# Assignment 1
# 2015313754 TaehyungGil (길태형)

My code consists of three parts.

**1. Get the whole input data from 'hw1_input.txt',**
      **and map key attribute string to integer values between 0 and 63.**

**2. Sort only key attribute values using optimized radix sort.**
      **(Optimized by adding two mapping lists)**

**3. Write 'hw1_output.txt' in order according to map_list.**

Below are detailed explanations and performance analysis of each part.


**1. Get the whole input data from 'hw1_input.txt',**
      **and map key attribute string to integer values between 0 and 63.**

```
53        input_state = fgets(input_buffer,BUFFER_SIZE,fd_input);
54        //n is stored in input buffer
55        while(input_buffer[idx]!='\n')
56        {
57            data_size*=10;
58            data_size += input_buffer[idx]-'0';
59            idx++;
60        }
61        // gets $
62        input_state = fgets(input_buffer,BUFFER_SIZE,fd_input);
63        // gets list of attributes
64        input_state = fgets(input_buffer,BUFFER_SIZE,fd_input);
65        idx=0;
66        while(input_buffer[idx]!='\n')
67        {
68            if(input_buffer[idx]==DELIMITER)
69                attribute_num++;
70            else if(input_buffer[idx]==OPEN_BR)
71                key_attribute_idx=attribute_num-1;
72            idx++;
73        }
74        // gets $
75        input_state = fgets(input_buffer,BUFFER_SIZE,fd_input);
```

It gets the value of n and t, and gets what is the key attribute.

```
89        for(int data_idx=0;data_idx<data_size;data_idx++)
90        {
91            input_state = fgets(input_buffer,BUFFER_SIZE,fd_input);
92            int deli_cnt=0;
93            int buff_idx=0;
94            int data_m_idx=0;
95            while(deli_cnt!=key_attribute_idx)
96            {
97                if(input_buffer[buff_idx]==DELIMITER)
98                {
99                    deli_cnt++;
100                   buff_idx++;
101               }
102               else
103                   buff_idx++;
104           }
105           // printf("deli_cnt : %d key_attribute_idx : %d \n",deli_cnt,key_attribute_idx);
106           while(input_buffer[buff_idx]!=DELIMITER && input_buffer[buff_idx]!=NEW_LINE && input_buffer[buff_idx]!=EOF && input_buffer[buff_idx]!=0 )
107           {
108               unsigned char data=0;
109               char b_data = input_buffer[buff_idx];
110               if(b_data<NUM_0) //space
111                   data_block[data_idx][data_m_idx]=1;
112               else if(b_data<AL_A) //numerical
113                   data_block[data_idx][data_m_idx]=b_data-NUM_0+2;
114               else if(b_data<AL_a) //Upper letters
115                   data_block[data_idx][data_m_idx]=2*(b_data-AL_A)+13;
116               else                 //Lower letters
117                   data_block[data_idx][data_m_idx]=2*(b_data-AL_a)+12;
118               // printf("%c -> %d \n",b_data,data_block[data_idx][data_m_idx]);
119               buff_idx++;
120               data_m_idx++;
121           }
122           buff_idx=-1;
123           do
124           {
125               buff_idx++;
126               // printf("buff_idx : %d\n",buff_idx);
127               org_data[data_idx][buff_idx]=input_buffer[buff_idx];
128           } while(input_buffer[buff_idx]!=NEW_LINE&&input_buffer[buff_idx]!=0);
129           org_data_len[data_idx]=buff_idx+1;
130       }
```

And gets each input line. Here, it maps each key attribute character into a integer value as below.

Blank → 0 (Blank is not space. It means there is no character.)
Space → 1
'0'     → 2
'1'     → 3
.
.
.
'a'     → 12
'A'     → 13
.
.
.
'z'     → 62
'Z'     → 63

It stores these mapped integer values into 'data_block' and 'temp_block'.
And it stores the whole attributes to 'org_data'.

```
20    static char org_data[MAX_BLOCK_N][(MAX_M)*(MAX_T)];
21    static int org_data_len[MAX_BLOCK_N];
22    static unsigned char data_block[MAX_BLOCK_N][MAX_M];
23    static unsigned char temp_block[MAX_BLOCK_N][MAX_M];
```

I declared global static arrays, 'data_block', 'temp_block', and 'org_data'.

This code block takes $O(nmt)$ , where n is the number of objects, m is the maximum length of each attribute value, and t is the number of attributes for each object. Since each maximum values of n,m, and t are 1,000,000, 15, 10, the execution count is as similar as 150,000,000. (not sec, just representation of counts the code should run.) It does $O(nmt)$ twice because it first maps into integer
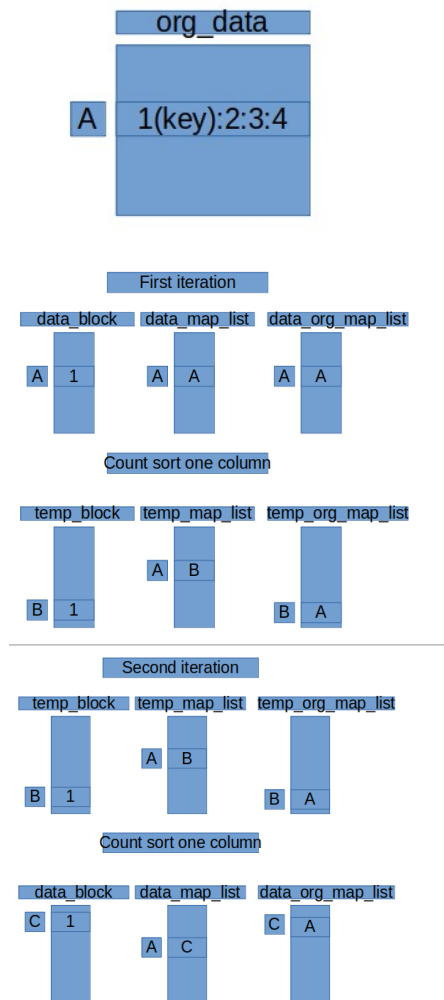
values and then store the whole attributes. I implemented like this for some optimization for storing whole attributes.

Anyway, **this code blocks costs O(nmt)**

## 2. Sort only key attribute values using optimized radix sort.
   **(Optimized by adding two mapping lists)**

```
152    void cnt_sort(int data_size, int col_idx) //count sort the directed index
153    {
154        for(int i=0;i<MAX_R;i++)
155            cnt_list[i]=0;
156        if(now_block_num==0)
157        {
158            for(int data_idx=0;data_idx<data_size;data_idx++)
159            {
160                unsigned char d_data = data_block[data_idx][col_idx];
161                cnt_list[d_data]++;
162            }
163            for(int i=1;i<MAX_R;i++)        //getting prefix sum
164                cnt_list[i]+=cnt_list[i-1];
165            for(int data_idx=data_size-1;data_idx>=0;data_idx--)
166            {
167                unsigned char d_data = data_block[data_idx][col_idx];
168                int next_idx = cnt_list[d_data]-1;
169                for(int m_idx=0;m_idx<MAX_M;m_idx++)
170                {
171                    temp_block[next_idx][m_idx]=data_block[data_idx][m_idx];
172                }
173                int org_idx = data_org_map_list[data_idx];
174                temp_map_list[org_idx]=next_idx;
175                temp_org_map_list[next_idx]=org_idx;
176                cnt_list[d_data]-=1;
177            }
178            now_block_num=1;
179        }
180        else
181        {
182            for(int data_idx=0;data_idx<data_size;data_idx++)
183            {
184                unsigned char d_data = temp_block[data_idx][col_idx];
185                cnt_list[d_data]++;
186            }
187            for(int i=1;i<MAX_R;i++)        //getting prefix sum
188                cnt_list[i]+=cnt_list[i-1];
189
190            for(int data_idx=data_size-1;data_idx>=0;data_idx--)
191            {
192                unsigned char d_data = temp_block[data_idx][col_idx];
193                int next_idx = cnt_list[d_data]-1;
194                for(int m_idx=0;m_idx<MAX_M;m_idx++)
195                {
196                    data_block[next_idx][m_idx]=temp_block[data_idx][m_idx];
197                }
198                int org_idx = temp_org_map_list[data_idx];
199                data_map_list[org_idx]=next_idx;
200                data_org_map_list[next_idx]=org_idx;
201                cnt_list[d_data]-=1;
202            }
203            now_block_num=0;
204        }
205
206        return;
207    }
208
```

I implemented count sorting function. The column index which should be sorted is given to the function. As I said above, there is only 64 integer values. So, it takes O(n+64) to count sort one column. As count sorting doesn't sort in place, I implemented one more data block called 'temp_block'. The code uses both blocks alternatively. And two mapping lists are given to each block. One is 'map_list' and another one is 'org_map_list'. 'map_list' records where the original data is in block. And, 'org_map_list' records which original data is stored in block. It is illustrated below.



It iterates this for 15 times. So after 15 iterations, the sorted information is stored in 'temp_block', 'temp_map_list', and 'temp_org_map_list'.

This count sorting has some overhead. This one needs to manage two map_list, and it takes 2*n for one iteration. So, the time complexity for count sorting once is O(n+64) + O(2*n) = O(n).

And this radix sort does count sorting for 15 times. So, it takes O(15*n) = O(n).

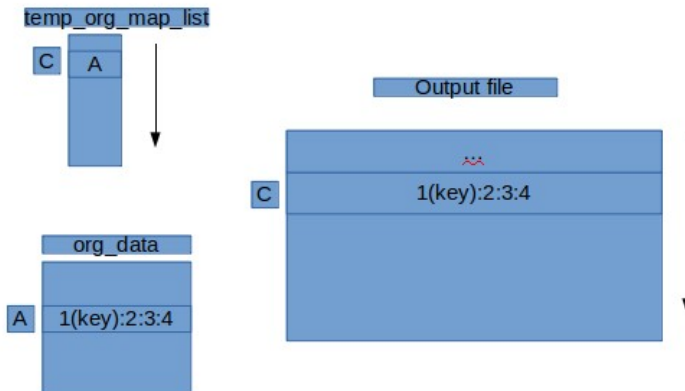**In conclusion, this sorting is done in O(n).**

**3. Write 'hw1_output.txt' in order according to map_list.**

```
143        for(int i=0;i<data_size-1;i++)
144        {
145            org_idx = temp_org_map_list[i];
146            fwrite(org_data[org_idx],1,org_data_len[org_idx]-1,fd_output);
147            fputc(NEW_LINE,fd_output);
148        }
149        org_idx = temp_org_map_list[data_size-1];
150        fwrite(org_data[org_idx],1,org_data_len[org_idx]-1,fd_output);
```



In 'temp_org_map_list', the original object index is stored. So, just writing the proper original object in order is enough to make 'hw1_output.txt'. It can decide to write which object in the following place. It just goes through the 'temp_org_map_list' from 0~(n-1) index, writing the corresponding original object to output file. **This code block takes O(nmt).**

**In summary,**
**first part takes O(nmt), second part takes O(n), last part takes O(nmt).**

**So, the total complexity is O(nmt) + O(n) + O(nmt) = O(nmt)**

**Thank you for reading.**