# Project 2: Understanding Cache Memories

518021910609　Yifan Liu　2602843473@qq.com

## 1. Introduction

This project is aimed to help us understand the impact that cache memories can have on the performance of our C programs. It contains two parts. In part A, we will implement a cache simulator that simulates the behavior of a cache memory. In part B, we will optimize a small matrix transpose function, with the goal of minimizing the number of cache misses in some particular size of matrices.

In part A, I implement a C program that simulates the behavior of cache memory. It could correctly detect the number of misses, hits and evictions. In part B, I implement a small matrix transpose function that could decrease the misses by changing the logic when transposing.

## 2. Experiments

### 2.1　Part A

#### 2.1.1 Analysis

In Part A I write a cache simulator in *csim.c* that takes a *valgrind* memory trace as input, simulates the behavior of a cache memory, and outputs the total number of hits, misses, and evictions.

I define a struct *Scache*, which contains variables that represent the number of sets, the number of lines and many sets. In a set, there are many lines,. In a line, there are three variables: valid bit, LRU number and tag number. The structure of *Scache* is similar to the real cache.

In *main* function, the program calls *get_Opt* to get variables from command line, then calls *init_cache* to initialize the cache and open the file. Then it continuously reads the line in the trace file to get the tag and set of the current data by calling *get_tag* and *get_set* and execute the corresponding load, store or modify operation. At last, it call *printSummary* to output the total number of hits, misses and evictions.

In *get_Opt* function, the program determines the parameters by checking the values after s, E, b, and t and store them to corresponding variables. It uses *getopt* function and *optarg* to get these values. I also define a check function *checkOptarg* to check whether these values are valid.

In *init_cache* function, it first checks the exceptions to avoid some errors. Then the function uses variables *s* and *E* to set the number of sets and lines. After that, it calls *malloc* function to initial each set and each line.

When initializing each line, it sets the LRU number and valid bit to 0.

Since the cache behavior of store operation is just like load operation, and modify operation is actually one load operation and one store operation, so we just need to implement load operation *load_op*.

In *load_op* function, it calls *miss_checking* function to check whether there is a miss in current operation. If not, the number of hits plus 1. If so, the number of misses plus 1. Then it calls *cache_full* function to check whether the cache is full and need to do eviction. If so, the number of evictions plus 1.

In *miss_checking* function, it locates the target line in the set by *tag* and check the valid bit. If the valid bit is equal to 1, it means that there are some data in this address. So there does not exist a miss of this operation. Then, it needs to call *update_LRU* to update LRU number of this address. If we could not find the target data, it means that there is a miss and return 1.

In *cache_full* function, it checks whether all valid bits are 1. If so, it means that we need to replace some data with new one. Therefore, it find the smallest LRU number by calling *find_smallest_LRU*. After that, it updates the valid bit and tag of the line that has the smallest LRU number and calls *update_LRU* to update LRU number. If the cache is not full, we simply update valid bit, tag and LRU number of one empty line.

In *find_smallest_LRU* function, it traverses the whole set to find the line that has the smallest LRU number. This line is just the one needs to be replaced.

In *update_LRU* function, it sets the LRU number of current line to the maximum number (I set it to 1000) and "-1" with LRU number of other lines.

## 2.1.2 Code

```
1.  // Name: Yifan Liu
2.  // ID:   518021910609
3.  #include <getopt.h>
4.  #include <stdlib.h>
5.  #include <unistd.h>
6.  #include "cachelab.h"
7.  #include <stdio.h>
8.  #include <strings.h>
9.  #include <string.h>
10.
11. #define LRUMAX 1000
```

```c
12.
13. typedef struct{
14.     int valid;//valid bit
15.     int LRU;//lru number
16.     int tag;//tag number
17. } Line;//the struct of one line in cache
18.
19. typedef struct{
20.     Line* lines;//many lines in a set
21. } Set;//the struct of one set in cache
22.
23. typedef struct{
24.     int set_cnt;   //the number of sets
25.     int line_cnt; //the number of lines
26.     Set* sets;     //many sets
27. } Scache;//struct of simulate cache
28.
29.
30. //output variables: misses, hits and evictions
31. int miss_count;
32. int hit_count;
33. int eviction_count;
34.
35. //print help messages
36. void printHelp()
37. {
38.     printf("Usage: ./csim-ref [-hv] -s <s> -E <E> -b <b> -
    t <tracefile>\n");
39.     printf("-h: Optional help flag that prints usage info\n");
40.     printf("-v: Optional verbose flag that displays trace info\n");
41.     printf("-s <s>: Number of set index bits\n");
42.     printf("-E <E>: Associativity (number of lines per set)\n");
43.     printf("-b <b>: Number of block bits\n");
44.     printf("-t <tracefile>: Name of the valgrind trace to replay\n");
45. }
46.
47. //check whether the command is valid
48.
49. void checkOptarg(char *curOptarg){
50.     if(curOptarg[0]=='-'){
51.         printf("Wrong command format\n");
52.         printHelp();
53.         exit(0);
54.     }
```

```
55. }
56.
57.
58. //get variables from command
59.
60. int get_Opt(int argc,char **argv,int *s,int *E,int *b,char *traceName,int *V
    erbose_mode){
61.     int c;
62.     while((c = getopt(argc,argv,"hvs:E:b:t:"))!=-1)
63.     {
64.         switch(c)
65.         {
66.         case 'v':
67.             *Verbose_mode = 1;
68.             break;
69.         case 's':
70.             checkOptarg(optarg);
71.             *s = atoi(optarg);
72.             break;
73.         case 'E':
74.             checkOptarg(optarg);
75.             *E = atoi(optarg);
76.             break;
77.         case 'b':
78.             checkOptarg(optarg);
79.             *b = atoi(optarg);
80.             break;
81.         case 't':
82.             checkOptarg(optarg);
83.             strcpy(traceName,optarg);
84.             break;
85.         case 'h':
86.         default:
87.             printHelp();
88.             exit(0);
89.         }
90.     }
91.     return 1;
92. }
93.
94. //initial cache
95. void init_cache(int s,int E,int b, Scache *cache)
96. {
97.     //exceptions
```

```c
98.     if(s<0)
99.     {
100.         printf("sets number must larger than 0\n!");
101.             exit(0);
102.     }
103.
104.     cache->set_cnt=2<<s;
105.     cache->line_cnt=E;
106.     cache->sets=(Set *)malloc(cache->set_cnt * sizeof(Set));//initial cache

107.
108.     int i,j;
109.     for(i=0;i<cache->set_cnt;++i)
110.     {
111.         cache->sets[i].lines=(Line *)malloc(E*sizeof(Line));//allocate line
    space
112.
113.         for(j=0;j<E;++j)//initial lines
114.         {
115.             cache->sets[i].lines[j].valid = 0;
116.                 cache->sets[i].lines[j].LRU = 0;
117.         }
118.     }
119.     return ;
120. }
121.
122. //get set number of current address
123. int get_set(int addr,int s,int b)
124. {
125.     addr=addr>>b;
126.     int k=(1<<s)-1;
127.     return addr &k;
128. }
129.
130. //get tag number of current address
131. int get_tag(int addr,int s,int b)
132. {
133.     int x=s+b;
134.     return addr>>x;
135. }
136.
137. //update LRU number, current LRU becomes the max number, others -1
138. void update_LRU(Scache *cache, int setaddr, int current_line)
139. {
```

```c
140.        cache->sets[setaddr].lines[current_line].LRU=LRUMAX;
141.        int i;
142.        for(i=0;i<cache->line_cnt;++i)//update other LRU
143.        {
144.            if(i!=current_line)
145.            cache->sets[setaddr].lines[i].LRU--;
146.        }
147. }
148.
149. //check whether current data is missing in cache
150. int miss_checking(Scache *cache,int setaddr,int tagaddr)
151. {
152.        int miss_flag=1;
153.        int i;
154.        for(i=0;i<cache->line_cnt;++i)
155.        {
156.            if(cache->sets[setaddr].lines[i].valid==1 && cache->sets[setaddr].l
    ines[i].tag==tagaddr)//not miss
157.            {
158.                miss_flag=0;
159.                update_LRU(cache, setaddr, i);
160.            }
161.        }
162.        return miss_flag;
163. }
164.
165. //find smallest LRU to do replacement
166. int find_smallest_LRU(Scache *cache,int setaddr)
167. {
168.        int target_line=0;//the smallest LRU's index
169.        int i;
170.        int min_LRU=LRUMAX;//smallest LRU
171.        for(i=0;i<cache->line_cnt;++i)//find
172.        {
173.            if(cache->sets[setaddr].lines[i].LRU < min_LRU)//update smallest LR
    U
174.            {
175.                target_line=i;
176.                min_LRU=cache->sets[setaddr].lines[i].LRU;
177.            }
178.        }
179.        return target_line;
180. }
181.
```

```c
182. //check whether the cache is full and need replacement
183. int cache_full(Scache *cache, int setaddr, int tagaddr)
184. {
185.     int full_flag=1;
186.     int i;
187.     //check full or not
188.     for(i=0;i<cache->line_cnt;++i)
189.     {
190.         if(cache->sets[setaddr].lines[i].valid==0)//there is an empty space

191.         {
192.             full_flag=0;
193.             break;
194.         }
195.     }
196.
197.     if(full_flag==0)//not full
198.     {
199.         cache->sets[setaddr].lines[i].valid=1;//update valid bits
200.         cache->sets[setaddr].lines[i].tag=tagaddr;//update tag
201.         update_LRU(cache, setaddr,i);
202.
203.     }
204.     else//cache is full
205.     {
206.         int eviction_line;
207.         eviction_line=find_smallest_LRU(cache,setaddr);//find the least rec
    ent used data
208.         cache->sets[setaddr].lines[eviction_line].valid=1;//update valid bi
    ts
209.         cache->sets[setaddr].lines[eviction_line].tag=tagaddr;//update tag

210.         update_LRU(cache,setaddr,eviction_line);
211.     }
212.     return full_flag;
213. }
214.
215. //load operation
216. void load_op(Scache *cache,int addr,int sizee,int setaddr,int tagaddr,int V
    erbose_mode)
217. {
218.     if(miss_checking(cache,setaddr,tagaddr)==1)//there is a miss
219.     {
220.         miss_count++;
```

```c
221.          if(Verbose_mode==1)
222.              printf("miss ");
223.          if(cache_full(cache,setaddr,tagaddr)==1)//the cache is full and nee
    d eviction
224.          {
225.              eviction_count++;
226.              if(Verbose_mode==1)
227.                  printf("eviction ");
228.          }
229.      }
230.      else//hit
231.      {
232.          hit_count++;
233.              if(Verbose_mode==1) printf("hit ");
234.      }
235. }
236.
237. //store operation
238. void store_op(Scache *cache,int addr,int sizee,int setaddr,int tagaddr,int
    Verbose_mode)
239. {
240.      load_op(cache,addr,sizee,setaddr,tagaddr,Verbose_mode);
241. }
242.
243. //modify operation
244. void modify_op(Scache *cache,int addr,int sizee,int setaddr,int tagaddr,int
    Verbose_mode)
245. {
246.      load_op(cache,addr,sizee,setaddr,tagaddr,Verbose_mode);
247.      store_op(cache,addr,sizee,setaddr,tagaddr,Verbose_mode);
248. }
249.
250. int main(int argc,char **argv)
251. {
252.      int s,E,b=0;
253.      int Verbose_mode=0;
254.      miss_count=0;
255.      hit_count=0;
256.      eviction_count=0;
257.      int addr,sizee;
258.
259.      Scache cache;
260.
261.      char cur_filename[100];
```

```
262.      char opt[10];
263.
264.      get_Opt(argc,argv,&s,&E,&b,cur_filename,&Verbose_mode);//get variables
    from command
265.      init_cache(s,E,b,&cache);//initial cache sim
266.      FILE *trace=fopen(cur_filename,"r");//open file and execute
267.
268.      while(fscanf(trace,"%s %x,%d",opt,&addr,&sizee)!=EOF)
269.      {
270.          if(strcmp(opt,"I")==0) continue;//jump I line
271.          int setaddr = get_set(addr,s,b);
272.              int tagaddr = get_tag(addr,s,b);
273.          if(Verbose_mode==1)
274.              printf("%s %x,%d ",opt,addr,sizee);
275.          if(strcmp(opt,"S")==0)//store operation
276.              store_op(&cache,addr,sizee,setaddr,tagaddr,Verbose_mode);
277.          if(strcmp(opt,"L")==0)//load operation
278.              load_op(&cache,addr,sizee,setaddr,tagaddr,Verbose_mode);
279.          if(strcmp(opt,"M")==0)//modify operation
280.              modify_op(&cache,addr,sizee,setaddr,tagaddr,Verbose_mode);
281.          if(Verbose_mode==1)
282.              printf("\n");
283.      }
284.      printSummary(hit_count, miss_count, eviction_count);
285.      return 0;
286. }
```

### 2.1.3 Evaluation

I use the *test-csim* program to automatically check the correctness of Part A. From the figure below, it could be concluded that I successfully solve Part A.

```
sjtulyf123@ubuntu:~/Architecture_Lab2/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.trac
e
     3 (2,1,3)     167      71      67     167      71      67  traces/trans.tra
ce
     3 (2,2,3)     201      37      29     201      37      29  traces/trans.tra
ce
     3 (2,4,3)     212      26      10     212      26      10  traces/trans.tra
ce
     3 (5,1,5)     231       7       0     231       7       0  traces/trans.tra
ce
     6 (5,1,5)  265189   21775   21743  265189   21775   21743  traces/long.trac
e
    27

TEST_CSIM_RESULTS=27
sjtulyf123@ubuntu:~/Architecture_Lab2/project2-handout$
```

## 2.2 Part B

### 2.2.1 Analysis

In Part B I modify the function ***transpose_submit*** in ***trans.c*** to implement transpose function of matrix and decrease the number of misses.

Let me briefly analyze the hits and misses when transposing the matrix. Since the size of block is 32 byte and one int element is 4 byte, so one block contains 8 elements in a matrix. We know that one cache contains 32 blocks, so it contains 256 elements in a matrix. If we simply visit each row and transpose these elements, it would cause a lot of conflict misses because the visit of matrix B would cause the eviction of some block in matrix A. However, this block needs to be visited in later operation. So we need to add this block to the cache again, which causes many misses. To reduce this kind of misses, I split the original 32*32 matrix into sixteen 8*8 small matrices and transpose each small matrix at a time, other than transpose the whole 32 elements in a row at a time. This could reduce the misses obviously, but still larger than 300.

After that, I notice that when it transposes the elements on the diagonal, it would also cause misses because the elements in A and B are mapped to the same block. Therefore, I store the element on the diagonal in block ***b*** to a variable and do transpose of this diagonal element after the program executes the other 7 elements in block ***b***. In this way, the number of misses is fewer than 300.

In 32*32 matrix, I **separate it into 8*8 small matrices**, **execute one small matrix at a time** and **optimize the operation on diagonal elements**; In 61*67 matrix, I **separate it into 16*16 small matrices**, **execute one small matrix at a time** and **optimize the operation on**
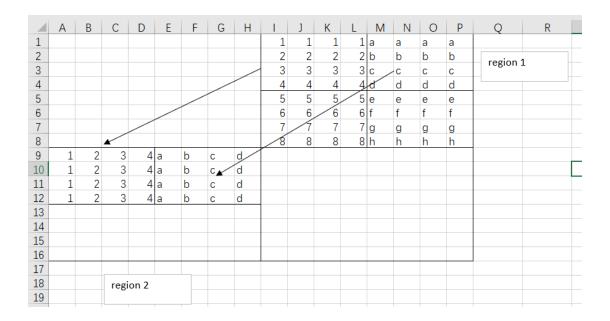
**diagonal elements**. The number of misses of these two matrix sizes is fewer than the requirements. However, the number of misses in 64*64 matrix is still large.

In 64*64 matrix, 4 rows of 64 elements would occupy the whole cache. Therefore, as shown in the figure below, if we transpose the elements in the second small 8*8 matrix and simply visit A by row and visit B by column, it would cause many misses: After visiting 9A, 10A, 11A 12A in matrix B, it wants to visit 13A, 14A, 15A, 16A in matrix B. However, the first 4 elements and last 4 elements in the same column are mapped onto the same blocks. Therefore, when we visit blocks in row 9,10,11,12, the blocks in row 13,14,15,16 would be replaced; when we visit blocks in row 13,14,15,16, the blocks in row 9,10,11,12 would also be replaced. So visiting every element in matrix B would cause a miss.
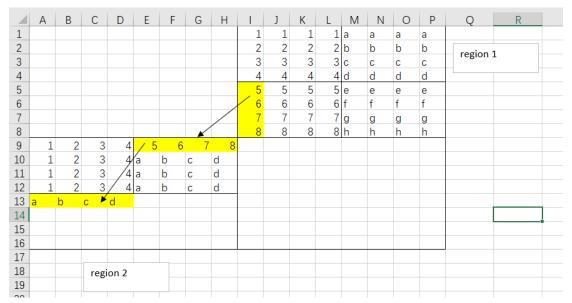
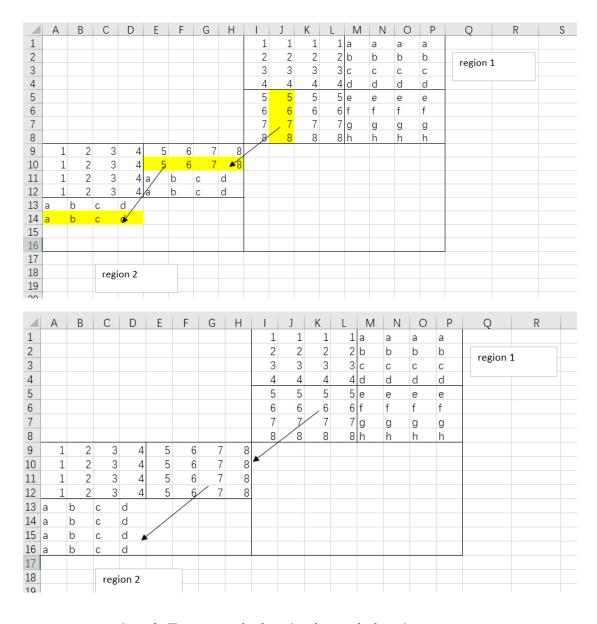| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 2 | | | | | | | | | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | |
| 3 | | | | | | | | | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | |
| 4 | | | | | | | | | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | |
| 5 | | | | | | | | | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | |
| 6 | | | | | | | | | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | |
| 7 | | | | | | | | | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | |
| 8 | | | | | | | | | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | |
| 9 | 1 | 2 | | | | | | | | | | | | | | | |
| 10 | 1 | 2 | | | | | | | | | | | | | | | |
| 11 | 1 | 2 | | | | | | | | | | | | | | | |
| 12 | 1 | 2 | | | | | | | | | | | | | | | |
| 13 | 1 | 2 | | | | | | | | | | | | | | | |
| 14 | 1 | 2 | | | | | | | | | | | | | | | |
| 15 | 1 | 2 | | | | | | | | | | | | | | | |
| 16 | 1 | 2 | | | | | | | | | | | | | | | |
| 17 | | | | | | | | | | | | | | | | | |

I optimize the operation in the following way (region 1 is in matrix A and region 2 is in matrix B). The key idea is to reduce the repeatedly visit of first 4 rows and last 4 rows in matrix B:

Step 1: Store first 4 rows in region 1 into region 2. During this process, the first 4 columns in first 4 rows are transposed and put into the correct places in region 2, the last 4 columns in first 4 rows are transposed but put into other empty places in region 2 in order not to visit the same block later.

Step 2: Put the last 4 columns in first 4 rows into correct places in region 2 by each row. Transpose the first 4 columns in last 4 rows and put them into the correct places in region 2. After this step, 3 small 4*4 matrices is transposed correctly.

Step 3: Transpose the last 4 columns in last 4 rows.

After executing the processes, 1 block in region 2 (matrix B) would just cause 1 miss. This kind of operation decreases the number of misses dramatically.

## 2.2.2 Code

```
1.  //Name: Yifan Liu
2.  //ID:   518021910609
3.
4.  /*
5.   * trans.c - Matrix transpose B = A^T
6.   *
7.   * Each transpose function must have a prototype of the form:
8.   * void trans(int M, int N, int A[N][M], int B[M][N]);
9.   *
10.  * A transpose function is evaluated by counting the number of misses
11.  * on a 1KB direct mapped cache with a block size of 32 bytes.
12.  */
13. #include <stdio.h>
14. #include "cachelab.h"
15.
16. int is_transpose(int M, int N, int A[N][M], int B[M][N]);
17.
18. /*
19.  * transpose_submit - This is the solution transpose function that you
20.  *     will be graded on for Part B of the assignment. Do not change
21.  *     the description string "Transpose submission", as the driver
22.  *     searches for that string to identify the transpose function to
23.  *     be graded.
24.  */
25. char transpose_submit_desc[] = "Transpose submission";
26. void transpose_submit(int M, int N, int A[N][M], int B[M][N])
27. {
28.     //12 variables in total
29.     //row is the row block
30.     //clo is the column block
31.     int i,j,row,col;
32.     int x1,x2,x3,x4,x5,x6,x7,x8;
33.
34.     if(M==32)//seperate 32*32 into 8*8, so the miss is smaller
35.     {
36.         for(row=0;row<N;row=row+8)
37.         for(col=0;col<M;col=col+8)
38.         for(i=row;i<row+8;++i)
```

```
39.          {
40.            for(j=col;j<col+8;++j)
41.             {
42.                if(i!=j) B[j][i]=A[i][j];
43.                else
44.                {
45.                    //i==j,the diagonal element,it need to save to B later,other
     wise the miss will increase
46.                    x1=A[i][j];
47.                    x2=i;
48.                }
49.             }
50.            if(col==row)
51.             {
52.                B[x2][x2]=x1;//save element to B in this time. This would no
     t cause more miss
53.             }
54.          }
55.        }
56.      else if(M==64)
57.      {
58.            for(row=0;row<N;row=row+8)
59.            for(col=0;col<M;col=col+8)
60.            {
61.                for(i=row;i<row+4;i++)
62.                {
63.                    x1=A[i][col];x2=A[i][col+1];x3=A[i][col+2];x4=A[i][col+3];
64.                    x5=A[i][col+4];x6=A[i][col+5];x7=A[i][col+6];x8=A[i][col+7];
     //first line
65.                    //first 4 element in first 4 lines:do transpose
66.                    B[col][i]=x1;B[col+1][i]=x2;B[col+2][i]=x3;B[col+3][i]=x4;
67.                    //last 4 elements in first 4 lines:transpose and move these
     to ''some empty B space''
68.                    B[col][i+4]=x5;B[col+1][i+4]=x6;B[col+2][i+4]=x7;B[col+3][i+
     4]=x8;
69.                }
70.                for(j=col;j<col+4;++j)
71.                {
72.
73.                    x1=A[row+4][j];x2=A[row+5][j];x3=A[row+6][j];x4=A[row+7][j];

74.                    x5=B[j][row+4];x6=B[j][row+5];x7=B[j][row+6];x8=B[j][row+7];

75.                    //last 4 elements in first 4 columns: do transpose
```

```
76.              B[j][row+4]=x1;B[j][row+5]=x2;B[j][row+6]=x3;B[j][row+7]=x4;

77.              //transpose data in ''some empty B space'' to correct space

78.              B[j+4][row]=x5;B[j+4][row+1]=x6;B[j+4][row+2]=x7;B[j+4][row+3]=x8;
79.          }
80.        for(i=row+4;i<row+8;++i)
81.        {
82.              //other data:do transpose per line
83.              x1=A[i][col+4];x2=A[i][col+5];x3=A[i][col+6];x4=A[i][col+7];

84.              B[col+4][i]=x1;B[col+5][i]=x2;B[col+6][i]=x3;B[col+7][i]=x4;

85.          }
86.        }
87.    }
88.    else//67*61,separate it into 16*16
89.    {
90.        for(row=0;row<N;row=row+16)
91.        for(col=0;col<M;col=col+16)
92.        for(i=row;i<row+16 &&(i<N);++i)
93.        {
94.        for(j=col;j<col+16&&(j<M);++j)
95.         {
96.            if(i!=j) B[j][i]=A[i][j];
97.            else
98.            {
99.                //i==j,the diagonal element,it need to save to B later,other
   wise the miss will increase
100.               x1=A[i][j];
101.               x2=i;
102.            }
103.          }
104.        if(col==row)
105.          {
106.              B[x2][x2]=x1;//save element to B in this time. This would n
   ot cause more miss
107.          }
108.        }
109.    }
110. }
111.
112. /*
```

```c
113.   * You can define additional transpose functions below. We've defined
114.   * a simple one below to help you get started.
115.   */
116.
117. /*
118.   * trans - A simple baseline transpose function, not optimized for the cach
       e.
119.   */
120. char trans_desc[] = "Simple row-wise scan transpose";
121. void trans(int M, int N, int A[N][M], int B[M][N])
122. {
123.       int i, j, tmp;
124.
125.       for (i = 0; i < N; i++) {
126.           for (j = 0; j < M; j++) {
127.               tmp = A[i][j];
128.               B[j][i] = tmp;
129.           }
130.       }
131.
132. }
133.
134. /*
135.   * registerFunctions - This function registers your transpose
136.   *     functions with the driver.  At runtime, the driver will
137.   *     evaluate each of the registered functions and summarize their
138.   *     performance. This is a handy way to experiment with different
139.   *     transpose strategies.
140.   */
141. void registerFunctions()
142. {
143.       /* Register your solution function */
144.       registerTransFunction(transpose_submit, transpose_submit_desc);
145.
146.       /* Register any additional transpose functions */
147.       registerTransFunction(trans, trans_desc);
148.
149. }
150.
151. /*
152.   * is_transpose - This helper function checks if B is the transpose of
153.   *     A. You can check the correctness of your transpose by calling
154.   *     it before returning from the transpose function.
155.   */
```

```
156. int is_transpose(int M, int N, int A[N][M], int B[M][N])
157. {
158.     int i, j;
159.
160.     for (i = 0; i < N; i++) {
161.         for (j = 0; j < M; ++j) {
162.             if (A[i][j] != B[j][i]) {
163.                 return 0;
164.             }
165.         }
166.     }
167.     return 1;
168. }
```

### 2.2.3 Evaluation

I use driver.py to check both Part A and Part B. From the figure below, it could be concluded that I successfully decrease the number of misses to the goal in 32*32, 64*64 and 61*67 matrices. I solve Part B successfully.

```
sjtulyf123@ubuntu:~/Architecture_Lab2/project2-handout$ ./driver.py
Part A: Testing cache simulator
Running ./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


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       1179
Trans perf 61x67        10.0        10       1985
      Total points      53.0        53
sjtulyf123@ubuntu:~/Architecture_Lab2/project2-handout$
```

# 3. Conclusion

## 3.1 Problems

In part A, I meet a problem of how to correctly get the variables s, E, and b from the command line. After searching for the use of **getopt** function and **optarg** variable, I solve this problem.

When writing the store operation, I am confused with what I need to do in this operation. To solve this, I refer to the CSAPP book to check what the cache does when it needs to store some data. It turns out that the store operation is just like the load operation, especially when counting the

number of hits, misses and evictions. Similarly, the modify operation is just one load and one store operation.

I also overcome a problem when I try to test my Part A program by using *valgrind*. The command says: "valgrind: not found". I thought that there are some problems with my program or segmentation fault, but I could not figure out what happens. After carefully check the hints in the project guide, I notice that I forget to install *valgrind* on my virtual machine. After installing it, the problem is solved.

In part B, it is easy to discover that separating matrix into some small matrices is useful, but determining the size of small matrix takes me a lot of time. I tried 4*4, 8*8 and 16*16 and compare the number of misses. Finally, I decide to use 8*8 size in 32*32 and 64*64 matrix, use 16*16 size in 61*67 matrix. But the number of misses is still larger that the requirement. After analyzing the behavior of transpose operation, I discover that the diagonal elements would also cause misses.

When optimizing 64*64 transpose matrix, the former method is not enough. After searching for some matrix optimizing methods, I finally realize that storing every element in B would cause a miss, so I continue to optimize the operation. The difficulty of Part B is to find a valid method to reduce the number of misses.

## 3.2   Achievements

1. I successfully implement the simulated cache and it could correctly report the number of misses, hits and evictions. I also successfully implement the matrix transpose function to reduce the number of misses when transposing a matrix.

2. Both of the programs have good robustness. They could detect the invalid input and send the corresponding error messages.

3. I have a better understanding of how cache works and how the misses happen when transposing the matrix.