# ECE785 – Project 3A: Code Optimization

Name: Zhiyu Chen

UnityID: zchen45

1. Overview

In this project, I am supposed to profile the execution time for the battery modeling program from Dr. Chow's research group and then optimize it within a limited amount of time. Obtain and analyze the execution time profile of the CoEst_offline.c program.
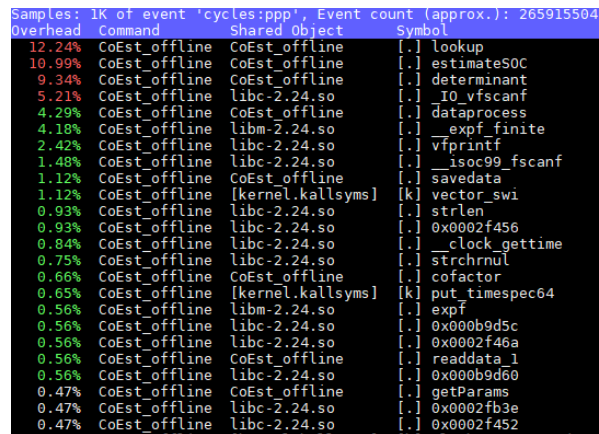
2. Optimization

2.1. Initial profile

In this project, I use the clock_gettime to measure the time duration. And the duration to be measured is between "fopen" and "fclose". By the command "gcc CoEst_offline.c –o CoEst_offline –lm; ./CoEst_offline", the results are listed below.

Time: 210890.500 us

The Figure 1 is the perf report. As we can see, about 33 percent of overheads occurs in the function lookup, estimateSOC and determinant.



```
Samples: 1K of event 'cycles:ppp', Event count (approx.): 265915504
Overhead  Command       Shared Object        Symbol
  12.24%  CoEst_offline  CoEst_offline       [.] lookup
  10.99%  CoEst_offline  CoEst_offline       [.] estimateSOC
   9.34%  CoEst_offline  CoEst_offline       [.] determinant
   5.21%  CoEst_offline  libc-2.24.so        [.] _IO_vfscanf
   4.29%  CoEst_offline  CoEst_offline       [.] dataprocess
   4.18%  CoEst_offline  libm-2.24.so        [.] __expf_finite
   2.42%  CoEst_offline  libc-2.24.so        [.] vfprintf
   1.48%  CoEst_offline  libc-2.24.so        [.] __isoc99_fscanf
   1.12%  CoEst_offline  CoEst_offline       [.] savedata
   1.12%  CoEst_offline  [kernel.kallsyms]   [k] vector_swi
   0.93%  CoEst_offline  libc-2.24.so        [.] strlen
   0.93%  CoEst_offline  libc-2.24.so        [.] 0x0002f456
   0.84%  CoEst_offline  libc-2.24.so        [.] __clock_gettime
   0.75%  CoEst_offline  libc-2.24.so        [.] strchrnul
   0.66%  CoEst_offline  CoEst_offline       [.] cofactor
   0.65%  CoEst_offline  [kernel.kallsyms]   [k] put_timespec64
   0.56%  CoEst_offline  libm-2.24.so        [.] expf
   0.56%  CoEst_offline  libc-2.24.so        [.] 0x000b9d5c
   0.56%  CoEst_offline  libc-2.24.so        [.] 0x0002f46a
   0.56%  CoEst_offline  CoEst_offline       [.] readdata_1
   0.56%  CoEst_offline  libc-2.24.so        [.] 0x000b9d60
   0.47%  CoEst_offline  CoEst_offline       [.] getParams
   0.47%  CoEst_offline  libc-2.24.so        [.] 0x0002fb3e
   0.47%  CoEst_offline  libc-2.24.so        [.] 0x0002f452
```

Figure 1    Perf Report 1

2.2. Compiler flags

The first step is to set the compiler flags. This is a simple but effective way to improve the performance. Also it is necessary for vectorizing the code in the future. The follows show the flags implemented.

Flags:

-mfpu=neon   -mfloat-abi=hard   -mcpu=cortex-a8   -mtune=cortex-a8   -ftree-vectorize   -mvectorize-with-neon-quad -O3 -g    -pg -ffast-math

Time:        149865.708 us

Perf report:



Figure 2      Perf Report 2

2.3. Lookup function

As we can see in Figure 1, lookup function generates the most overheads. Also by annotating the main function in the perf report 2, the lookup function produces more than 50 percent of overheads as Figure 3 shown below. So the next step is to optimize the lookup function.



Figure 3      Annotate main in perf report 2

In the annotation, it is obvious that the while loop generates lots of overhead.

```
1.  int i = 0; // index of look-up table
2.  float slope = 0.0;
3.  // find the index
4.  while (estVoc > table_Voc[i] && i < 100){
5.      i++;
```

```
6.  }
```

Every time when lookup function is called, it will compare estVoc with the table of voc, from low index to high index, until the estVoc is bigger than the voc or it compares 100 times. According to the figure 3, there is a red line leading to 55 percent of overheads. It is the conditional branch of the while loop. The optimized code is shown below.

```
1.   float32x4_t n_index, n_index_2, n_index_3, t_voc;
2.   uint32x4_t mask_i = vmovq_n_u32(0);
3.   uint32x4_t change_or_not;
4.   float32x4_t index_add = {1.0f, 2.0f, 3.0f, 4.0f};
5.   float32x4_t estVoc_n = vld1q_dup_f32(&estVoc);
6.
7.   for(int a = 0; a < 100; a = a + 4){
8.     index = (float) a;
9.     n_index = vdupq_n_f32(index);
10.    n_index = vaddq_f32(n_index,index_add);
11.
12.    t_voc = vld1q_f32(&table_Voc[a]);
13.    mask_i = vcgtq_f32(estVoc_n, t_voc);
14.    n_index_3 = n_index_2;
15.    n_index_2 = vbslq_f32(mask_i, n_index, n_index_2);
16.    n_index_2 = vmaxq_f32(n_index_2, n_index_3);
17. }
```

As the conditional branch produces too much overheads, I decide to vectorize the code to unroll the while loop. At here, I compare the estVoc with all the elements in voc table and record the element with the smallest index.

Time:      134308.416 us

Perf report:



Figure 4    perf report 3

## 2.4. Matrix calculation optimization

In the original code, it implements a lot of matrix multiplication. The follow is an example.

```
1.  // Step 7: Calculate Ad^2 =  A * A
2.  for (int i = 0; i < tforder; i++){
3.      for (int j = 0; j < tforder; j++){
4.          AA[i][j] = 0;
5.          for (int k = 0; k < tforder; k++)
6.              AA[i][j] += Ad[i][k] * Ad[k][j];
7.      }
8.  }
```

It implements 3 for-loops to achieve the multiplication of 2 2x2 matrixes. Actually how we calculate it on paper?

$$\begin{bmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{bmatrix} \times \begin{bmatrix} b_{00} & b_{01} \\ b_{10} & b_{11} \end{bmatrix} = \begin{bmatrix} a00b00 + a01b10 & a00b01 + a01b11 \\ a10b00 + a11b10 & a10b01 + a11b11 \end{bmatrix}$$
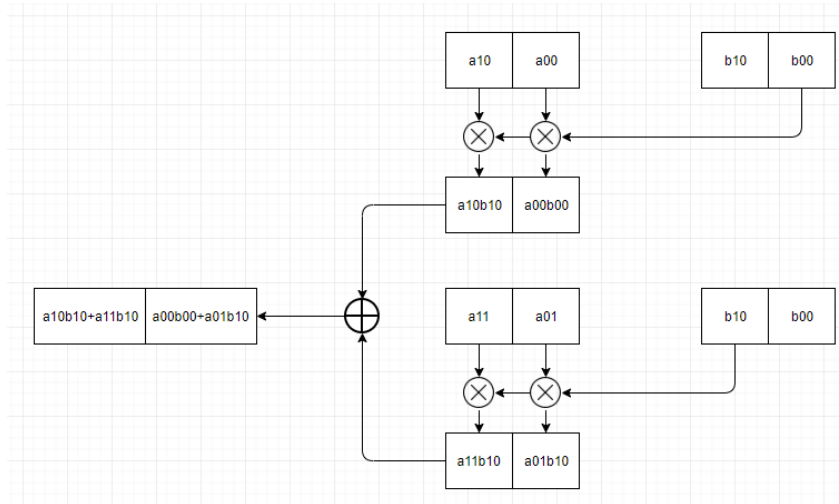


Figure 5    Matrix Multiplication

If each column of matrixes is a vector in NEON register, we can use the vector-by-scalar-multiplication to calculate each result column efficiently. The optimized code is shown below.

```
1.  // Step 7: Calculate Ad^2 =  A * A
2.  float temp[4] = {Ad[0][0], Ad[1][0], Ad[0][1], Ad[1][1]};
3.  float32x2_t Ad_n0 = vld1_f32(&temp[0]);
4.  float32x2_t Ad_n1 = vld1_f32(&temp[2]);
5.
6.  float32x2_t AA_n0 = vadd_f32(vmul_n_f32(Ad_n0, Ad_n0[0]), vmul_n_f32(Ad_n1,
    Ad_n0[1]));
7.  float32x2_t AA_n1 = vadd_f32(vmul_n_f32(Ad_n0, Ad_n1[0]), vmul_n_f32(Ad_n1,
    Ad_n1[1]));
```

Time:        136094.291 us

Perf report:



Figure 6    Perf Report 4

But the result doesn't meet my expectation. The time of duration increases. Maybe it is influenced by the size of the for-loop. As it is only a 2x2 matrix, the benefit from vectorizing is small. Also, it is said that the NEON is in-order executing and has limited hit-under-miss leading to limited memory latency.