

Trabajo Práctico Nro 2

Organización de Computadoras

Lambre, Juan Manuel 95978
juanmlambre@gmail.com
Israel, Pablo 95849
pabloisrael94@gmail.com

Resumen

El siguiente informe explica y detalla el trabajo práctico número 2 de la materia.

Índice

1. Objetivo	2
2. Introducción	2
2.1. Memoria Cache	2
2.1.1. Política de ubicación	2
2.1.2. Política de reemplazo	2
3. Código	3
3.1. Diseño del código	3
3.2. Especificaciones del cache	3
3.3. Repositorio	4
3.4. Instructivo de instalación	4
3.4.1. Cómo compilar el código	4
3.4.2. Cómo utilizar el programa	4
4. Casos de prueba	4
4.1. Prueba 1	4
4.2. Prueba 2	5
4.3. Prueba 3	6
4.4. Prueba 4	7
4.5. Prueba 5	8
5. Conclusiones	8
6. Anexo - Código fuente	9

1. Objetivo

Familiarizarse con el funcionamiento de la memoria cache implementando una simulación de una cache dada.

2. Introducción

2.1. Memoria Cache

Una caché es un componente de hardware o software que almacena datos para que las solicitudes futuras de esos datos se puedan atender con mayor rapidez, los datos almacenados en un caché pueden ser el resultado de un cálculo anterior o el duplicado de datos almacenados en otro lugar, generalmente, de velocidad de acceso más rápido. Se produce un acierto de caché cuando los datos solicitados se pueden encontrar en esta, mientras que un error de caché ocurre cuando no están dichos datos y se debe ir a buscar el dato a memoria principal. La lectura de la caché es más rápido que volver a calcular un resultado o leer desde un almacén de datos más lento. Por lo tanto, cuantas más solicitudes se puedan atender desde la memoria caché, más rápido funcionará el sistema.

2.1.1. Política de ubicación

Decide dónde debe colocarse un bloque de memoria principal que entra en la memoria caché. Las más utilizadas son:

Directa Al bloque i -ésimo de memoria principal le corresponde la posición i módulo n , donde n es el número de bloques de la memoria caché. Cada bloque de la memoria principal tiene su posición en la caché y siempre en el mismo sitio. Su inconveniente es que cada bloque tiene asignada una posición fija en la memoria caché y ante continuas referencias a palabras de dos bloques con la misma localización en caché, hay continuos fallos habiendo sitio libre en la caché.

Asociativa Los bloques de la memoria principal se alojan en cualquier bloque de la memoria caché, comprobando solamente la etiqueta de todos y cada uno de los bloques para verificar acierto. Su principal inconveniente es la cantidad de comparaciones que realiza.

Asociativa por conjuntos Cada bloque de la memoria principal tiene asignado un conjunto de la caché, pero se puede ubicar en cualquiera de los bloques que pertenecen a dicho conjunto. Ello permite mayor flexibilidad que la correspondencia directa y menor cantidad de comparaciones que la totalmente asociativa.

2.1.2. Política de reemplazo

Aleatoria el bloque es reemplazado de forma aleatoria.

Fifo se usa el algoritmo First In First Out (FIFO) (primero en entrar primero en salir) para determinar qué bloque debe abandonar la caché. Este algoritmo generalmente es poco eficiente.

Usado menos recientemente (LRU) Sustituye el bloque que hace más tiempo que no se ha usado en la caché, traeremos a caché el bloque en cuestión y lo modificaremos ahí.

Usado con menor frecuencia (LFU) Sustituye el bloque que ha experimentado menos referencias.

3. Código

3.1. Diseño del código

La memoria a simular es una cache asociativa por conjuntos de cuatro vías, de 4KB de capacidad, bloques de 64 bytes, política de reemplazo LRU y política de escritura WB/WA. Se asume que el espacio de direcciones es de 16 bits, y hay entonces una memoria principal a simular con un tamaño de 64KB. Estas memorias pueden ser implementadas como variables globales. Cada bloque de la memoria cache deber a contar con su metadata, incluyendo el bit D, el tag, y un campo que permita implementar la política de LRU.

El diseño del código fue realizado orientado a pseudo objetos mediante la utilización de structs dado que C no soporta el uso de objetos. Se modelo el dominio del problema con las siguientes entidades:

- **Command:** Esta entidad es la encargada de parsear las instrucciones que provienen del archivo de entrada.
- **Cache:** Esta entidad es el cache que simulamos, el cual contiene cuatro vías (cuatros vectores), un contador de la cantidad de accesos, otro de cantidad de misses y por último una referencia a la memoria para que pueda acceder a ella cuando lo necesite.
- **Block:** Esta entidad es simplemente un bloque de la memoria cache. Cuando se tiene que ir a buscar un dato al cache se busca mediante el índice en las 4 vías y se obtiene el bloque el cual es un vector y para obtener finalmente el dato en cuestión se utiliza el offset.
- **Memory:** Esta entidad simula la memoria principal a la cuál se accede directamente con la dirección del dato. También contiene un vector para almacenar la información.

3.2. Especificaciones del cache

Las especificaciones del cache se definen en el archivo *config.h* como macros del preprocesador. Estas son:

- Tamaño de cache: 4KB
- Tamaño de bloque: 64B
- Cantidad de vías: 4
- Tamaño de la memoria ppal: 64KB

De estas especificaciones se desprenden algunas propiedades tenidas en cuenta:

- Bits de dirección: 16
- Bits para el tag: 6
- Bits para el índice: 4
- Bits para el offset: 6

3.3. Repositorio

<https://github.com/JuanmaLambre/6620-tp2>

3.4. Instructivo de instalación

3.4.1. Cómo compilar el código

Ejecutar

```
$ bash build
```

Esto compila el código y genera el ejecutable *cache*. Se puede agregar el parámetro *debug* para compilar con flags de debugging

3.4.2. Cómo utilizar el programa

Uso:

```
$ ./cache direccion/del/archivo.mem
```

4. Casos de prueba

4.1. Prueba 1

```
1 W 0, 255
2 W 1024, 254
3 W 2048, 248
4 W 4096, 096
5 W 8192, 192
6 R 0
7 R 1024
8 R 2048
9 R 8192
10 MR
```

prueba1.mem

```
1 ~/6620-tp2 $ ./cache test/inputs/prueba1.mem
2 255
3 254
4 248
5 96
6 192
7 255
8 254
9 248
10 192
11 MissRate 88%
12 ~/6620-tp2 $
```

Output para prueba1.mem

Para la ejecución del programa con prueba1.mem, todos los comandos de *write* son misses compulsivos, dado que cada bloque es de 64B. Además, para todas las direcciones los bits de index (6-9) son 0, por lo que los bloques se almacenarán en el mismo set. Cuando el programa llegue a la escritura en 8192 el set 0 ya estará lleno, y el último bloque usado sera el bloque 0 en la vía 0, que será reemplazado.

Debido a este reemplazo, la primera lectura dará miss, pero reemplazará al último bloque usado, que será el de la vía 1 (bloque que comienza en 1024). Consecuentemente, la siguiente lectura dará miss, y reemplazará al bloque de la dirección 2048 en la vía 2. Lo mismo sucederá con la lectura de dirección 2048, que se ubicará en la vía 3. Sin embargo, la última lectura de la dirección 8192 dará hit dado que se encontrará en la vía 0.

En total se tienen entonces 9 accesos de los cuales 8 son misses, lo que resulta en un 88 % de miss rate.

4.2. Prueba 2

```
1 R 0
2 R 31
3 W 64, 10
4 R 64
5 W 64, 20
6 R 64
7 MR
```

prueba2.mem

```
1 ~/6620-tp2 $ ./cache test/inputs/prueba2.mem
2 0
3 0
4 10
5 10
6 20
7 20
8 MissRate 33%
9 ~/6620-tp2 $
```

Output para prueba2.mem

La primer lectura en la dirección 0 será un miss compulsivo, pero traerá a cache los primeros 64B de memoria. Por esta razón, la siguiente lectura en la dirección 31 será hit. Ambas devuelven 0 porque así se inicializa la memoria.

La primera escritura es en la dirección 64, por lo que será otro miss compulsivo. Sin embargo, todas las siguientes operaciones se hacen sobre esta misma dirección, por lo que serán hits.

En definitiva, hay 2 misses en 6 accesos, resultando en un 33 % de miss rate.

4.3. Prueba 3

```
1 W 128, 1
2 W 129, 2
3 W 130, 3
4 W 131, 4
5 R 1152
6 R 2176
7 R 3200
8 R 4224
9 R 128
10 R 129
11 R 130
12 R 131
13 MR
```

prueba3.mem

```
1 ~/6620-tp2 $ ./cache test/inputs/prueba3.mem
2 1
3 2
4 3
5 4
6 0
7 0
8 0
9 0
10 1
11 2
12 3
13 4
14 MissRate 50%
15 ~/6620-tp2 $
```

Output para prueba3.mem

En esta prueba todas las direcciones comparten en mismo set (index=2), por lo que compiten constantemente por el espacio en la cache.

La primera escritura es un miss compulsivo y se almacena en la vía 0. Sin embargo las siguientes tres perteneces al mismo bloque y son hit.

Las primeras cuatro lecturas son misses compulsivos también, y se van almacenando en las vías libres. Para la lectura en la dirección 4224, todas las vías están ocupadas, y la menos recientemente usada será la primera de todas (vía 0).

Cuando llega la lectura de la dirección 128, el bloque no estará presente en la vía 0 porque acaba de ser reemplazada, lo que provoca un miss. Pero, al igual que al principio, las siguientes y últimas tres lecturas perteneces al mismo bloque, resultando en hit.

En total, de 12 accesos a memoria hay 6 que son hit, y el miss rate es de 50 %

4.4. Prueba 4

```
1 W 0, 256
2 W 1, 2
3 W 2, 3
4 W 3, 4
5 W 4, 5
6 R 0
7 R 1
8 R 2
9 R 3
10 R 4
11 R 4096
12 R 8192
13 R 0
14 R 1
15 R 2
16 R 3
17 R 4
18 MR
```

prueba4.mem

```
1 ~/6620-tp2 $ ./cache test/inputs/prueba4.mem
2 0
3 2
4 3
5 4
6 5
7 0
8 2
9 3
10 4
11 5
12 0
13 0
14 0
15 2
16 3
17 4
18 5
19 MissRate 17%
20 ~/6620-tp2 $
```

Output para prueba4.mem

La primera escritura de la prueba 4 inevitablemente es un miss compulsivo, pero tanto las escrituras que le siguen como las primeras 5 lecturas se realizan sobre el mismo bloque, por lo que son hits.

Las lecturas sobre las direcciones 4096 y 8192 son misses compulsivos y se guardan en el mismo conjunto que el bloque 0, pero habiendo 4 vías disponibles no se reemplaza ningún bloque en cache. Esto permite que las últimas cinco lecturas sean hits, dado que el bloque 0 ya estaba almacenado en cache.

Se tienen entonces 17 lecturas, de las cuales sólo 3 son misses, por lo que el miss rate es de 17%.

4.5. Prueba 5

```
1 R 131072
2 R 4096
3 R 8192
4 R 4096
5 R 0
6 R 4096
7 MR
```

prueba5.mem

```
1 ~/6620-tp2 $ ./cache test/inputs/prueba5.mem
2 Address 131072 overflows
3 0
4 0
5 0
6 0
7 0
8 MissRate 60%
9 ~/6620-tp2 $
```

Output para prueba5.mem

En estas pruebas, la primera lectura es en una dirección que excede el tamaño de la memoria principal, por lo que se imprime por stderr un mensaje de error y se descarta el comando. Vale aclarar que no cuenta como un acceso a memoria.

Las direcciones de los comandos siguientes son sólo tres, por lo que entran en la cache aunque pertenezcan al mismo conjunto (cosa que sucede). Los misses a las direcciones 0, 4096 y 8192 son compulsivos, y las demás lecturas son a la dirección 4096, que está en cache.

En definitiva, hay 5 accesos a memoria de los cuales sólo 2 son hits. El miss rate para estas pruebas es de 60 %

5. Conclusiones

El trabajo realizado nos permitió conocer y estudiar en un nivel práctico como funciona una memoria cache y el proceso necesario para obtener la información requerida. Pudimos reproducir todo el proceso que se realiza desde que llega una dirección de memoria para obtener hasta obtenerlo pasando por el cache y yendo a memoria principal de ser necesario (Miss).

Al ser una simulación en vez de ser un cache real, no se obtiene el resultado que uno busca cuando agrega una memoria cache el cuál es reducir los tiempos significativamente. En esta simulación tanto cuando se va a buscar un dato al cache o cuando se va a buscar a memoria principal siempre accede al mismo lugar físico (no hay dos tipos de memorias). Y si bien existe un tiempo de miss penalty cuando se quiere traer un bloque de memoria (dado que se copian los datos), la simulación no imita el proceso físico que el hardware realiza. Es por esto que la simulación no es fiel a la hora de medir tiempos.

Lo que sí permite la simulación es poder comparar los miss rates de diferentes arquitecturas que tienen especificaciones diferentes, aunque dicha medición y comparación escape al alcance del trabajo.

6. Anexo - Código fuente

```
1  #include <stdlib.h>
2  #include <stdio.h>
3
4  #include "Cache.c"
5  #include "Memory.c"
6  #include "Command.c"
7
8
9  char* read_line(FILE* input) {
10     char* str = (char*) malloc(sizeof(char));
11     *str = '\0';
12     int len = 0;
13     char c;
14
15     if (feof(input)) return NULL;
16
17     while (EOF != (c=fgetc(input)) && c != '\n') {
18         ++len;
19         str = (char*) realloc(str, sizeof(char)*(len+1));
20         str[len-1] = c;
21         str[len] = '\0';
22     }
23
24     if (c == EOF && len == 0) {
25         free(str);
26         return NULL;
27     }
28
29     return str;
30 }
31
32 int load_commands(FILE* input, Command** ptr) {
33     Command* cmds = NULL;
34     int count = 0;
35     char* line = NULL;
36
37     while (line = read_line(input)) {
38         ++count;
39         cmds = (Command*) realloc(cmds, sizeof(Command)*count);
40         Command_init(&cmds[count-1], line);
41         free(line);
42     }
43
44     *ptr = cmds;
45     return count;
46 }
47
48
49 int main(int argc, char** argv) {
50     if (argc < 2) {
51         fprintf(stderr, "No file specified\n");
52         return 1;
53     }
54
55     init();
56
57     // Get commands
58     Command* commands;
59     int cmd_count = 0;
```

```

60 FILE* input = fopen(argv[1], "r");
61 if (0) {
62     fprintf(stderr, "Error opening file %s", argv[1]);
63     return 1;
64 }
65 cmd_count = load_commands(input, &commands);
66
67 // Execute commands
68 for (int i = 0; i < cmd_count; ++i) {
69     Command_execute(commands+i);
70 }
71
72 for (int i = 0; i < cmd_count; ++i) {
73     Command_destroy(commands+i);
74 }
75
76 return 0;
77 }

```

src/main.c

```

1  #ifndef COMMAND_C
2  #define COMMAND_C
3
4  #include <string.h>
5
6  typedef unsigned int uint;
7  typedef unsigned char uchar;
8
9  typedef struct {
10     char* op;
11     uint* argv;
12     int argc;
13 } Command;
14
15
16 int is_eoarg_char(char c) {
17     return c == ' ' || c == '\n' || c == '\t' || c == ',';
18 }
19
20
21 void Command_init(Command* self, char* cmdstr) {
22     int start = 0, end = 0, cmdlen = strlen(cmdstr);
23
24     // Default values
25     self->op = NULL;
26     self->argv = NULL;
27     self->argc = 0;
28
29     // Get the operation
30     int eoarg = 0;
31     while (end < cmdlen && !eoarg) {
32         eoarg = is_eoarg_char(cmdstr[++end]);
33     }
34     self->op = malloc(end-start+1);
35     strncpy(self->op, cmdstr, end-start);
36     self->op[end-start] = '\0';
37
38     // Now we get the arguments
39     while (end < cmdlen) {
40         start = end;
41         while (start < cmdlen && eoarg) {
42             eoarg = is_eoarg_char(cmdstr[++start]);
43         }
44         end = start;
45         while (end < cmdlen && !eoarg) {
46             eoarg = is_eoarg_char(cmdstr[++end]);
47         }
48         if (start < cmdlen) {
49             ++self->argc;
50             self->argv = (uint*) realloc(self->argv, sizeof(uint)*self->argc);
51             char buffer[16];
52             memset(buffer, '\0', 16);
53             strncpy(buffer, cmdstr+start, end-start);
54             self->argv[self->argc-1] = atoi(buffer);
55         }
56     }
57 }
58
59 void Command_execute(Command* self) {
60     if (strcmp(self->op, "R") == 0) {
61         int value = (uchar) read_byte(self->argv[0]);

```

```

62         if (value >= 0)
63             printf("%d\n", value);
64     } else if (strcmp(self->op, "W") == 0) {
65         int res = write_byte(self->argv[0], self->argv[1]);
66         printf("%d\n", res);
67     } else if (strcmp(self->op, "MR") == 0) {
68         printf("MissRate %d%%\n", get_miss_rate());
69     } else {
70         fprintf(stderr, "Unknown command: %s\n", self->op);
71     }
72 }
73
74 void Command_destroy(Command* self) {
75     free(self->op);
76     if (self->argc > 0)
77         free(self->argv);
78 }
79
80
81 #endif

```

src/Command.c

```

1  #include "config.h"
2  #include "Block.c"
3  #include "Memory.c"
4
5  #include "math.h"
6
7  #define NO_BLOCKS (CACHE_SIZE/NO_WAYS/BLOCK_SIZE)
8  #define INDEX_SIZE (log(NO_BLOCKS)/log(2))
9  #define OFFSET_SIZE (log(BLOCK_SIZE)/log(2))
10 #define TAG_SIZE ((int)(log(MEMORY_SIZE)/log(2))-INDEX_SIZE-OFFSET_SIZE)
11
12
13 typedef unsigned char uchar;
14
15 typedef struct {
16     Memory memory;
17     Block blocks[NO_WAYS][NO_BLOCKS];
18     int access_count;
19     int misses_count;
20 } Cache;
21
22 Cache self;
23
24
25 void init() {
26     self.access_count = 0;
27     self.misses_count = 0;
28     for (int w = 0; w < NO_WAYS; ++w) {
29         for (int b = 0; b < NO_BLOCKS; ++b) {
30             Block_init(&self.blocks[w][b]);
31         }
32     }
33     Memory_init(&self.memory);
34 }
35
36 int get_index(int addr) {
37     return (int)(addr / pow(2, OFFSET_SIZE)) % (int)(pow(2, INDEX_SIZE));
38 }
39
40 int get_offset(int addr) {
41     return addr % BLOCK_SIZE;
42 }
43
44 int get_tag(int addr) {
45     return (int)(addr / pow(2, INDEX_SIZE + OFFSET_SIZE));
46 }
47
48 void read_block(int blocknum) {
49     // Doesnt check if block is dirty
50     int tag = get_tag(blocknum * BLOCK_SIZE);
51     int set = get_index(blocknum * BLOCK_SIZE);
52     int way = find_lru(set);
53     Block_read(&self.blocks[way][set], &self.memory.blocks[blocknum], tag);
54 }
55
56 Block* cache_address(int addr) {
57     int set = find_set(addr);
58     int way = find_lru(set);
59     if (is_valid(way, set) && is_dirty(way, set)) {
60         write_block(way, set);
61     }

```

```

62
63     read_block(get_blocknum(addr));
64     return &self.blocks[way][set];
65 }
66
67 int read_byte(int address) {
68     if (address >= MEMORY_SIZE) {
69         fprintf(stderr, "Address %d overflows\n", address);
70         return -1;
71     }
72     int tag = get_tag(address);
73     int offset = get_offset(address);
74     ++self.access_count;
75     int set = find_set(address);
76     for (int way = 0; way < NO_WAYS; ++way) {
77         Block block = self.blocks[way][set];
78         if (block.valid && block.tag == tag) {
79             Block_update_lru(&block);
80             return block.data[offset];
81         }
82     }
83     // Miss in cache
84     ++self.misses_count;
85     Block *b = cache_address(address); // This updates lru
86     return b->data[offset];
87 }
88
89 int find_lru(int setnum) {
90     int minway = 0;
91     for (int w = 0; w < NO_WAYS; ++w) {
92         Block b = self.blocks[w][setnum];
93         if (!b.valid)
94             return w;
95         else if (b.last_used < self.blocks[minway][setnum].last_used)
96             minway = w;
97     }
98     return minway;
99 }
100
101 int is_valid(int way, int set) {
102     return self.blocks[way][set].valid;
103 }
104
105 int is_dirty(int way, int set) {
106     return self.blocks[way][set].bit_D;
107 }
108
109 int get_blocknum(int address) {
110     return (int)(address/BLOCK_SIZE);
111 }
112
113 void write_block(int way, int set) {
114     Block block = self.blocks[way][set];
115     int block_num = (int)(block.tag * pow(2, INDEX_SIZE) + set);
116     Memory_write_block(&self.memory, &block, block_num);
117 }
118
119 int write_byte(int address, char value) {
120     // Write-Allocate politic => update cache to write
121     if (address >= MEMORY_SIZE) {
122         fprintf(stderr, "Address %d overflows\n", address);
123         return -1;

```

```

124     }
125     Block* block = NULL;
126     int tag = get_tag(address);
127     int offset = get_offset(address);
128     int set = find_set(address);
129     int way = 0;
130     ++self.access_count;
131
132     for (; way < NO_WAYS && !block; ++way) {
133         Block *b = &self.blocks[way][set];
134         if (b->valid && b->tag == tag)
135             block = b;
136     }
137
138     if (!block) {
139         ++self.misses_count;
140         block = cache_address(address);
141     }
142
143     block->data[offset] = value;
144     block->bit_D = 1;
145     Block_update_lru(&block);
146     return (unsigned char) value;
147 }
148
149 int find_set(int addr) {
150     return get_index(addr);
151 }
152
153 int get_miss_rate() {
154     return (int)(100.0*self.misses_count/self.access_count);
155 }

```

src/Cache.c

```

1 #define CACHE_SIZE (4*1024)
2 #define NO_WAYS 4
3 #define BLOCK_SIZE 64
4 #define MEMORY_SIZE (64*1024)

```

src/config.h

```

1  #include "config.h"
2
3  int LRU_COUNT = 1;
4
5  typedef struct {
6      int last_used;
7      char valid;
8      char data[BLOCK_SIZE];
9      int tag;
10     char bit_D;
11 } Block;
12
13 void Block_init(Block *self) {
14     self->last_used = 0;
15     self->tag = -1;
16     self->bit_D = 0;
17     self->valid = 0;
18     for (int i = 0; i < BLOCK_SIZE; ++i) {
19         self->data[i] = 0;
20     }
21 }
22
23 void Block_update_lru(Block *self) {
24     self->last_used = LRU_COUNT++;
25 }
26
27 void Block_read(Block *self, Block *src, int tag) {
28     Block_init(self);
29     self->valid = 1;
30     Block_update_lru(self);
31     self->tag = tag;
32     memcpy(self->data, src->data, BLOCK_SIZE);
33 }

```

src/Block.c

```

1  #ifndef MEMORY_C
2  #define MEMORY_C
3
4  #include "config.h"
5
6  #define BLOCKS_COUNT MEMORY_SIZE/BLOCK_SIZE
7
8  typedef struct {
9      Block blocks[BLOCKS_COUNT];
10 } Memory;
11
12 void Memory_init(Memory *self) {
13     for (int b = 0; b < BLOCKS_COUNT; ++b) {
14         Block_init(&self->blocks[b]);
15     }
16 }
17
18 void Memory_write_block(Memory *self, Block *src, int blockNo) {
19     Block *memblock = self->blocks + blockNo;
20     memcpy(memblock->data, src->data, BLOCK_SIZE);
21 }
22
23 #endif

```

src/Memory.c