

Clase 27. Uso de la memoria Cache

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

Cache

Introducción

La memoria se copia por bloques de la memoria principal (RAM) a la cache L3 compartida por todos los cores, de esta al L2 y de ahí al L1. Las diferencias de velocidad de la memoria son de entre uno y dos ordenes de magnitud (uno con el L3, 2 con el L1). Utilizar variables que ya están en la cache L1, puede significar un ahorro considerable de tiempo. Todas las técnicas de optimización de la cache se basan en 2 cosas:

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

Introducción

La memoria se copia por bloques de la memoria principal (RAM) a la cache L3 compartida por todos los cores, de esta al L2 y de ahí al L1. Las diferencias de velocidad de la memoria son de entre uno y dos ordenes de magnitud (uno con el L3, 2 con el L1). Utilizar variables que ya están en la cache L1, puede significar un ahorro considerable de tiempo. Todas las técnicas de optimización de la cache se basan en 2 cosas:

- Cercanía temporal de la memoria que se usa. Ej. Reutilizar lo mas pronto posible una variable. Si ocupo x en alguna operación, lo mejor sería reutilizar todo lo que pudiera y realizar todas las operaciones posible con x antes de realizar operaciones con y , ya que el cache es una memoria muy pequeña, y y puede ocupar el lugar de x , lo que requeriria bajar x de la RAM cada que se utilice.

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

Introducción

La memoria se copia por bloques de la memoria principal (RAM) a la cache L3 compartida por todos los cores, de esta al L2 y de ahí al L1. Las diferencias de velocidad de la memoria son de entre uno y dos ordenes de magnitud (uno con el L3, 2 con el L1). Utilizar variables que ya están en la cache L1, puede significar un ahorro considerable de tiempo. Todas las técnicas de optimización de la cache se basan en 2 cosas:

- Cercanía temporal de la memoria que se usa. Ej. Reutilizar lo mas pronto posible una variable. Si ocupo x en alguna operación, lo mejor sería reutilizar todo lo que pudiera y realizar todas las operaciones posible con x antes de realizar operaciones con y , ya que el cache es una memoria muy pequeña, y y puede ocupar el lugar de x , lo que requeriria bajar x de la RAM cada que se utilice.
- Cercanía espacial. Ej. Si $x[4]$ es un vector, dado que la memoria se copia por bloques o líneas, lo mas eficiente sería utilizar $x[0]$, $x[1]$, $x[2]$ y $x[3]$ consecutivamente (en el caso de paralelo tambien sería eficiente usarlos en el mismo core).

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

Línea de caché

La memoria se copia por líneas de y hacia la caché, si se utilizan datos en la misma línea no es necesario descargar de la RAM otra vez la copia de la memoria. Mientras más alejado esté un dato en la memoria más tarda la operación, es decir, no es lo mismo leer/escribir sobre $x[1]$, $x[2]$ que sobre $x[1]$, $x[1260]$. Si los datos están muy alejados no dejará de importar su distancia ya que no se aprovechará para nada la línea de cache.

Ejemplo:

```
1 const int n=200000000;  
2 double *x= (double*) malloc (n* sizeof (double));  
3 double iniTime=get_time ();  
4     for (register int i=0; i<n; i++)  
5         x[i]=5;  
6     printf ("Tiempo de llenado consecutivo=%lf\n",  
7         get_time ()-iniTime);  
8     free (x);
```

El tiempo de ejecución de este código es de aproximadamente 0.27 seg.

Ej. saltos

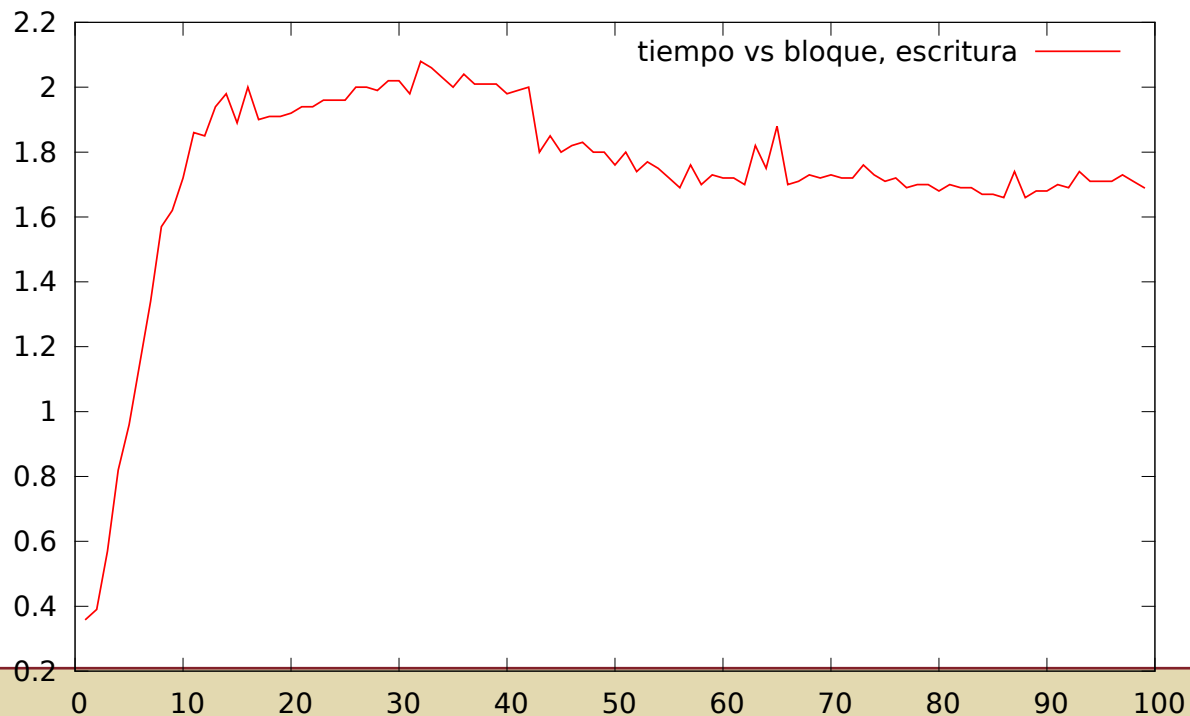
Se hacen las mismas operaciones (más o menos) dando **saltos**, para no acceder a memoria continua. Ejemplo de escritura.

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

```
1  for (register int nb=1; nb<100; nb++){
2      double iniTime=get_time();
3      for (register int j=0; j<nb; j++)
4          for (register int i=0; (i+j)<n; i+=nb)
5              x[i+j]=5;
6      printf("%d %lf \n",nb, get_time()-iniTime);
7  }
```



Ej. saltos

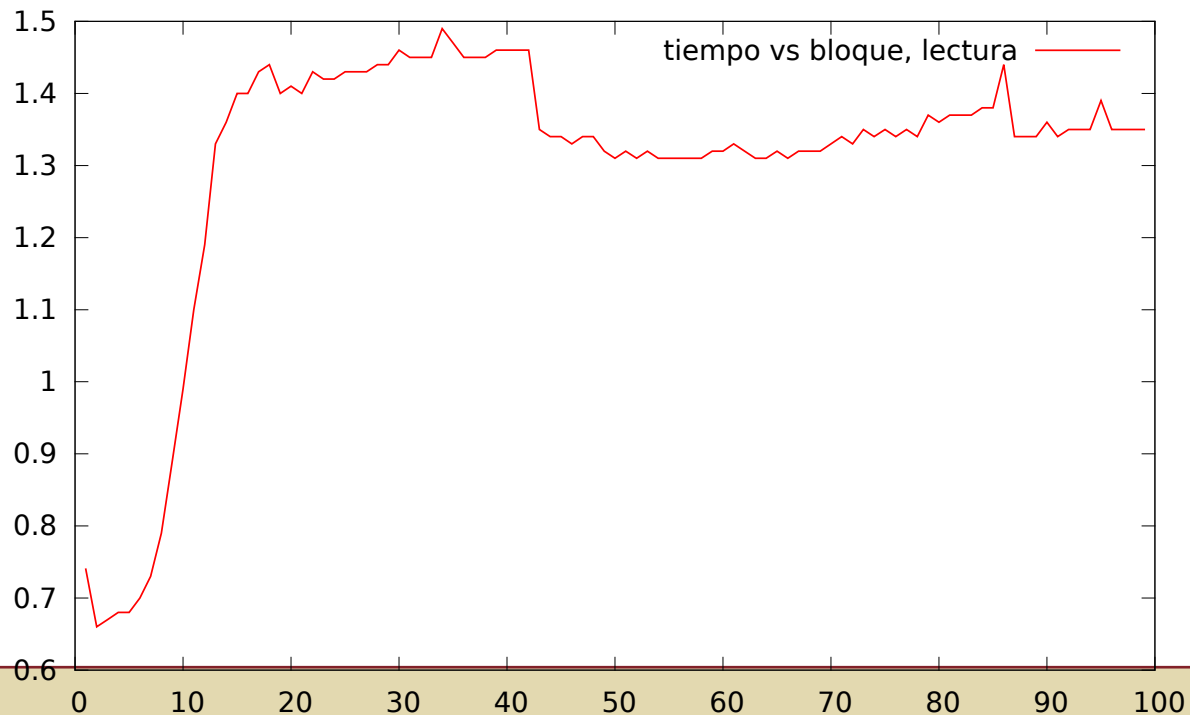
Se hacen las mismas operaciones (más o menos) dando **saltos**, para no acceder a memoria continua. Ejemplo de **lectura**.

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

```
1 for (register int nb=1; nb<100; nb++){
   double iniTime=get_time();      sum=0.0;
3   for (register int j=0; j<nb; j++)
      for (register int i=0; (i+j)<n; i+=nb)
5         sum+=x[i+j];
      printf("%d %lf \n",nb, get_time()-iniTime);
7 }
```



Ej. saltos

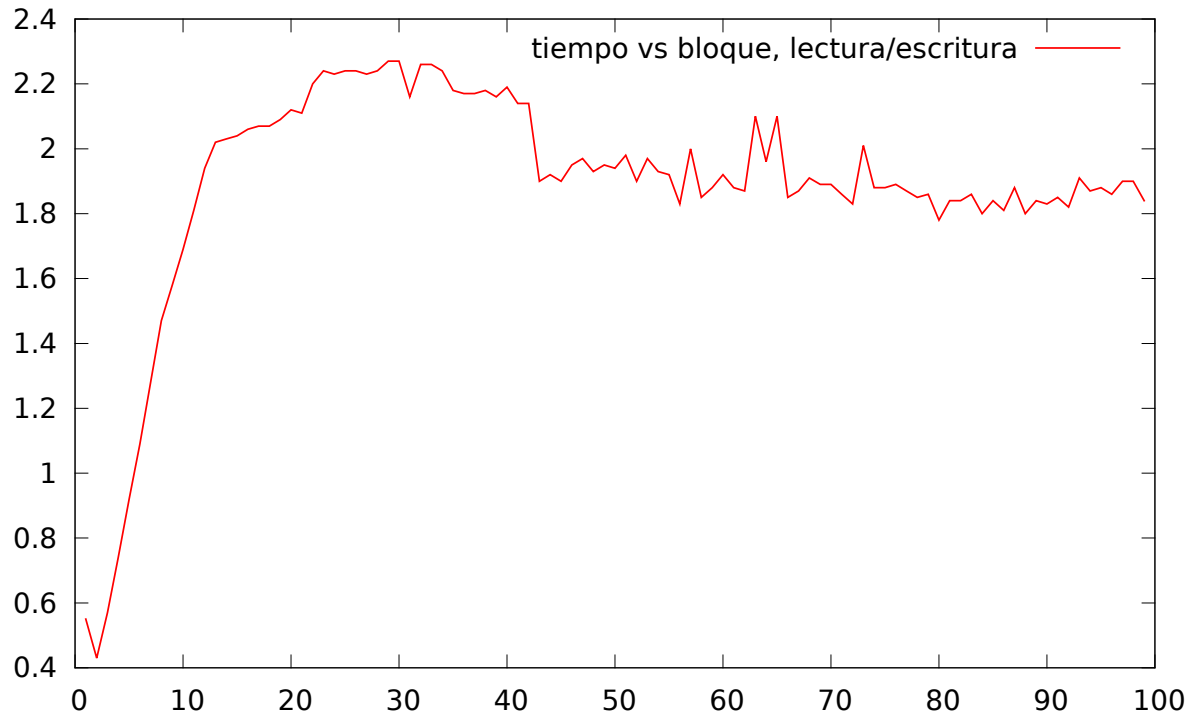
Ejemplo de lectura/escritura.

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

```
1 for (register int nb=1; nb<100; nb++){
2     double iniTime=get_time();      sum=0.0;
3     for (register int j=0; j<nb; j++)
4         for (register int i=0; (i+j)<n; i+=nb)
5             sum+=x[i+j];
6     printf("%d %lf \n",nb, get_time()-iniTime);
7 }
```



Ej. cercanía temporal

El calculo de los productos en Forma 1 tarda 0.95 seg. la Forma 2 0.41

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

```
1 #Forma 1
  for (register int i=0; i<n; i++)
3     px+=x[i]*x[i];
  for (register int i=0; i<n; i++)
5     py+=y[i]*y[i];
  for (register int i=0; i<n; i++)
7     pxy+=x[i]*y[i];
  #Forma 2
9  for (register int i=0; i<n; i+=4){
  px+=x[i]*x[i]+x[i+1]*x[i+1]+x[i+2]*x[i+2]+x[i+3]*x[i
    +3];
11 py+=y[i]*y[i]+y[i+1]*y[i+1]+y[i+2]*y[i+2]+y[i+3]*y[i
    +3];
  pxy+=x[i]*y[i]+x[i+1]*y[i+1]+x[i+2]*y[i+2]+x[i+3]*y[i
    +3];
13 }
```

Cache

cachegrind

Técnicas mas
usuales para
mejorar el uso del
cache

cachegrind

cachegrind

Cache

cachegrind

Técnicas mas
usuales para
mejorar el uso del
cache

Valgrind contiene una herramienta para analizar el comportamiento del cache, básicamente los cache hits y cache miss. El uso es el siguiente:

```
valgrind --tool=cachegrind prog
```

Donde prog es el ejecutable. Cachegrind simulara el programa y los accesos al primer y último cache (usualmente L1 y L3). Como se quiere simular el comportamiento real, se recomienda que prog ya sea la versión (compilado) optimizado del programa. Usualmente la simulación de cachegrind se tarda mucho mas que la ejecución normal del programa, se recomienda ejecutar con un ejemplo pequeño.

cachegrind

Cache

cachegrind

Técnicas mas
usuales para
mejorar el uso del
cache

Significado:

- I1= Cache de instrucciones del primer nivel (L1).
- D1= Cache de datos del primer nivel (D1).
- LL= Cache de último nivel (instrucciones y datos).
- Ir= lectura de instrucciones, I1mr= misses de lectura de intrucciones del cache 1. ILmr= misses de lectura de instrucciones del cache de último nivel.
- Dr= Lectura de datos, D1mr= misses de lectura de datos en el primer cache, DLmr= misses de lectura de datos en el último cache.
- Dw= Escritura de datos... lo mismo que los de arriba.

cachegrind

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

Para algún programa entrega algo como esto.

```
==12035== I    refs:          1,975,017,358
==12035== I1  misses:              760
==12035== LLi misses:          753
==12035== I1  miss rate:         0.00%
==12035== LLi miss rate:        0.00%
==12035==
==12035== D    refs:          713,229,506 (497,049,338 rd + 216,18
==12035== D1  misses:          27,023,905 ( 2,023,359 rd + 25,00
==12035== LLd misses:          27,023,482 ( 2,022,972 rd + 25,00
==12035== D1  miss rate:         3.7% ( 0.4% +
==12035== LLd miss rate:        3.7% ( 0.4% +
==12035==
==12035== LL  refs:          27,024,665 ( 2,024,119 rd + 25,00
==12035== LL  misses:          27,024,235 ( 2,023,725 rd + 25,00
==12035== LL  miss rate:         1.0% ( 0.0% +
```

cg_annotate

Cache

cachegrind

Técnicas mas
usuales para
mejorar el uso del
cache

Además de la información en pantalla, la cual es un resumen de la ejecución, cachegrind genera un archivo con información mas detallada, que tiene mas o menos la siguiente sintaxis:

```
cachegrind.out.12136
```

Donde 12136 es el ID del proceso.

Este archivo se puede visualizar con otra aplicación:
`cg_annotate`.

cg_annotate

Cache

cachegrind

Técnicas mas
usuales para
mejorar el uso del
cache

Se pueden extraer datos de este archivo por medio de **cg_annotate**

```
cg_annotate cachegrind.out.12136
```

Lo primero que despliega es información del procesador:

```
I1 cache:          32768 B, 64 B, 8-way associative  
D1 cache:          32768 B, 64 B, 8-way associative  
LL cache:         6291456 B, 64 B, 12-way associative
```

Tamaño, línea, y asociatividad. Un cache es completamente asociativo si una dirección de la RAM puede copiarse en cualquier posición del cache. Si solo se puede copiar en N posiciones, entonces el cache es N-asociativo.

cg_annotate

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

ejercicio1a: productos puntos

```
-----  
      Ir  I1mr  ILmr      Dr  D1mr  DLmr      Dw  D1mw  DLmw  
191,490 1,035 1,024 44,745 2,467 1,952 17,745 819 745 PROGRAM T
```

```
-----  
      Ir I1mr  ILmr      Dr  D1mr  DLmr      Dw  D1mw  DLmw  file:function  
9,029   3    3      3    1    0 2,009 248 248  ???:ejercicio1  
9,029   3    3      3    0    0 2,009 6  0  ???:ejercicio1
```

Linea de ejecución para que solo muestre cierta información (no corresponde a lo de arriba):

```
cg_annotate --show=Dr,Dw --threshold=0.8 cachegrind.out.12136
```

KCacheGrind

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

KCacheGrind es una herramienta para linux/KDE para visualizar la salida de cachegrind, que facilita muchisimo el análisis del código.

The screenshot shows the KCacheGrind application interface. The top menu bar includes File, View, Go, Settings, and Help. Below the menu is a toolbar with navigation buttons (Open, Back, Forward, Up) and analysis options like Relative, Cycle Detection, Relative to Parent, and Shorten Templates. The main window is titled 'ejercicio1a' and displays a 'Flat Profile' view. On the left, a list of functions is shown with their respective costs and locations. The right pane shows a table of event types with their inclusion and self costs, and a formula for each. The status bar at the bottom indicates the total LL Data Write Miss Cost.

Self	Function	Loca
49.94	ejercicio1a	eje01
49.94	ejercicio1b	eje01
0.05	_dl_relocate_object	(unkr
0.01	memset	(unkr
0.01	_dl_allocate_tls_storage	(unkr
0.01	malloc_init_state	(unkr
0.01	_dl_map_object_from_fd	(unkr
0.00	_dl_new_object	(unkr
0.00	vfprintf	(unkr
0.00	_dl_start	(unkr
0.00	dl_main	(unkr
0.00	__printf_fp	(unkr
0.00	_dl_lookup_symbol_x	(unkr
0.00	_dl_check_map_versions	(unkr
0.00	do_lookup_x	(unkr
0.00	_int_malloc	(unkr
0.00	memcpy	(unkr
0.00	_IO_file_overflow@@GLIBC...	(unkr
0.00	_dl_catch_error	(unkr
0.00	_dl_fini	(unkr
0.00	_dl_map_object	(unkr
0.00	check_match.9325	(unkr
0.00	mempcpy	(unkr

Event Type	Incl.	Self	Short	Formula
Instruction Fetch	53.16	53.16	Ir	
L1 Instr. Fetch Miss	0.83	0.83	ILmr	
LL Instr. Fetch Miss	0.83	0.83	ILmr	
Data Read Access	56.35	56.35	Dr	
L1 Data Read Miss	66.43	66.43	D1mr	
LL Data Read Miss	65.36	65.36	DLmr	
Data Write Access	64.40	64.40	Dw	
L1 Data Write Miss	49.94	49.94	D1mw	
LL Data Write Miss	49.94	49.94	DLmw	
L1 Miss Sum	59.79	59.79	L1m	$L1m = ILmr + D1mr + D1mw$
Last-level Miss Sum	59.03	59.03	LLm	$LLm = ILmr + DLmr + DLmw$
Cycle Estimation	56.21	56.21	CEst	$CEst = Ir + 10 L1m + 100 LLm$

cachegrind.out.9099 [1] - Total LL Data Write Miss Cost: 500 582

KCacheGrind

Para desplegar toda la información esto si debe de ser compilado con -g.

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache

The screenshot shows the KCacheGrind application interface. The top menu includes File, View, Go, Settings, and Help. Below the menu is a navigation bar with buttons for Open, Back, Forward, Up, and a Relative percentage selector. The main window is titled 'ejercicio1a' and has tabs for Types, Callers, All Callers, Callee Map, and Source Code. The Source Code tab is active, displaying the following code:

```
# : Ir : Source ('/home/ivvan/Dropbox/curso_programacion2014/ej
59      }
60
61      void ejercicio1a(int argc, char *argv[])
62      {
63          0.00  const int n=1000000;
64
65          0.00  double *x= (double*)malloc(n*sizeof(double));
66          0.00  double *y= (double*)malloc(n*sizeof(double));
67          0.00  double ppx=0.0, ppy=0.0, ppxy=0.0;
68
69          2.35  for (register int j=0; j<n; j++){
70              4.69      x[j]=j;
71              6.25      y[j]=n-j;
72          }
73
74          0.00  double iniTime=get time();
75          2.35  for (register int i=0; i<n; i++)
76          10.94  ppx+=x[i]*x[i];
77          2.35  for (register int i=0; i<n; i++)
78          10.94  ppy+=y[i]*y[i];
79          2.35  for (register int i=0; i<n; i++)
80          10.94  ppxy+=x[i]*y[i];
81
```

On the left side, there is a 'Flat Profile' section with a search bar and a '(No Grouping)' dropdown. Below it is a table listing functions and their costs:

Self	Function	Loca
53.16	ejercicio1a	eje0:
46.71	ejercicio1b	eje0:
0.05	_dl_addr	(unki
0.02	do_lookup_x	(unki
0.02	_dl_relocate_object	(unki
0.01	_dl_lookup_symbol_x	(unki
0.01	strcmp	(unki
0.00	check_match.9325	(unki
0.00	bsearch	(unki
0.00	_dl_map_object_from_fd	(unki
0.00	strsep	(unki
0.00	__printf_fp	(unki
0.00	dl_main	(unki
0.00	_dl_check_map_versions	(unki
0.00	_dl_name_match_p	(unki
0.00	_dl_map_object_deps	(unki
0.00	bcmp	(unki
0.00	vfprintf	(unki
0.00	_dl_cache_libcmp	(unki
0.00	intel_check_word	(unki
0.00	_int_malloc	(unki
0.00	_dl_start	(unki
0.00	_IO_file_overflow@@GLIBC...	(unki
0.00	malloc_init_state	(unki
0.00	ptmalloc_init	(unki
0.00	_dl_next_ld_env_entry	(unki
0.00	intel_02_known_compare	(unki
0.00	strlen	(unki

At the bottom of the window, there is a 'Parts' section with tabs for Parts, Callees, Call Graph, All Callees, Caller Map, and Machine Code. The status bar at the bottom indicates 'cachegrind.out.9749 [1] - Total Instruction Fetch Cost: 127 927 822'.

Cache

cachegrind

Técnicas mas
usuales para
mejorar el uso del
cache



Técnicas mas usuales para mejorar el uso del cache

Intercambio de loops

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache



Kowarschik, Markus, and Christian Weiß.

"An overview of cache optimization techniques and cache-aware numerical algorithms." Algorithms for Memory Hierarchies. Springer Berlin Heidelberg, 2003. 213-232.

Intercambiar los loops para acceder a los datos en el sentido del bloque de datos:

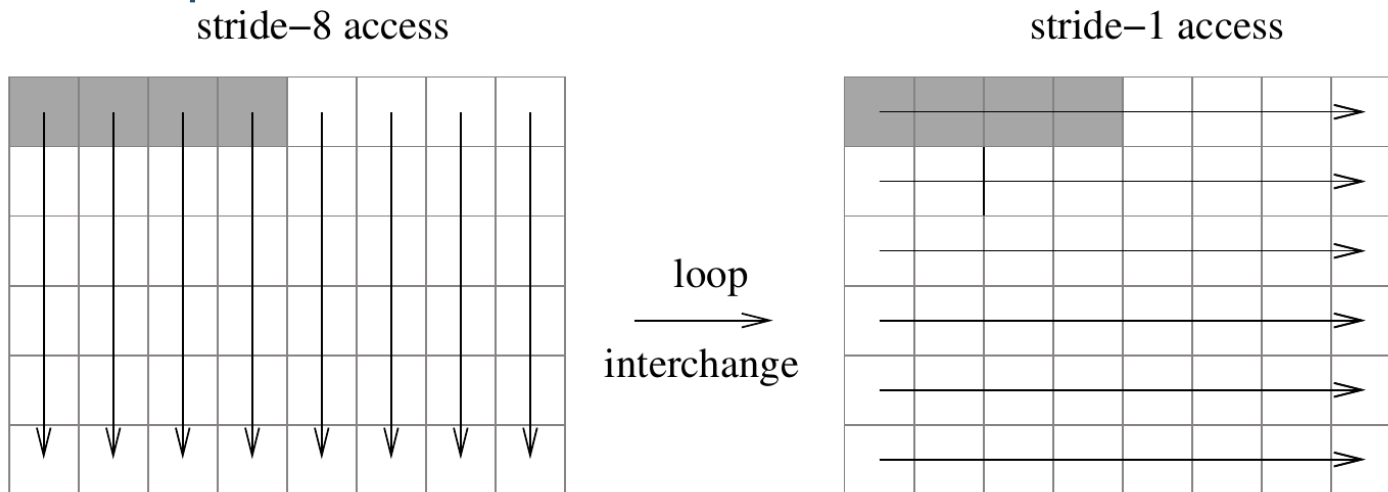


Fig. 2. Access patterns for interchanged loop nests.

Intercambio de loops

```
1 void ejercicio2 (int argc, char *argv [])
  {
3   const int nrow=20000, ncol=20000;
   double **M=matrix (nrow, ncol) ;;
5   double sum=0.0;

7   double iniTime=get_time ();
   for (register int i=0; i<nrow; i++)
9     for (register int j=0; j<ncol; j++) sum+=M[i][j];
   printf ("Tiempo renglon %lf \n", get_time ()-iniTime);

11
   iniTime=get_time ();
13   for (register int j=0; j<ncol; j++)
     for (register int i=0; i<nrow; i++) sum+=M[i][j];
15   printf ("Tiempo columnas %lf \n", get_time ()-iniTime);

17   free_matrix (M, nrow);
   }
19
```

Diferencia de tiempos: renglones=1.19 seg, columnas 5.62 seg.

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache



Fusión de loops

Mejora el loop overhead, la paralelización , y la espacialidad de los datos, reduciendo los misses.

Fission=0.38seg. Fussion=0.23seg.

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache



```
2  const int n=50000000;
3  double *x= (double*) malloc (n* sizeof (double) );
4  double *y= (double*) malloc (n* sizeof (double) );
5  double *z= (double*) malloc (n* sizeof (double) );
6  double iniTime=get_time ();
7  for (int i=0; i<n; i++)
8      x[i]=y[i]+1.0;
9  for (int i=0; i<n; i++)
10     z[i]+=x[i]*4.0;
11  printf ("loop fission %lf \n", get_time ()-iniTime );
12  iniTime=get_time ();
13  for (int i=0; i<n; i++){
14     x[i]=y[i]+1.0;
15     z[i]+=x[i]*4.0;
16 }
17 printf ("loop fussion %lf \n", get_time ()-iniTime );
18 free (x); free (y); free (z);
```

Loop Blocking

Se intenta aprovechar los datos que ya han bajado a la cache. normal =0.26seg. por bloques=0.02seg.

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache



```
1  const int n=5000;
   const int B=12;
3  double **a= matrix(n,n);
   double **b= matrix(n,n);
5
   double iniTime=get_time();
7   for (int i=0; i<n; i++)
      for (int j=0; j<n; j++)
9         a[i][j]=b[j][i];
   printf("transpuesta normal %lf\n",get_time()-iniTime);
11
   iniTime=get_time();
13   for (int ii=0; ii<n; ii+=B)
      for (int jj=0; jj<n; jj+=B)
15         for (int i=ii; i<min(ii+B-1,n); i++)
   for (int j=jj; j<min(jj+B-1,n); j+=B)
17         a[i][j]=b[j][i];
   printf("transpuesta bloques %lf\n",get_time()-iniTime);
19 free_matrix(a,n); free_matrix(b,n);
```


Organización de los datos

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache



-
- Array padding: Si dos arreglos son mapeados a la misma línea de cache y se acceden alternadamente, se pueden producir un gran número de misses, debido a que un arreglo sobrescribe los datos del otro. Una solución es insertar memoria intermedia en la RAM para que los arreglos mapeen a diferentes locaciones del cache.

Organización de los datos

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache



-
- Array padding: Si dos arreglos son mapeados a la misma línea de cache y se acceden alternadamente, se pueden producir un gran número de misses, debido a que un arreglo sobrescribe los datos del otro. Una solución es insertar memoria intermedia en la RAM para que los arreglos mapeen a diferentes locaciones del cache.
- Array merging. Supongamos que almaceno una matriz rala como un conjunto de valores e indices. Por lo general, para operaciones con matrices utilizaría el valor y el indice al mismo tiempo, entonces es lógico bajarlos a la misma linea de cache, esto se puede lograr con una estructura que contenga un valores y los indices y solicitando un arreglo de esta estructura.

Organización de los datos

Cache

cachegrind

Técnicas más usuales para mejorar el uso del cache



-
- Transposición del arreglo. Similar al intergambio de loop, se puede transponer una matriz (por ejemplo para la multiplicación) y mejorar el acceso desde el arreglo de los datos.

Organización de los datos

Cache

cachegrind

Técnicas mas usuales para mejorar el uso del cache



-
- Transposición del arreglo. Similar al intergambio de loop, se puede transponer una matriz (por ejemplo para la multiplicación) y mejorar el acceso desde el arreglo de los datos.
- Copia de datos. Similar al problema del primer punto. En algoritmos por bloques (como la transpuesta), como usualmente RAM contigua mapea al mismo bloque de cache, puede ser que dos partes del mismo bloque definido en la RAM estén mapeando al mismo bloque de cache produciendo misses, una posible solución es tener una copia del bloque y usar una parte de uno y una parte de otro.

Cache

cachegrind

Técnicas mas
usuales para
mejorar el uso del
cache



Terminamos.

¡Gracias!