>
#include <condition_variable>

bool done = false;
std::promise<void> done_promise;
std::future<void> done_future = done_promise.get_future();
std::mutex m;
std::condition_variable c;
```

Paso 3: Definición de la función del hilo hijo

En la función del hilo hijo, después de imprimir `child`, establecemos `done = true` y establecemos el valor de la promesa.

```
auto child = []() -> void {
    std::cout << "child" << std::endl;
    {
        std::lock_guard<std::mutex> lock(m);
        done = true;
    }
    done_promise.set_value();
};
```

Paso 4: Implementación del método join personalizado

El método join personalizado usará `done_future` para esperar hasta que el hilo hijo complete su ejecución.

```
void custom_join() {
    done_future.wait();
}
```

Paso 5: Implementación de la función main

En la función main, creamos el hilo hijo y luego llamamos a nuestro método `custom_join`.

```
int main() {
    std::thread thread(child);
    custom_join();
    std::cout << "parent" << std::endl;
    if (thread.joinable()) {
        thread.join();
    }
    return 0;
}
```

```
1 #include <iostream>
2 #include <thread>
3 #include <future>
4 #include <mutex>
5 #include <condition_variable>
6
7 bool done = false;
8 std::promise<void> done_promise;
9 std::future<void> done_future = done_promise.get_future();
10 std::mutex m;
11 std::condition_variable c;
12
13 auto child = []() -> void {
14     std::cout << "child" << std::endl;
15
16     {
17         std::lock_guard<std::mutex> lock(m);
18         done = true;
19     }
20
21     done_promise.set_value();
22 };
23
24 void custom_join() {
25     done_future.wait();
26 }
27
28 int main() {
29     std::thread thread(child);
30
31     custom_join();
32
33     std::cout << "parent" << std::endl;
34
35     if (thread.joinable()) {
36         thread.join();
37     }
38
39     return 0;
40 }
```

Implementation

Explicación:

1. Definición de variables globales:

- **done**: Booleano que indica si el hilo hijo ha terminado.
- **done_promise** y **done_future**: Usados para la sincronización de una sola vez.
- **m** y **c**: Mutex y variable de condición, aunque no son estrictamente necesarios para esta versión usando **promise** y **future**.

2. Función del hilo hijo:

- Imprime `child`.
- Establece `done = true` dentro de un `lock_guard` para garantizar la sincronización.
- Establece el valor de la promesa usando `done_promise.set_value()`, lo que satisface el futuro asociado.

3. Método `join` personalizado:

- Usa `done_future.wait()` para esperar hasta que el hilo hijo complete su ejecución.

4. Función `main`:

- Crea y lanza el hilo hijo.
- Llama a `custom_join` para esperar a que el hilo hijo termine.
- Imprime `"parent"`.
- Verifica si el hilo es joinable y llama a `thread.join()` para limpiar.

Esta implementación garantiza que siempre se imprimirá `child` seguido de `"parent"`, independientemente del orden de ejecución de los hilos.

4. Ejercicio 7

Recuerda la distribución cíclica de bloques programada dinámicamente para el cálculo de la matriz de distancias de todos los pares presentada en la Sección 4.4. El esquema propuesto selecciona dinámicamente bloques de tamaño fijo uno tras otro comenzando en la fila cero. Invierte el esquema de indexación de tal manera que el programa comienza con los bloques en la parte inferior de la matriz que contienen las filas más largas. ¿Observas una mejora en el tiempo de ejecución? Explica el resultado.

4.1. Declaración de Variables Globales:

```
1 index_t global_lower = 0;
2 std::mutex mutex;
```

4.2. Función Lambda para el Cálculo Dinámico:

La función lambda `dynamic_block_cyclic` se encargará de distribuir las tareas entre los hilos.

4.3. Detalles de la función `dynamic_all_pairs_rev`:

4.3.1. Parámetros

- `mnist`: Un vector que contiene los datos de las imágenes.
- `all_pair`: Un vector que almacenará las distancias de pares calculadas.
- `rows`: Número de filas (imágenes) en la matriz `mnist`.
- `cols`: Número de columnas (características o píxeles por imagen) en la matriz `mnist`.
- `num_threads`: Número de hilos a utilizar para el cálculo paralelo. Por defecto es 64.
- `chunk_size`: Tamaño del bloque de datos que cada hilo procesará. Por defecto es `64/sizeof(value_t)`.

4.3.2. Variables Globales y Bloques Cíclicos Dinámicos

- **global_lower**: Una variable compartida que mantiene el índice inferior actual que necesita ser procesado. Se utiliza para coordinar los hilos.
- **mutex**: Un mutex para proteger el acceso a **global_lower** y asegurar que solo un hilo a la vez lo modifique.

4.3.3. Lambda `dynamic_block_cyclic`

La función lambda `dynamic_block_cyclic` es la unidad de trabajo que cada hilo ejecutará. Aquí están los pasos detallados:

- **Inicialización**: **lower** se inicializa a 0, representando el índice inferior del bloque que el hilo procesará.
- **Bucle de Trabajo**: Mientras **lower** sea menor que **rows**, el hilo continuará procesando.
- **Sección Crítica**: Se utiliza un `std::lock_guard` para bloquear el acceso a **global_lower**. Dentro de la sección crítica:
 - **lower** se actualiza con el valor de **global_lower**.
 - **global_lower** se incrementa por **chunk_size** para asignar el siguiente bloque.
- **Cálculo de Límites del Bloque**: Se determina el límite superior del bloque como el menor valor entre **lower** + **chunk_size** y **rows**. Se calcula **LOWER** como **upper** - **chunk_size** si **upper** es mayor o igual a **chunk_size**, de lo contrario se asigna a 0.
- **Cálculo de Distancias**: Para cada par de imágenes (*i*, *I*) dentro del bloque:
 - Se calcula la distancia Euclidiana cuadrada entre las imágenes *i* e *I*.
 - El resultado se almacena en **all_pair** tanto en la posición $i \times \text{rows} + I$ como en $I \times \text{rows} + i$.

4.3.4. Creación y Unión de Hilos

- Se crea un vector de hilos.
- Se lanzan **num_threads** hilos, cada uno ejecutando `dynamic_block_cyclic`.
- Se espera a que todos los hilos terminen su ejecución con `thread.join()`.

4.3.5. Diferencias con `dynamic_all_pairs`

La principal diferencia entre `dynamic_all_pairs_rev` y `dynamic_all_pairs` es la forma en que se calculan los límites **lower** y **upper** para el procesamiento de bloques. En `dynamic_all_pairs_rev`, los bloques se procesan en orden inverso, desde la parte superior de la matriz hacia abajo, mientras que en `dynamic_all_pairs`, los bloques se procesan en orden natural desde la parte inferior hacia arriba. Esta variación puede ayudar en ciertos contextos de balance de carga y eficiencia de caché.

4.4. Implementación

```

1 void dynamic_all_pairs_rev(
2     std::vector<value_t>& mnist,
3     std::vector<value_t>& all_pair,
4     index_t rows,
5     index_t cols,
6     index_t num_threads=64,
7     index_t chunk_size=64/sizeof(value_t)) {
8
9     // declare mutex and current lower index
10    index_t global_lower = 0;
11
12    auto dynamic_block_cyclic = [&] (const index_t& id ) -> void {
13
14        // assume we have not done anything
15        index_t lower = 0;
16
17        // while there are still rows to compute
18        while (lower < rows) {
19
20            // update lower row with global lower row
21            {
22                std::lock_guard<std::mutex> lock_guard(mutex);
23                lower = global_lower;
24                global_lower += chunk_size;
25            }
26
27            // compute the upper border of the block (exclusive)
28            const index_t upper = rows >= lower ? rows-lower : 0;
29            const index_t LOWER = upper >= chunk_size ? upper-chunk_size : 0;
30
31            // for all entries below the diagonal (i'=I)
32            for (index_t i = LOWER; i < upper; i++) {
33                for (index_t I = 0; I <= i; I++) {
34
35                    // compute squared Euclidean distance
36                    value_t accum = value_t(0);
37                    for (index_t j = 0; j < cols; j++) {
38                        value_t residue = mnist[i*cols+j]
39                            - mnist[I*cols+j];
40                        accum += residue * residue;
41                    }
42
43                    // write Delta[i,i'] = Delta[i',i]
44                    all_pair[i*rows+I] =
45                    all_pair[I*rows+i] = accum;
46                }
47            }
48        }
49    };
50
51    // business as usual
52    std::vector<std::thread> threads;
53
54    for (index_t id = 0; id < num_threads; id++)
55        threads.emplace_back(dynamic_block_cyclic, id);

```

```
56
57     for (auto& thread : threads)
58         thread.join();
59 }
```

Implementation

4.5. Explicación de los resultados

Los resultados obtenidos muestran que la función `dynamic_all_pairs_rev` es ligeramente más rápida que la función `dynamic_all_pairs` original. Aquí está el análisis detallado:

Resultados

`dynamic_all_pairs`:

- Carga de datos desde el disco: 0.023642s
- Cálculo de distancias: 23.707s
- Guardado de resultados en el disco: 0.0070773s

`dynamic_all_pairs_rev`:

- Carga de datos desde el disco: 0.0227069s
- Cálculo de distancias: 22.859s
- Guardado de resultados en el disco: 0.0080092s

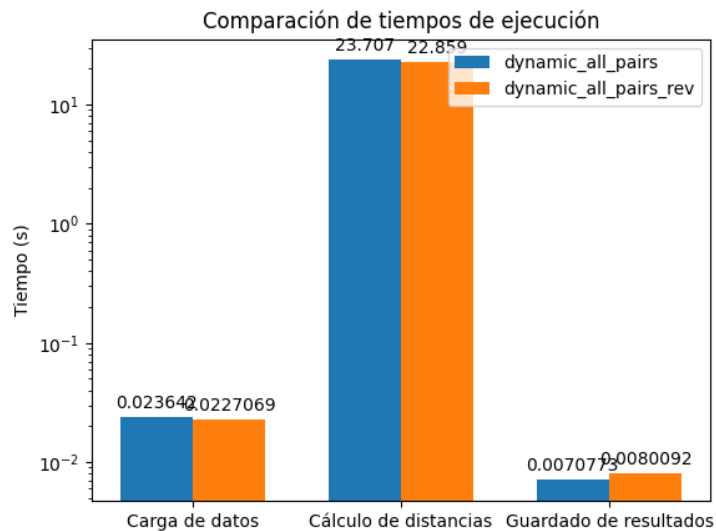


Figura 1: Grafico

Análisis del Tiempo de Ejecución

La diferencia clave se encuentra en el tiempo de cálculo de distancias:

- Tiempo de cálculo de distancias original: 23.707s
- Tiempo de cálculo de distancias invertido: 22.859s

La función `dynamic_all_pairs_rev` ha reducido el tiempo de cálculo de distancias en aproximadamente 0.848s, lo que representa una mejora de alrededor del 3.58 %.

Detalles adicionales de los resultados

- **Distribución de Carga:** Al comenzar desde la parte inferior de la matriz, donde las filas pueden tener más elementos relevantes (especialmente en matrices donde las filas inferiores tienden a ser más densas o grandes), la carga de trabajo puede distribuirse de manera más eficiente entre los hilos. Esto podría resultar en un mejor balanceo de carga y menor tiempo de espera para los hilos, ya que están procesando fragmentos más sustanciales desde el principio.
- **Localidad de Datos:** Acceder a filas consecutivas de manera invertida puede mejorar la localidad de los datos, aprovechando mejor las cachés de la CPU y reduciendo los tiempos de acceso a la memoria.
- **Optimización de la Sincronización:** Con menos hilos quedándose inactivos esperando por nuevas tareas, se puede optimizar el uso de la CPU y minimizar el overhead de sincronización.

La implementación de la función `dynamic_all_pairs_rev` ha demostrado ser más eficiente que la implementación original `dynamic_all_pairs`. Aunque la mejora en el tiempo de ejecución no es drástica, es significativa en un contexto de procesamiento paralelo y cálculo de matrices de distancias de gran tamaño.

5. link del repositorio:

[GitHub](#).