

Pruebas sobre el Comportamiento de la Memoria Cache: BUCLES ANIDADOS

Iñigo Manuel Diez Canseco Fuentes

Marzo 2022

1. Introducción

Este documento se realizó para resolver la pregunta propuesta en el aula el cual trata de los *FOR* anidados mostrados en la sesión 2 del día Miércoles 17 de Marzo de la clase de Computación Paralela y Distribuida.

Se implementó y comparó los 2-bucles anidados presentados en clase, además de contar con un análisis explicando los resultados y por último tener un enlace a un repositorio de *Github* donde se encontrará el trabajo.

En la siguiente figura podemos observar los 2-bucles que fueron mostrados en clase.

```
double A[MAX][MAX], x[MAX], y[MAX];
. . .
/* Initialize A and x, assign y = 0 */
. . .
/* First pair of loops */
for (i = 0; i < MAX; i++)
    for (j = 0; j < MAX; j++)
        y[i] += A[i][j]*x[j];
. . .
/* Assign y = 0 */
. . .
/* Second pair of loops */
for (j = 0; j < MAX; j++)
    for (i = 0; i < MAX; i++)
        y[i] += A[i][j]*x[j];
```

Figura 1: Pseudocódigo de los "for" del libro *An Introduction to Parallel Programming*



2. Implementación

El *link* del repositorio se encontrará en el siguiente enlace¹.

En la implementación hay una variación en los bucles que es vital para el funcionamiento y la forma en el cual se comportarían y que al ejecutarse iban a tener tiempos diferentes, lo cual se tuvo que utilizar librerías de tiempo para poder saber cuál bucle era más rápido que el otro.

A continuación está la implementación en el lenguaje de programación C++:

```
1  #include <ctime>
2  #include <iomanip>
3  #include <chrono>
4  #include <cstdlib>
5  #include <string>
6  #include <iostream>
7
8  using namespace std;
9  using namespace std::chrono;
10
11 #define MAX 15000
12 double A[MAX][MAX], x[MAX], y[MAX];
13
14 void primer_bucle() {
15     for (size_t i = 0; i < MAX ; i++) {
16         for (size_t j = 0; j < MAX ; j++) {
17             y[i] += A[i][j] * x[j];
18         }
19     }
20 }
21
22 void segundo_bucle() {
23     for (size_t j = 0; j < MAX ; j++) {
24         for (size_t i = 0; i < MAX ; i++) {
25             y[i] += A[i][j] * x[j];
26         }
27     }
28 }
29
30 void relleno(){
31     for (size_t i = 0; i < MAX; ++i) {
32         for (size_t j = 0; j < MAX; ++j) {
33             A[i][j] = rand() % MAX;
34         }
35     }
36     for (size_t i = 0; i < MAX; ++i) {
37         x[i] = rand() % MAX;
38         y[i] = rand() % MAX;
39     }
40 }
41
42
43
```

¹https://github.com/inigomanuel/ComputacionParalela_y_Distribuida



```
44 void Tiempo_Primer_Bucle() {
45     std::chrono::time_point<std::chrono::high_resolution_clock>
46     start, end;
47     start = std::chrono::high_resolution_clock::now();
48     primer_bucle();
49     end = std::chrono::high_resolution_clock::now();
50     int64_t duration = std::chrono::duration_cast<std::chrono::
51     microseconds>(end - start).count();
52     std::cout << std::setw(30) << "Microseg: " + std::to_string(
53     duration) + " us\n";
54 }
55
56 void Tiempo_Segundo_Bucle() {
57     std::chrono::time_point<std::chrono::high_resolution_clock>
58     start, end;
59     start = std::chrono::high_resolution_clock::now();
60     segundo_bucle();
61     end = std::chrono::high_resolution_clock::now();
62     int64_t duration = std::chrono::duration_cast<std::chrono::
63     microseconds>(end - start).count();
64     std::cout << std::setw(30) << "Microseg: " + std::to_string(
65     duration) + " us\n";
66 }
67
68 int main() {
69     relleno();
70     cout << "Primer Bucle: " << endl;
71     Tiempo_Primer_Bucle();
72     cout << endl;
73     cout << "Segundo Bucle: " << endl;
74     Tiempo_Segundo_Bucle();
75     return 0;
76 }
```

Listing 1: Laboratorio₁.cpp

3. Resultados

En la siguiente imagen donde se encuentra la tabla de resultados, el primer bucle es más rápido debido a que se obtiene el valor avanzando por filas en una columna, en cambio, en el segundo obtiene el valor avanzando por columnas en una fila. En estos movimientos los índices toman un papel importante.

Valor de Elementos	Primer Bucle		Segundo Bucle	
	Microsegundos	Segundos	Microsegundos	Segundos
5000	58399	0.058399	145858	0.145858
10000	59597	0.059597	157669	0.157669
15000	521570	0.52157	2105226	2.105226

Figura 2: Tabla de tiempo de los resultados de la comparación en microsegundos y segundos.



Para lograr estos resultados se utilizó entre 5000 a 15000 elementos para los arreglos y la matriz.

4. Análisis de la Ejecución

Primero hay que mencionar que nosotros al programar no podemos determinar que *data* o cuál instrucción está en la *cache*. Al tomar en cuenta esto, primero explicamos que sucede con ambos bucles, primero tenemos que el primer bucle obtiene el valor avanzando por filas en una columna, esto significa, que si estamos en $A[0][0]$ y como este dato no está en la *cache* entonces el resultado será un *cache miss* y el sistema leerá la línea que consta en la primera fila de A, estos serían $A[0][0]$, $A[0][1]$, $A[0][2]$, $A[0][3]$ en la *cache*, debido a esto el primer bucle accederá a $A[0][1]$, $A[0][2]$, $A[0][3]$, que ya se encontrarían en la *cache*.

En el segundo bucle, necesitará acceder a $A[1][0]$, $A[2][0]$, $A[3][0]$, ya que primero se encuentra en $A[0][0]$, entonces ninguno de estos está en la *cache* entonces, los próximos tres accesos de A también resultarán en *misses*. Desde que $A[2][0]$ está en la *cache* en la línea 2, leer su línea desalojará la línea 0 y leer $A[3][0]$ desalojará la línea 1. Después a través del bucle exterior, se necesita acceder a $A[0][1]$ el cual ya había sido desalojado antes, lo que implica que si cada vez que se lee un elemento de A, tendremos un *miss*.

En resumen, en el segundo bucle habrá más *cache misses* y accederá más a memoria que en el primer bucle, debido a esto el primer bucle es mucho más rápido que el segundo.