# Project 2: Understanding Cache Memories

519021911174, Haoyuan Tian, thy0210@sjtu.edu.cn

June 3, 2021

## 1 Introduction

The objective of the project is to help us develop an understanding of how the cache memory works and how it affects the performance of our programs. The first part of the project requires writing a simple cache simulator. The second part requires optimizing the performance of a matrix transpose program, where the main target is to minimize the number of cache misses.

## 2 Experiments

### 2.1 Part A

#### 2.1.1 Analysis

In **Part A**, we are required to write a C program to simulate the mechanism of the cache memory. The program takes input from the *.trace* file and the three parameters, respectively S, E and B are read from the *valgrind* command. The replacement strategy is LRU, and the program output the number of hits, misses, and evictions.
The main difficulties here are listed below:

1. Analyze the input of instructions properly.

2. Arrange the miss and hit situation properly on the basis.

#### 2.1.2 Functions in detail

- *getopt()*

**Prototype:**

```
1   int getopt(int argc, char * const argv[], const char * optstring);
```

As shown in the complete code below, the function is used to analyze command line arguments.

- *readCache()*

The function is used to imitate the behaviour of the cache memory. According to the reference, we divide the situations that the cache may encounter into 5 cases, each corresponding to a different number of hits, misses and evictions.

### 2.1.3 Code

```
1  // Haoyuan Tian, 519021911174
2
3  #include "cachelab.h"
4  #include <getopt.h>
5  #include <unistd.h>
6  #include <math.h>
7  #include <limits.h>
8  #include <memory.h>
9  #include <stdio.h>
10 #include <stdlib.h>
11 #include <assert.h>
12 #include <time.h>
13 #include <math.h>
14
15 typedef struct
16 {
17     int valid;
18     long address;
19     int time;
20 } Cache;
21
22 Cache *cache;
23
24 int S, E, B;
25 int hits = 0;
26 int misses = 0;
27 int evictions = 0;
28
29 int readCache(Cache *cache, char type, long address)
30 {
31     int flag = 0;
32     int line = (int)((address / B) % S);
33     int pos = 0;
34     int minTime = INT_MAX;
35     int current = -1;
36     int exist = 0;
37     int insert = 0;
38     for (int i = 0; i < E; i++)
39     {
40         Cache *element = (cache + line * E + i);
41         if (insert == 0 && flag == 0)
42         {
43             if (!element->valid)
44             {
45                 element->valid = 1;
46                 element->address = address;
47                 insert = 1;
48                 exist = 0;
49                 flag = 2;
50                 pos = i;
51                 if (type == 'M')
52                 {
53                     flag = 3;
```

```
54                    }
55                }
56                else
57                {
58                    if ((element->address / B) == (address / B))
59                    {
60                        exist = 1;
61                        pos = i;
62                        flag = 1;
63                    }
64                    if (element->time < minTime)
65                    {
66                        minTime = element->time;
67                        if (exist == 0)
68                        {
69                            pos = i;
70                        }
71                    }
72                }
73            }
74        if (element->time > current)
75            current = element->time;
76    }
77    Cache *element = (cache + line * E + pos);
78    if (insert == 0 && exist == 0)
79    {
80        element->address = address;
81        flag = 4;
82        if (type == 'M')
83        {
84            flag = 5;
85        }
86    }
87    element->time = current + 1;
88    return flag;
89 }
90
91 void printHelp()
92 {
93 printf("Usage: ./csim [-hv] -s <s> -E <E> -b <b> -t <tracefile>\n");
94 }
95
96 int main(int argc, char *argv[])
97 {
98     opterr = 0;
99     int verbose = 0;
100    int para;
101    char *name = NULL;
102    while ((para = getopt(argc, argv, "s:E:b:t:hv")) != -1)
103    {
104        switch (para)
105        {
106        case 's':
107            S = (int)pow(2, atoi(optarg));
108            break;
```

```c
        case 'E':
            E = atoi(optarg);
            break;
        case 'b':
            B = (int)pow(2, atoi(optarg));
            break;
        case 't':
            name = optarg;
            break;
        case 'v':
            verbose = 1;
            break;
        case 'h':
            printHelp();
            break;
        }
    }

    FILE *file = fopen(name, "r");

    if (file == NULL)
    {
        printf("NO SUCH FILE!\n");
        return -1;
    }

    cache = (Cache *)malloc(S * E * sizeof(Cache));
    if (cache == NULL)
    {
        printf("MEMORY ALLOCATION FAILED!\n");
        return -1;
    }

    memset(cache, 0, S * E * sizeof(Cache));
    char type;
    int size;
    long address;
    while (!feof(file))
    {
        int tmp = fscanf(file, " %c %lx,%x", &type, &address, &size);
        if (tmp != 3)
            continue;
        if (type != 'I')
        {
            int ret = readCache(cache, type, address);
            char *tmp1 = NULL;
            switch (ret)
            {
            case 1:
                hits++;
                tmp1 = "hit";
                if (type == 'M')
                    hits++;
                break;
            case 2:
```

```c
164                     misses++;
165                     tmp1 = "miss";
166                     break;
167
168                 case 3:
169                     misses++;
170                     hits++;
171                     tmp1 = "miss";
172                     break;
173                 case 4:
174                     misses++;
175                     evictions++;
176                     tmp1 = "miss";
177                     break;
178                 case 5:
179                 default:
180                     misses++;
181                     evictions++;
182                     hits++;
183                     tmp1 = "miss";
184                     break;
185                 }
186                 if (verbose)
187                 {
188                     printf("%c %lx,%x %s", type, address, size, tmp1);
189                     if (ret == 4)
190                     {
191                         printf(" eviction");
192                     }
193                     else if (ret == 3)
194                     {
195                         printf(" hit");
196                     }
197                     else if (ret == 5)
198                     {
199                         printf(" eviction hit");
200                     }
201                     else if (ret == 1 && type == 'M')
202                     {
203                         printf(" hit");
204                     }
205
206                     printf("\n");
207                 }
208             }
209         }
210     free(cache);
211     printSummary(hits, misses, evictions);
212 }
```

### 2.1.4 Evaluation



Figure 1: Use *make* to compile the project

In the project, the warnings will be considered errors when compiling the project. **Figure 1** shows both **Part A** and **Part B** are free of warnings.



Figure 2: Use *./test-csim* to test the result

Figure 2 shows that the simulator we realized produces exactly the same results as the reference simulator. It reaches the full mark which is 27 here.



Figure 3: Test the *verbose* output

The verbose output displays the hits, misses, and evictions that occur as a result of each memory access. The feature is not required but extremely helpful when debugging.

## 2.2 Part B

### 2.2.1 Analysis

In **Part B**, we are required to write a matrix-transpose function that causes as few cache misses as possible. The most crucial part here is understanding how cache handles the elements inside a matrix. In simple term, when a cache miss happens, the particular element will be stored into cache together with the surrounding elements on the same row. Due to the given specifications of cache and the different matrix sizes, we have to make special arrangement for each of the three conditions.

- M=32, N=32
  According to the cache size, each line of the cache can store 8 integers, which means that the cache will be filled with every 8 rows of the matrix A. Therefore, we can divide B into $8 \times 8$ matrices. Arrange the elements on the diagonal line will further reduce the misses.

- M=64, N=64
  The idea used here will be discussed later in **3.1 Problems**.

- M=61, N=67
  Neither M nor N is 2 to the power. However, the criteria is not strict here. Dividing A into $16 \times 16$ matrices is a feasible scheme.

### 2.2.2 Code

```c
char transpose_submit_desc[] = "Transpose submission";
void transpose_submit(int M, int N, int A[N][M], int B[M][N])
{
    int i, j, n, m;
    int tmp, tmp1, tmp2, tmp3, tmp4, tmp5, tmp6, tmp7;
    if (M == 32)
    {
        for (n = 0; n < N; n += 8)
        {
            for (m = 0; m < M; m += 8)
            {
                for (i = n; i < n + 8; i++)
                {
                    for (j = m; j < m + 8; j++)
                    {
                        if (i != j)
                        {
                            B[j][i] = A[i][j];
                        }
                        else
                        {
                            tmp = A[i][j];
                        }
                    }
                    if (m == n)
                    {
                        B[i][i] = tmp;
                    }
                }
            }
        }
    }
    else if (M == 64)
```

```
     {
         for (m = 0; m < 64; m += 8)
         {
             for (n = 0; n < 64; n += 8)
             {
                 for (i = n; i < n + 4; ++i)
                 {
                     tmp = A[i][m];
                     tmp1 = A[i][m + 1];
                     tmp2 = A[i][m + 2];
                     tmp3 = A[i][m + 3];
                     tmp4 = A[i][m + 4];
                     tmp5 = A[i][m + 5];
                     tmp6 = A[i][m + 6];
                     tmp7 = A[i][m + 7];

                     B[m][i] = tmp;
                     B[m][i + 4] = tmp4;
                     B[m + 1][i] = tmp1;
                     B[m + 1][i + 4] = tmp5;
                     B[m + 2][i] = tmp2;
                     B[m + 2][i + 4] = tmp6;
                     B[m + 3][i] = tmp3;
                     B[m + 3][i + 4] = tmp7;
                 }
                 for (i = m; i < m + 4; ++i)
                 {
                     tmp = B[i][n + 4];
                     tmp1 = B[i][n + 5];
                     tmp2 = B[i][n + 6];
                     tmp3 = B[i][n + 7];
                     tmp4 = A[n + 4][i];
                     tmp5 = A[n + 5][i];
                     tmp6 = A[n + 6][i];
                     tmp7 = A[n + 7][i];


                     B[i][n + 4] = tmp4;
                     B[i][n + 5] = tmp5;
                     B[i][n + 6] = tmp6;
                     B[i][n + 7] = tmp7;
                     B[i + 4][n] = tmp;
                     B[i + 4][n + 1] = tmp1;
                     B[i + 4][n + 2] = tmp2;
                     B[i + 4][n + 3] = tmp3;

                 }
                 for (i = m + 4; i < m + 8; ++i)
                 {
                     tmp = A[n + 4][i];
                     tmp1 = A[n + 5][i];
                     tmp2 = A[n + 6][i];
                     tmp3 = A[n + 7][i];

                     B[i][n + 4] = tmp;
```

```
89                        B[i][n + 5] = tmp1;
90                        B[i][n + 6] = tmp2;
91                        B[i][n + 7] = tmp3;
92                    }
93                }
94            }
95        }
96        else
97        {
98            for (n = 0; n < N; n += 16)
99            {
100               for (m = 0; m < M; m += 16)
101               {
102                   for (i = n; (i < n + 16) && (i < N); i++)
103                   {
104                       for (j = m; (j < m + 16) && (j < M); j++)
105                       {
106                           if (i != j)
107                           {
108                               B[j][i] = A[i][j];
109                           }
110                           else
111                           {
112                               tmp = A[i][j];
113                           }
114                       }
115                       if (m == n)
116                       {
117                           B[i][i] = tmp;
118                       }
119                   }
120               }
121           }
122       }
123 }
```

### 2.2.3 Evaluation

Since the correctness during the compile stage has been proved in **Part A**, we won't bother to show it again here.



Figure 4: M=32, N=32



Figure 5: M=64, N=64



Figure 6: M=61, N=67

1. M=32, N=32: 287 misses in total.

2. M=64, N=64: 1219 misses in total.

3. M=61, N=67: 1989 misses in total.

All reach the full mark!

# 3   Conclusion

## 3.1   Problems

### 3.1.1   Problems in Part A

The first obstacle comes from the unfamiliarity with the function *getopt()*. After reading some materials on the Internet, I manage to apply it properly.

The second obstacle is the arrangement of different cache situations. Once I figured out the logic of the cache and the impact of the instructions from the *.trace* file, all left to be done is debugging. Using the reference simulator *csim-ref* and verbose feature wisely helps a lot.

### 3.1.2   Problems in Part B

This part is particularly focused on the $M = 64, N = 64$ situation.

The cache will be filled with every 4 rows of the matrix A, so the initial idea was to divide B into $4 \times 4$ matrices. The total number of misses using such method is 1699, which is not ideal.

Next, I tried the division plan of first $8 \times 8$, then $4 \times 4$. By arranging the sequence properly, the number of misses went down to 1451. Effective, but still not satisfying enough.

The final idea was inspired by the reading material provided and blogs on the Internet. The previous ideas only focused on dividing A and B, but ignored that we could make changes to matrix B. The figure below illustrated the idea in a brief way.
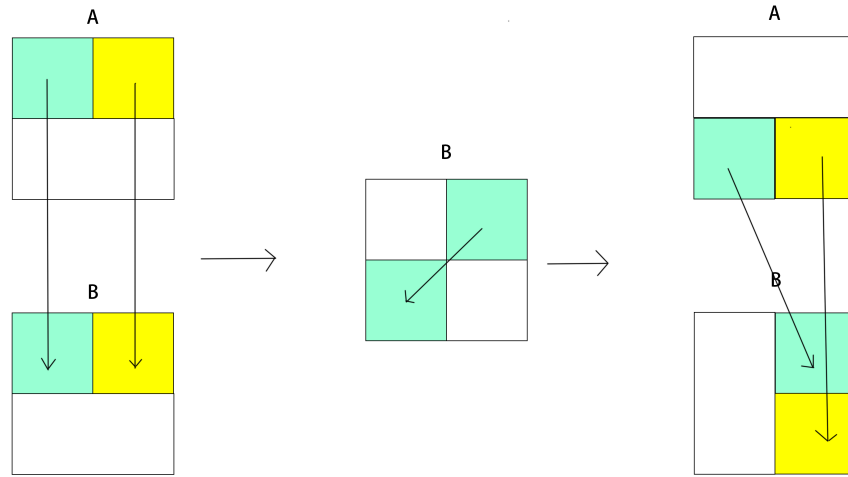


Figure 7: The final idea

The yellow block of A doesn't go to the right place in the first round of transposition, and will be dealt with later. In this way, the number of misses was eventually cut down to 1219.

## 3.2   Achievements

In **Part A**, I managed to realize a cache simulator in only 212 lines. The handling of the LRU replacement strategy might seem confusing but it worked just well. The situations were divided precisely and tackled successfully. I also paid effort to improve the robustness of the program by rejecting the wrong input and avoiding memory leaks.

In **Part B**, I strictly followed the restriction about the number of local variables. In the mean time, the variables defined locally were used wisely and none of them was redundant. The loop structure is clear and easy to understand.

## 3.3 Feelings

The project puts emphasis on cache and cache alone. It provides a golden opportunity to put knowledge learnt in the course into practice, and therefore deepen our understanding of the mechanism of cache. Without exaggeration, I can say that I've realized a significant optimization and produced a marvellous result. The process of trial and error might be miserable, especially when the result of **Part B** didn't live up to my expectation. But finally accomplishing the full mark with my effort brought more delight. Cache is a very interesting and important part of the computer architecture, and I hope that I will have the chance to explore more about it later.

Great thanks to teacher Shen and the teaching assistants for your patience and guidance.