# Project 2: Understanding Cache Memories

519021910366, Dou Yiming, douyiming@sjtu.edu.cn

June 5, 2021

## 1   Introduction

This project aims at making us understand cache memories better. The whole project can be divided into two parts.

Firstly, we are required to simulate a cache including sets, ways, tags and valid bits. It also covers cache replacement algorithm LRU that we have learned in operating system class. I complete this part with reference to slides of the course and CSAPP book.

Secondly, on the basis of the cache simulator, we are required to implement an optimization. When transposing the matrix, the cache may not play to its full capability if without enough design. Therefore, we need to improve the cache hit rate by blocking the matrix and transfer each block with the help of cache in order to speed up.

## 2   Experiments

### 2.1   Part A

#### 2.1.1   Analysis

The cache we are required to simulate contains $2^s$ sets, $E$ ways and $2^b$ block offsets, whose structure can be shown in Fig. 1.

In the process of a single entry, we begin with traversing each set, checking whether it's valid or not and whether the tag is the same as the data we want. If a qualified set is found, we define this as a cache hit. Otherwise, if the cache is full, we define this situation as eviction, in which we should swap a new element with the existing element in the cache.

The next problem is how we choose the place we want to evict. In this project, we use the Least-Recently-Used (LRU) strategy, meaning that we need to seek the data that we have not used for the longest time. Therefore, we attach a time stamp for each set, each time the set is visited, we update the time stamp of it to the current time, and when eviction happens, we choose the set whose time stamp is the smallest.
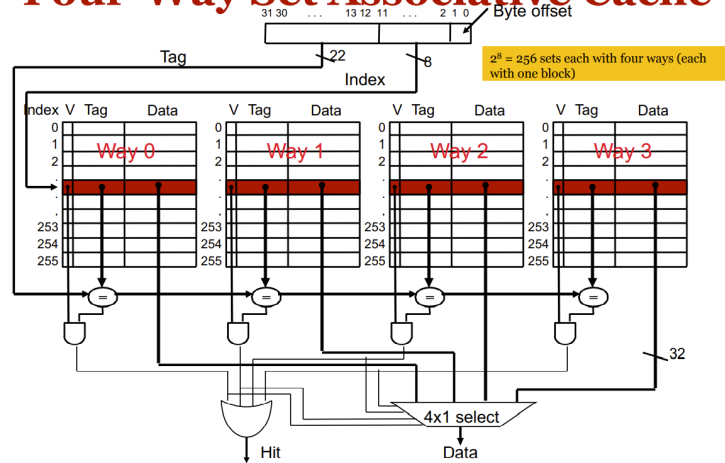
Figure 1: Cache Structure

## 2.1.2 Code

```c
/*
 * Name: DouYiming
 * ID: 519021910366
 */
#include <getopt.h>
#include <stdlib.h>
#include <stdio.h>
#include "cachelab.h"

const int m = 64;

typedef struct Arg
{
    int s; // Number of set index bits
    int E; // number of lines per set
    int b; // Number of block bits
    FILE *f; // valgrind trace to replay
} Arg;

typedef struct summary
{
    int hits;
    int misses;
    int evictions;
} summary;

typedef struct Line
{
    int valid;
    int tag;
    int time_stamp;
```

2

```c
32  } Line;
33
34  typedef struct Set
35  {
36      Line* lines;
37  } Set;
38
39  Arg arg_parser(int argc, char ** argv)
40  {
41      Arg arg;
42      if(getopt(argc, argv, "s:E:b:t:") == 's')
43          arg.s = atoi(optarg);
44      if(getopt(argc, argv, "s:E:b:t:") == 'E')
45          arg.E = atoi(optarg);
46      if(getopt(argc, argv, "s:E:b:t:") == 'b')
47          arg.b = atoi(optarg);
48      if(getopt(argc, argv, "s:E:b:t:") == 't')
49          arg.f = fopen(optarg, "r");
50      return arg;
51  }
52
53  Set* init_cache(Arg arg)
54  {
55      int set_num = 1 << arg.s;
56      Set* caches=(Set*)malloc(set_num * sizeof(Set));
57      for(int i = 0;i < set_num; ++i)
58      {
59          caches[i].lines = (Line*)malloc(arg.E * sizeof(Line));
60          for(int j = 0; j < arg.E; ++j)
61          {
62              caches[i].lines[j].valid = 0;
63              caches[i].lines[j].time_stamp = 0;
64          }
65      }
66      return caches;
67  }
68
69  summary simulator(Arg arg, Set* caches)
70  {
71      // information to display
72      summary ans;
73      ans.hits = 0;
74      ans.misses = 0;
75      ans.evictions = 0;
76      // compute tag bits
77      int t = m - arg.b - arg.s;
78      // simulator time
79      int cur_time = 0;
80      // information of each instruction
81      char op;
82      long addr;
83      int len;
84      int cycle_time;
85
86      // read from file
87      while(!feof(arg.f))
88      {
```

```c
89          if (fscanf(arg.f, " %c %lx,%d", &op, &addr, &len) != 3)
90              continue;
91
92          if (op == 'L' || op == 'S')
93              cycle_time = 1;
94          else if (op == 'M')
95              cycle_time = 2;
96          else
97              continue;
98          while(cycle_time--)
99          {
100             ++cur_time;
101             int set_num = (addr >> arg.b) & ((1 << arg.s) - 1);
102             int tag = (addr >> (arg.s + arg.b)) & ((1 << t) - 1);
103             int is_hit = 0;
104             int invalid_line = -1;
105             int LRU_line = -1;
106             int min_time = 0x7fffffff;
107             // search each way
108             for(int i = 0; i < arg.E; ++i)
109             {
110                 Line* line = &caches[set_num].lines[i];
111                 if(line->valid && line->tag == tag) // hit
112                 {
113                     is_hit = 1;
114                     line->time_stamp = cur_time;
115                     break;
116                 }
117                 else if(!line->valid)
118                 {
119                     invalid_line = i;
120                     break;
121                 }
122                 else if(line->valid && line->time_stamp<min_time)
123                 {
124                     min_time = line->time_stamp;
125                     LRU_line = i;
126                 }
127             }
128             if(is_hit)//hit
129                 ++ans.hits;
130             else
131             {
132                 ++ans.misses;
133                 if(invalid_line != -1)//miss
134                 {
135                     Line*l=&caches[set_num].lines[invalid_line];
136                     l->tag = tag;
137                     l->time_stamp = cur_time;
138                     l->valid = 1;
139                 }
140                 else//evict
141                 {
142                     ++ans.evictions;
143                     Line*l=&caches[set_num].lines[LRU_line];
144                     l->tag = tag;
145                     l->time_stamp = cur_time;
```

```c
146                    }
147                }
148            }
149        }
150        return ans;
151 }
152
153 void terminate(Arg arg, Set* caches)
154 {
155        fclose(arg.f);
156        for (int i = 0; i < (1 << arg.s); ++i)
157            free(caches[i].lines);
158        free(caches);
159 }
160
161 int main(int argc, char **argv)
162 {
163        // arg parsing
164        Arg arg = arg_parser(argc, argv);
165        // cache initialization
166        Set* caches = init_cache(arg);
167        // simulation
168        summary s = simulator(arg, caches);
169        // ending
170        terminate(arg, caches);
171        // print ans
172        printSummary(s.hits, s.misses, s.evictions);
173        return 0;
174 }
```

```
dodo@ubuntu:~/code/Computer_Architecture/project2-handout$ ./test-csim
                        Your simulator      Reference simulator
Points (s,E,b)    Hits  Misses  Evicts    Hits  Misses  Evicts
      3 (1,1,1)      9       8       6       9       8       6  traces/yi2.trace
      3 (4,2,4)      4       5       2       4       5       2  traces/yi.trace
      3 (2,1,4)      2       3       1       2       3       1  traces/dave.trace
      3 (2,1,3)    167      71      67     167      71      67  traces/trans.trace
      3 (2,2,3)    201      37      29     201      37      29  traces/trans.trace
      3 (2,4,3)    212      26      10     212      26      10  traces/trans.trace
      3 (5,1,5)    231       7       0     231       7       0  traces/trans.trace
      6 (5,1,5) 265189   21775   21743  265189   21775   21743  traces/long.trace
     27

TEST_CSIM_RESULTS=27
```

Figure 2: The result of part A

### 2.1.3 Evaluation

The result of part A is shown in Fig. 2.

## 2.2 Part B

### 2.2.1 Analysis

In this part, we need to optimize the process of transposing matrix of different sizes.

The basic idea is that we should use the elements in the cache as much as possible for each time. Therefore, simply traversing the rows and columns of a large matrix may not be a good idea due to large amounts of evictions. In order to make full use of the cache, we separate the matrix into blocks and move a block into the cache each time, ensuring that most elements in cache are used.

Furthermore, our method to separate the matrix should vary when facing different sizes of matrix by calculating the address of elements in the matrix. Therefore, we respectively implement a optimization method for each matrix transposing situation.

### 2.2.2 Code

```
1  /*
2   * Name: Dou Yi-ming
3   * ID: 519021910366
4   * trans.c - Matrix transpose B = A^T
5   *
6   * Each transpose function must have a prototype of the form:
7   * void trans(int M, int N, int A[N][M], int B[M][N]);
8   *
9   * A transpose function is evaluated by counting the number of
         misses
10  * on a 1KB direct mapped cache with a block size of 32 bytes.
11  */
12 #include <stdio.h>
13 #include "cachelab.h"
```

6

```
14
15  int is_transpose(int M, int N, int A[N][M], int B[M][N]);
16
17  /*
18   * transpose_submit − This is the solution transpose function that
         you
19   *      will be graded on for Part B of the assignment. Do not
        change
20   *      the description string "Transpose submission", as the driver
21   *      searches for that string to identify the transpose function
        to
22   *      be graded.
23   */
24  char transpose_submit_desc[] = "Transpose submission";
25  void transpose_submit(int M, int N, int A[N][M], int B[M][N])
26  {
27      // local variables
28      int ii, jj, i, j;
29      int tmp0, tmp1, tmp2, tmp3, tmp4, tmp5, tmp6, tmp7;
30      if (M == 32 && N == 32)
31      {
32          // 8 * 8 blocks
33          for(ii = 0;ii < 4;++ii)
34              for(jj = 0;jj < 4;++jj)
35                  for(i = 0;i < 8;++i)
36                  {
37                      for(j = 0;j < 8;++j)
38                          if(i != j)
39                              B[8 * jj + j][8 * ii + i] = A[8 * ii +
        i][8 * jj + j];
40                      B[8 * jj + i][8 * ii + i] = A[8 * ii + i][8 *
        jj + i];
41                  }
42      }
43      else if (M == 64 && N == 64)
44      {
45          for(ii = 0;ii < 8;++ii)
46              for(jj = 0;jj < 8;++jj)
47              {
48                  for(i = 0;i < 4;++i)
49                  {
50                      tmp0 = A[ii * 8 + i][jj * 8 + 0];
51                      tmp1 = A[ii * 8 + i][jj * 8 + 1];
52                      tmp2 = A[ii * 8 + i][jj * 8 + 2];
53                      tmp3 = A[ii * 8 + i][jj * 8 + 3];
54                      tmp4 = A[ii * 8 + i][jj * 8 + 4];
55                      tmp5 = A[ii * 8 + i][jj * 8 + 5];
56                      tmp6 = A[ii * 8 + i][jj * 8 + 6];
57                      tmp7 = A[ii * 8 + i][jj * 8 + 7];
58
59                      B[jj * 8 + 0][ii * 8 + i] = tmp0;
60                      B[jj * 8 + 1][ii * 8 + i] = tmp1;
61                      B[jj * 8 + 2][ii * 8 + i] = tmp2;
62                      B[jj * 8 + 3][ii * 8 + i] = tmp3;
63
64                      B[jj * 8 + 0][ii * 8 + i + 4] = tmp4;
65                      B[jj * 8 + 1][ii * 8 + i + 4] = tmp5;
```

```
66              B[jj * 8 + 2][ii * 8 + i + 4] = tmp6;
67              B[jj * 8 + 3][ii * 8 + i + 4] = tmp7;
68           }
69           for(j = 0; j < 4; ++j)
70           {
71              tmp0 = B[jj * 8 + j][ii * 8 + 4];
72              tmp1 = B[jj * 8 + j][ii * 8 + 5];
73              tmp2 = B[jj * 8 + j][ii * 8 + 6];
74              tmp3 = B[jj * 8 + j][ii * 8 + 7];
75
76              tmp4 = A[ii * 8 + 4][jj * 8 + j];
77              tmp5 = A[ii * 8 + 5][jj * 8 + j];
78              tmp6 = A[ii * 8 + 6][jj * 8 + j];
79              tmp7 = A[ii * 8 + 7][jj * 8 + j];
80
81              B[jj * 8 + j][ii * 8 + 4] = tmp4;
82              B[jj * 8 + j][ii * 8 + 5] = tmp5;
83              B[jj * 8 + j][ii * 8 + 6] = tmp6;
84              B[jj * 8 + j][ii * 8 + 7] = tmp7;
85
86              tmp4 = A[ii * 8 + 4][jj * 8 + j + 4];
87              tmp5 = A[ii * 8 + 5][jj * 8 + j + 4];
88              tmp6 = A[ii * 8 + 6][jj * 8 + j + 4];
89              tmp7 = A[ii * 8 + 7][jj * 8 + j + 4];
90
91              B[jj * 8 + j + 4][ii * 8 + 0] = tmp0;
92              B[jj * 8 + j + 4][ii * 8 + 1] = tmp1;
93              B[jj * 8 + j + 4][ii * 8 + 2] = tmp2;
94              B[jj * 8 + j + 4][ii * 8 + 3] = tmp3;
95              B[jj * 8 + j + 4][ii * 8 + 4] = tmp4;
96              B[jj * 8 + j + 4][ii * 8 + 5] = tmp5;
97              B[jj * 8 + j + 4][ii * 8 + 6] = tmp6;
98              B[jj * 8 + j + 4][ii * 8 + 7] = tmp7;
99           }
100         }
101      }
102      else if (M == 61 && N == 67)
103      {
104          for(ii = 0; ii < 5; ++ii)
105              for(jj = 0; jj < 4; ++jj)
106                  for(i = 0; i < 16 && ii * 16 + i < N; ++i)
107                      for(j = 0; j < 16 && jj * 16 + j < M; ++j)
108                          B[16 * jj + j][16 * ii + i] = A[16 * ii + i
     ][16 * jj + j];
109      }
110      else
111      {
112          int tmp;
113          for (i = 0; i < N; i++)
114              for (j = 0; j < M; j++)
115              {
116                  tmp = A[i][j];
117                  B[j][i] = tmp;
118              }
119      }
120
121 }
```

```
122
123  /*
124   * You can define additional transpose functions below. We've
           defined
125   * a simple one below to help you get started.
126   */
127
128  /*
129   * trans − A simple baseline transpose function, not optimized for
           the cache.
130   */
131  char trans_desc[] = "Simple row−wise scan transpose";
132  void trans(int M, int N, int A[N][M], int B[M][N])
133  {
134      int i, j, tmp;
135
136      for (i = 0; i < N; i++) {
137          for (j = 0; j < M; j++) {
138              tmp = A[i][j];
139              B[j][i] = tmp;
140          }
141      }
142
143  }
144
145  /*
146   * registerFunctions − This function registers your transpose
147   *     functions with the driver.  At runtime, the driver will
148   *     evaluate each of the registered functions and summarize
           their
149   *     performance. This is a handy way to experiment with
           different
150   *     transpose strategies.
151   */
152  void registerFunctions()
153  {
154      /* Register your solution function */
155      registerTransFunction(transpose_submit, transpose_submit_desc);
156
157      /* Register any additional transpose functions */
158      registerTransFunction(trans, trans_desc);
159
160  }
161
162  /*
163   * is_transpose − This helper function checks if B is the transpose
            of
164   *     A. You can check the correctness of your transpose by
           calling
165   *     it before returning from the transpose function.
166   */
167  int is_transpose(int M, int N, int A[N][M], int B[M][N])
168  {
169      int i, j;
170
171      for (i = 0; i < N; i++) {
172          for (j = 0; j < M; ++j) {
```

```
173            if (A[i][j] != B[j][i]) {
174                return 0;
175            }
176        }
177    }
178    return 1;
179 }
```

```
Part B: Testing transpose function
Running ./test-trans -M 32 -N 32
Running ./test-trans -M 64 -N 64
Running ./test-trans -M 61 -N 67

Cache Lab summary:
                        Points    Max pts       Misses
Csim correctness          27.0        27
Trans perf 32x32           8.0         8          287
Trans perf 64x64           8.0         8         1187
Trans perf 61x67          10.0        10         1992
        Total points      53.0        53
```

Figure 3: Result of part B

### 2.2.3 Evaluation

The result of the test is shown in Fig. 3

# 3 Conclusion

## 3.1 Problems

1. The name of each part of the cache is different from the name we use in class and slides, which cause lots of confusion. With reference to CSAPP book, we finally solve this problem.

2. The optimization of the matrix transposing problem is far from our intuition, which is simply traversing the rows and columns. After referring to lots of materials, we eventually understand the separation method.

## 3.2 Achievements

1. My results are excellent. All of the requirements have been satisfied and i get full score in each test.

2. When designing the cache simulator, i paid a lot of attention to the structure of the code, and optimize the structure to improve readability and conciseness for a long time. Finally, i reduced the length of main function to 6 lines and i am very proud of it.

3. I figured out and successfully implemented the optimization method of transposing the matrix, which is very complicated.