

6620. Organización de Computadoras

Trabajo Práctico 2: Optimización de Software

González, Juan Manuel, *Padrón Nro. 79.979*
juan0511@yahoo.com

Smocovich, Elizabeth, *Padrón Nro. 86.326*
liz.smocovich@gmail.com

Pereira, María Florencia, *Padrón Nro. 88.816*
mflorenciapereira@gmail.com

1er. Cuatrimestre de 2010

66.20 Organización de Computadoras – Práctica Jueves

Facultad de Ingeniería, Universidad de Buenos Aires

Resumen

El presente trabajo tiene como objetivo aplicar optimizaciones de software en ambientes reales y familiarizarse con la medición de comportamiento utilizando herramientas de *profiling*.

1. Introducción

Un algoritmo de la multiplicación de matrices en lenguaje C sencillo, pero ineficiente en cuanto a la utilización de la memoria cache es la siguiente:

```
for(i=0;i<N;i++)
    for(j=0;j<M;j++)
        for(k=0;k<L;k++)
            c[i][j]+=a[i][k]*[k][j];
```

Este programa multiplica las matrices $A \in \mathbb{R}^{N \times L}$ y $B \in \mathbb{R}^{L \times M}$, dando como resultado una matriz $C \in \mathbb{R}^{N \times M}$. Al recorrer una de las matrices por filas y la otra por columnas, para matrices mayores que la memoria cache se producirán ineficiencias a la hora de operar con ella.

En el presente trabajo se hará uso de las técnicas de optimización de uso de la memoria cache vistas en el curso, modificando un programa que resuelve una multiplicación de dos matrices. Se buscará reducir la cantidad de desaciertos (*misses*) en la memoria caché de datos de nivel 1 (L1) y el tiempo de ejecución. Para esto se utilizará el módulo de simulación de memorias caché, llamado Cachegrind, e integrado a Valgrind y la herramienta de profiling gprof. Las optimizaciones a aplicar serán *padding*, *blocking* y *prefetch*.

2. Consideraciones sobre el desarrollo e implementación

A continuación se detallan las consideraciones tenidas en cuenta para el desarrollo del trabajo práctico.

En esta sección se discutirán los métodos de optimización que fueron implementados en la versión final del trabajo práctico. También se expondrán los resultados obtenidos de las mejoras introducidas en el código, basados en las pruebas realizadas según se expone en la sección 3.

2.1. Elección del nivel de memoria a optimizar

Para reducir la brecha de velocidad entre el procesador y la memoria física, la arquitectura de las computadoras suele incluir varios niveles de memoria. Las memorias caches multi-nivel operan revisando la memoria L1 primero. Si se produce un miss, a continuación se revisa la cache L2 y así sucesivamente con los niveles siguientes hasta la memoria física. Las memorias cache grandes tienen un porcentaje alto de hit rate pero una latencia mayor. El primer nivel de memoria es pequeño para que los ciclos que se consuman sean comparables con los ciclos de reloj del CPU. La memoria cache de nivel 2 es de mayor tamaño para capturar gran parte de los accesos a memoria, disminuyendo así el *miss penalty*. Dado el tamaño de la memoria L1, es probable que su miss rate sea mayor que el de la memoria L2. Por esta razón, se decidió optimizar el código para la memoria cache L1, ya que si se redujera el miss rate de este nivel, el tiempo de ejecución sería menor. Para la memoria cache L2 esto no sería posible ya que el miss rate es bastante menor que el de la L1 y la latencia no es una característica que pueda optimizarse por software.

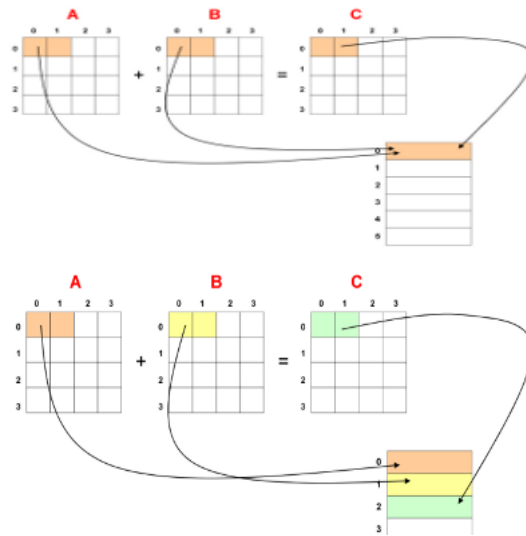
2.2. Optimización del algoritmo y resultados

A fin de mejorar el rendimiento, se ensayaron tres soluciones:

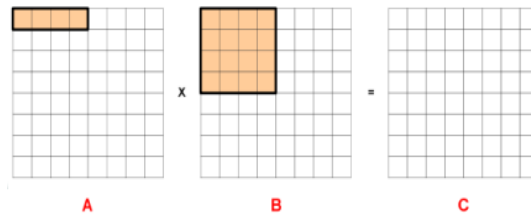
Padding

Se conoce como *Padding* al agregado de variables sin usar entre dos *arrays* para reducir la interferencia entre ellos y de esta forma mejorar la tasa de *miss*. Agregar un *pad* modifica el *offset* del segundo *array*, por lo que se puede lograr que estos queden mapeados en distintas partes de la cache, evitando así conflictos. El tamaño de los *pad* depende del esquema de *mapeo* de la memoria cache, su grado de asociatividad y el patrón de acceso a los datos del código. Valores típicos suelen ser múltiplos del tamaño de línea. De esta manera, se logra evitar el acceso a datos distintos ubicados en la misma posición de memoria cache (*thrashing*). Esta optimización cambia la forma en que los arreglos son almacenados en memoria principal, reduciendo los accesos conflictivos. Una desventaja es que utiliza mayor cantidad de memoria.

En cuanto a los resultados obtenidos mediante esta técnica, no se observó una mejora significativa en la ejecución del algoritmo. Esto puede deberse a que el tamaño de la matriz sobre la que estamos operando tiene dimensión de $N \times N$ donde $N=576$, que para la arquitectura donde se llevaron a cabo las pruebas, era múltiplo del tamaño de la línea de cache (64 bytes).



Loop blocking o tiling



Loop blocking es una transformación de los ciclos en los cuales se incrementa su profundidad. Si la profundidad original es n , el resultado podrá ser de cualquier profundidad entre $n+1$ y $2n$. Se busca mediante este método resolver mejor la localidad de los datos, para que la mayoría de ellos se encuentren en el cache en el momento y orden que se pretende leerlos de memoria. Esta técnica permite decrementar la cantidad de desajustes de capacidad cuando se opera con arreglos grandes y mejora la localidad temporal (mantiene en la cache, datos que se usarán en el corto plazo). Asimismo, se reducen los accesos conflictivos, ya que los bloques pequeños pueden ser mantenidos en la cache.

Para esta optimización existe un parámetro llamado blocking factor (B). Este factor deberá ser elegido para que las submatrices $B \times B$ entren completamente en la memoria cache. En la práctica, si bien al variar este parámetro se producía una reducción importante en la tasa de miss, esto no sucedía con el tiempo de ejecución. Una explicación para esto es que si bien se reducen la cantidad de miss en el total de accesos, al incorporar más lazos se ejecutarán mayor cantidad de instrucciones, lo cual hace incrementar los ciclos de ejecución.

Software prefetch

Software prefetch es una técnica de optimización que especula con los accesos futuros a datos e instrucciones. El objetivo es traer a la memoria cache los datos que se necesitarán en el futuro, al mismo tiempo que se procesan los datos actuales. En lenguaje C, es posible realizar una lectura adelantada mediante `__builtin_prefetch()`, si la arquitectura lo soporta. La operación de prefetch es de tipo *non-blocking*, razón por la cual para poder implementar software prefetching se deberá contar con una arquitectura que permita al procesador continuar funcionando mientras se realiza la operación. *Software prefetching* suele utilizarse frecuentemente cuando se tienen lazos de gran dimensión dado que exhiben un patrón predecible de accesos. Para aplicar esta técnica se dividió el lazo correspondiente a la matriz A (que se recorre por filas) en 3 partes: prólogo, procesamiento principal y epílogo. En el prólogo se traen tantos elementos como entren en una línea de cache. En el procesamiento principal, se traen a memoria los siguientes datos y se procesan los que ya se encuentran en ella. Finalmente, en el epílogo se termina de procesar los últimos datos. No se aplicó prefetch a la matriz B porque, al ser ésta recorrida por columnas, siempre producirá un miss dada la estructura de la memoria. Un parámetro configurable de esta técnica es la distancia de prefetch D, que es igual a la cantidad de bloques leídos en forma adelantada dividida la cantidad de bloques preprocesados. Otra definición dice que la distancia $D = l/s$ donde l es la latencia promedio de la memoria y s la cantidad de ciclos estimado de la iteración más

rápida del lazo. Tomamos $D=8$ dado que los resultados fueron mejores para la arquitectura donde se realizaron las pruebas. Un aspecto importante es que esta optimización permitió reducir el tiempo de ejecución significativamente, solucionando el problema mencionado anteriormente en la descripción de *Loop blocking*.

Otras mejoras posibles

Otras opciones para optimizar aún mas el algoritmo, fuera del alcance de nuestra implementación, incluyen por ejemplo el uso de alineación de datos (*data alignment*), que evita la fragmentación de los datos en la caché.

Otra posible optimización sería loop fusion que aplica a lazos con el mismo espacio de iteración una transformación que fusiona los cuerpos de cada lazo en uno solo.

3. Instalación

A continuación se detallan los pasos para instalar y ejecutar el trabajo práctico

1. Descomprimir el archivo TP2.tar.gz en una carpeta.
2. En un terminal, dentro de la carpeta del paso anterior, escribir `./instalar`.
3. En un terminal, dentro de la carpeta del paso anterior, escribir `./tp2 [opciones]`.

4. Corridas de prueba

En esta sección se presentan, a modo demostrativo, algunas de las corridas efectuada para probar el funcionamiento del trabajo práctico

También se comentan los resultados que se obtuvieron al variar los parámetros de cada optimización Las pruebas se realizaron para una memoria cache de datos, cuyas características son:

- Cantidad de vías: 8
- Tamaño: 32K
- Cantidad de sets: 64

Padding

- $\text{pad}=2 \times (\text{tamaño de la línea})$

```

florencia@florencia-laptop:~/tpsorga/TP2/pruebas$ ./tp2

Obteniendo el tiempo de ejecución en modo normal ...
Obteniendo la tasa de miss en modo normal ...

Obteniendo el tiempo de ejecución en modo optimizado ...
Obteniendo la tasa de miss en modo optimizado ...

Cache de datos L1:
#Vías: 8
Tamaño: 32K
Cantidad de sets: 64
Tamaño de la línea: 64 bytes

Tiempo de ejecución en modo normal: 2.70 ms
Tasa de miss en modo normal: 7.2
Tiempo de ejecución en modo optimizado: 2.95 ms
Tasa de miss en modo optimizado: 7.2

Speedup logrado: .91
Cociente de tasas de miss: 1.00

florencia@florencia-laptop:~/tpsorga/TP2/pruebas$ █

```

Padding

- $\text{pad} = 4 * (\text{tamaño de la línea})$

```

florencia@florencia-laptop:~/tpsorga/TP2/pruebas$ ./tp2

Obteniendo el tiempo de ejecución en modo normal ...
Obteniendo la tasa de miss en modo normal ...

Obteniendo el tiempo de ejecución en modo optimizado ...
Obteniendo la tasa de miss en modo optimizado ...

Cache de datos L1:
#Vías: 8
Tamaño: 32K
Cantidad de sets: 64
Tamaño de la línea: 64 bytes

Tiempo de ejecución en modo normal: 2.78 ms
Tasa de miss en modo normal: 7.2
Tiempo de ejecución en modo optimizado: 3.09 ms
Tasa de miss en modo optimizado: 7.2

Speedup logrado: .89
Cociente de tasas de miss: 1.00

florencia@florencia-laptop:~/tpsorga/TP2/pruebas$ █

```

Padding

- $\text{pad} = 8 * (\text{tamaño de la línea})$

```
florecia@florecia-
Archivo Editar Ver Terminal Ayuda
florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ ./tp2

Obteniendo el tiempo de ejecución en modo normal ...
Obteniendo la tasa de miss en modo normal ...

Obteniendo el tiempo de ejecución en modo optimizado ...
Obteniendo la tasa de miss en modo optimizado ...

Cache de datos L1:
#Vías: 8
Tamaño: 32K
Cantidad de sets: 64
Tamaño de la línea: 64 bytes

Tiempo de ejecución en modo normal: 2.77 ms
Tasa de miss en modo normal: 7.2
Tiempo de ejecución en modo optimizado: 3.15 ms
Tasa de miss en en modo optimizado: 7.2

Speedup logrado: .87
Cociente de tasas de miss: 1.00

florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ █
```

Como se mencionó anteriormente, no se observan cambios significativos al agregar padding a las matrices.

Padding y Blocking

- $\text{pad}=2 \times \text{tamaño de la línea}$
- $B=2$

```
florecia@florecia-
Archivo Editar Ver Terminal Ayuda
florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ ./tp2

Obteniendo el tiempo de ejecución en modo normal ...
Obteniendo la tasa de miss en modo normal ...

Obteniendo el tiempo de ejecución en modo optimizado ...
Obteniendo la tasa de miss en modo optimizado ...

Cache de datos L1:
#Vías: 8
Tamaño: 32K
Cantidad de sets: 64
Tamaño de la línea: 64 bytes

Tiempo de ejecución en modo normal: 2.75 ms
Tasa de miss en modo normal: 7.2
Tiempo de ejecución en modo optimizado: 6.13 ms
Tasa de miss en en modo optimizado: 0.6

Speedup logrado: .44
Cociente de tasas de miss: 12.00

florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ █
```

Padding y Blocking

- Con $\text{pad}=2 \times \text{tamaño de la línea}$
- $B=4$

```
florencia@florencia-laptop:~/tpsorga/TP2/pruebas$ ./tp2

Obteniendo el tiempo de ejecución en modo normal ...
Obteniendo la tasa de miss en modo normal ...

Obteniendo el tiempo de ejecución en modo optimizado ...
Obteniendo la tasa de miss en modo optimizado ...

Cache de datos L1:
#Vías: 8
Tamaño: 32K
Cantidad de sets: 64
Tamaño de la línea: 64 bytes

Tiempo de ejecución en modo normal: 2.68 ms
Tasa de miss en modo normal: 7.2
Tiempo de ejecución en modo optimizado: 4.42 ms
Tasa de miss en modo optimizado: 0.2

Speedup logrado: .60
Cociente de tasas de miss: 36.00

florencia@florencia-laptop:~/tpsorga/TP2/pruebas$
```

Padding y Blocking

- $\text{pad} = 2 \times \text{tamaño de la línea}$
- $B = (\text{cache_size}) / 8$

```
florencia@florencia-laptop:~/tpsorga/TP2/pruebas$ ./tp2

Obteniendo el tiempo de ejecución en modo normal ...
Obteniendo la tasa de miss en modo normal ...

Obteniendo el tiempo de ejecución en modo optimizado ...
Obteniendo la tasa de miss en modo optimizado ...

Cache de datos L1:
#Vías: 8
Tamaño: 32K
Cantidad de sets: 64
Tamaño de la línea: 64 bytes

Tiempo de ejecución en modo normal: 2.76 ms
Tasa de miss en modo normal: 7.2
Tiempo de ejecución en modo optimizado: 4.83 ms
Tasa de miss en modo optimizado: 0.1

Speedup logrado: .57
Cociente de tasas de miss: 72.00

florencia@florencia-laptop:~/tpsorga/TP2/pruebas$
```

Podemos observar que el mientras más grande sea B (siempre siendo menor al tamaño de la memoria cache), disminuye más la tasa de miss. Sin embargo, el tiempo de ejecución es mayor al original. Como se mencionó antes, esto puede deberse a que ahora se ejecutan más instrucciones por el agregado de más lazos. Para poder implementar prefetch, se tomó $B = (\text{cache_size}) / 8 \times D$

Prefetch, Padding y Blocking

- $\text{pad} = 2 \times \text{tamaño de la línea}$

- $D=2$
- $B=D*(\text{cache_size})/8$

```
florecia@florecia-
Archivo Editar Ver Terminal Ayuda
florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ ./tp2

Obteniendo el tiempo de ejecución en modo normal ...
Obteniendo la tasa de miss en modo normal ...

Obteniendo el tiempo de ejecución en modo optimizado ...
Obteniendo la tasa de miss en modo optimizado ...

Cache de datos L1:
#Vías: 8
Tamaño: 32K
Cantidad de sets: 64
Tamaño de la línea: 64 bytes

Tiempo de ejecución en modo normal: 2.70 ms
Tasa de miss en modo normal: 7.2
Tiempo de ejecución en modo optimizado: 2.06 ms
Tasa de miss en en modo optimizado: 0.5

Speedup logrado: 1.31
Cociente de tasas de miss: 14.40

florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ █
```

Prefetch, Padding y Blocking

- $\text{pad}=2*\text{tamaño de la línea}$
- $D=4$
- $B=D*(\text{cache_size})/8$

```
florecia@florecia-
Archivo Editar Ver Terminal Ayuda
florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ ./tp2

Obteniendo el tiempo de ejecución en modo normal ...
Obteniendo la tasa de miss en modo normal ...

Obteniendo el tiempo de ejecución en modo optimizado ...
Obteniendo la tasa de miss en modo optimizado ...

Cache de datos L1:
#Vías: 8
Tamaño: 32K
Cantidad de sets: 64
Tamaño de la línea: 64 bytes

Tiempo de ejecución en modo normal: 2.68 ms
Tasa de miss en modo normal: 7.2
Tiempo de ejecución en modo optimizado: 1.96 ms
Tasa de miss en en modo optimizado: 0.5

Speedup logrado: 1.36
Cociente de tasas de miss: 14.40

florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ █
```

Prefetch, Padding y Blocking

- $\text{pad}=2*\text{tamaño de la línea}$

- $D=8$
- $B=D*(\text{cache_size})/8$

```
florecia@florecia-
Archivo Editar Ver Terminal Ayuda
florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ ./tp2

Obteniendo el tiempo de ejecución en modo normal ...
Obteniendo la tasa de miss en modo normal ...

Obteniendo el tiempo de ejecución en modo optimizado ...
Obteniendo la tasa de miss en modo optimizado ...

Cache de datos L1:
#Vías: 8
Tamaño: 32K
Cantidad de sets: 64
Tamaño de la línea: 64 bytes

Tiempo de ejecución en modo normal: 2.71 ms
Tasa de miss en modo normal: 7.2
Tiempo de ejecución en modo optimizado: 1.92 ms
Tasa de miss en en modo optimizado: 0.5

Speedup logrado: 1.41
Cociente de tasas de miss: 14.40

florecia@florecia-laptop:~/tpsorga/TP2/pruebas$ █
```

Observamos que con el aumento de la distancia de prefetch, mejora el tiempo de ejecución y la tasa de miss se mantiene constante.

Como conclusión decidimos presentar esta última opción dado que presenta una mejora importante tanto en el tiempo de ejecución como en la tasa de miss.

5. Código fuente original

6. Código fuente modificado

A continuación se incluye el código fuente del programa, con las modificaciones incorporadas.

7. Código en shell script de instalar y tp2

A continuación se incluyen los archivos ejecutables que realizan la instalación y medición de tiempos de ejecución y miss rates.

8. Enunciado

A continuación se incluye el enunciado del trabajo práctico.

9. Conclusiones

El estudio de la problemática de la multiplicación entre matrices cuadradas generó diversos algoritmos para su optimización. Las condiciones que se presentan en la misma (saltos en la memoria principal, *thrashing*, gran cantidad de ciclos) la vuelven ideal para tomarla como caso de estudio. Utilizando *padding*, no logramos mejorar significativamente el rendimiento de la memoria cache al ejecutar la aplicación. Por otro lado, encarando la resolución del problema con la técnica de *blocking* se logró una importante mejora en el rendimiento, aunque no así en el tiempo de ejecución. Finalmente, incluir *software prefetch* produce una mejora en este último aspecto.

Al juntar todas las optimizaciones en un solo código fuente, pudimos observar que se obtiene una mejora global. Es importante seleccionar los parámetros de cada optimización de forma tal que se obtenga un mayor beneficio. En definitiva, la multiplicación de matrices es una situación ideal para ejercitar los distintos métodos de optimización investigados.

A modo de conclusión, podemos comentar que una vez concluidos los casos de prueba y luego de sendos testeos, notamos que se obtiene una mejora en el procesamiento y la aplicación hace mejor uso de los recursos del sistema.

Referencias

- [1] Herramientas de profiling Valgrind. <http://valgrind.org/>
- [2] The Cache Performace and Optimizations of Blocked Algorithms. Monica S. Lam an Edward E. Rothberg and Michael E. Wolf. Computer Systems Laboratory, Stanford University. <http://suif.stanford.edu/papers/lam91.ps>
- [3] GCC: Online Documentation. <http://gcc.gnu.org/onlinedocs/>
- [4] Oetiker, Tobias, “The Not So Short Introduction To LaTeX2”, <http://www.physics.udel.edu/~dubois/lshort2e/>
- [5] Kernighan, Brian W. / Ritchie, Dennis M. El Lenguaje De Prorgamación C. Segunda Edición, PRENTICE-HALL HISPANOAMERICANA SA, 1991.
- [6] Patterson / Hennesy. Computer Organization and Design; The Hardware-Software Interface. Second Edition. Morgan Kaufman, 1997.