# lab_B　Bloom filter　　　　108061173　邱崇喆

## Lab overview:

In this lab, we want to score all documents with a given user array and word weights, and every document will be inputted as an array of uint words, where the most significant 24 bits are word ID and the other 8 bits indicate its occurring frequency.

Only words (word ID) that occur in the given user array will count for the score of a document, namely their corresponding weights are not zero, and that is why we need a bloom filter to efficiently search if a word is inside the input array.

We roughly divide the implementation into three part:
- prepare input data, including documents array and bloom filter coefficient
- check the existence of all words in the documents by bloom filter
- evaluate every document

# Where do we need FPGA?

Let's see the software-only codes first:
////////////////////////////////////////////////////////////////////////////////////////////////////////////////////

MurmurHash2() can be used to calculate the hash of input doc words:

```
15 unsigned int MurmurHash2 ( const void * key, int len, unsigned int seed )
16 {
17   // 'm' and 'r' are mixing constants generated offline.
18   // They're not really 'magic', they just happen to work well.
19
20   const unsigned int m = 0x5bd1e995;
21 //  const int r = 24;
22
23   // Initialize the hash to a 'random' value
24
25   unsigned int h = seed ^ len;
26
27   // Mix 4 bytes at a time into the hash
28
29   const unsigned char * data = (const unsigned char *)key;
30
31
32   switch(len)
33   {
34   case 3: h ^= data[2] << 16;
35   case 2: h ^= data[1] << 8;
36   case 1: h ^= data[0];
37           h *= m;
38   };
39
40   // Do a few final mixes of the hash to ensure the last few
41   // bytes are well-incorporated.
42
43   h ^= h >> 13;
44   h *= m;
45   h ^= h >> 15;
46
47   return h;
48 }
49
```

Use the results of Murmur for bloom filtering, checking the existence of the words:

```
25      unsigned int size_offset=0;
26
27      chrono::high_resolution_clock::time_point t1 = chrono::high_resolution_clock::now();
28 |
29      unsigned char* inh_flags = (unsigned char*)aligned_alloc(4096, total_size*sizeof(char));
30
31
32      for(unsigned int doc=0;doc<total_num_docs;doc++)
33      {
34          profile_score[doc] = 0.0;
35          unsigned int size = doc_sizes[doc];
36
37          for (unsigned i = 0; i < size ; i++)
38          {
39              unsigned curr_entry = input_doc_words[size_offset+i];
40              unsigned word_id = curr_entry >> 8;
41              unsigned hash_pu =  MurmurHash2( &word_id , 3,1);
42              unsigned hash_lu =  MurmurHash2( &word_id , 3,5);
43              bool doc_end = (word_id==docTag);
44              unsigned hash1 = hash_pu&hash_bloom;
45              bool inh1 = (!doc_end) && (bloom_filter[ hash1 >> 5 ] & ( 1 << (hash1 & 0x1f)));
46              unsigned hash2 = (hash_pu+hash_lu)&hash_bloom;
47              bool inh2 = (!doc_end) && (bloom_filter[ hash2 >> 5 ] & ( 1 << (hash2 & 0x1f)));
48
49
50              if (inh1 && inh2) {
51                  inh_flags[size_offset+i]=1;
52              }else {
53                  inh_flags[size_offset+i]=0;
54              }
55          }
56
57          size_offset+=size;
58      }
59
```

If a word is existing in the bloom filter, we'll set the corresponding bool entry in the inh_flags array as 1, which will be used for computing score.
/////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////

a code snippet of setupData() sets up the bloom filter:

```
89    profile_weights.reserve( (1L << 24) );
90    for (unsigned i=0; i<(1L << bloom_size); i++) {
91        bloom_filter[i] = 0x0;
92    }
93    std::cout << "Creating profile weights" << endl;
94    std::cout << endl;
95
96    for (unsigned i=0; i<(1L << 24); i++) {
97        profile_weights[i] = 0;
98    }
99
00    for (unsigned i=0; i<16384; i++) {
01        unsigned entry = (rand()%(1<<24));
02
03        profile_weights[entry] = 10;
04        unsigned hash_pu = MurmurHash2(&entry,3,1);
05        unsigned hash_lu = MurmurHash2(&entry,3,5);
06
07        unsigned hash1 = hash_pu&hash_bloom;
08        unsigned hash2 = (hash_pu+hash_lu)&hash_bloom;
09
10        bloom_filter[ hash1 >> 5 ] |= 1 << (hash1 & 0x1f);
11        bloom_filter[ hash2 >> 5 ] |= 1 << (hash2 & 0x1f);
12    }
13
14 }
```

The following is a code snippet of compute_score:

```
61
62    for(unsigned int doc=0, n=0; doc<total_num_docs;doc++)
63    {
64        profile_score[doc] = 0.0;
65        unsigned int size = doc_sizes[doc];
66
67        for (unsigned i = 0; i < size ; i++,n++)
68        {
69            if(inh_flags[n])
70            {
71                unsigned curr_entry = input_doc_words[n];
72                unsigned frequency = curr_entry & 0x00ff;
73                unsigned word_id = curr_entry >> 8;
74                profile_score[doc]+= profile_weights[word_id] * (unsigned long)frequency;
75            }
76        }
77    }
```

Among the aforementioned codes, the best part we can choose to implement on FPGA is the two snippets in the red segment, because:
- MurmurHash2 consists of four XORs, three arithmetic shifts, and two multiplication operations.
- A shift of 1-bit in an arithmetic shift operation takes one clock cycle on the CPU, so we need at most 44 cycles.
- We can create custom architectures on FPGA, and therefore create an accelerator that will shift the data by an arbitrary number of bits in single clock cycle.
- FPGAs have DSP units that can perform faster multiplication than CPU.
- The second snippet in the red segment sequentially accesses the input_doc_array, which allows sequential accesses to DDR. It'll be efficient if FPGAs can do this.

The aforementioned compute_score snippet won't be implemented on FPGA because:
- The profile_weights array (weights) is not sequentially but randomly accessed.
  - The profile_weights array in this lab is 128MB, which means if implemented on FPGAs, it will be a big performance bottleneck.
- Another reason is that it only takes about 11% of running time, so we can keep this on the host side.

The following is the execution time of mere software:

```
03.bMXnEl@HLS03:~/workspace/LAB_B/Vitis-Tutorials/Hardware_Acceleration/Design_Tutorials/02-bloom/cpu_src$ make run
rm -rf temp_dir log_dir report_dir *log host runOnfpga* *.csv *summary .run .Xil vitis* *jou xilinx* gpofresult.txt gmon.out
g++ -D__USE_XOPEN2K8 -D__USE_XOPEN2K8 \
        -I. \
        -O3 -Wall -fmessage-length=0 -std=c++11\
        ./compute_score_host.cpp \
        ./MurmurHash2.c \
        ./main.cpp \
        -o ./host
./host 100000
Initializing data
Creating documents - total size : 1398.903 MBytes (349725824 words)
Creating profile weights

 Total execution time of CPU          |   1708.3256 ms
 Compute Hash processing time         |   1447.0368 ms
 Compute Score processing time        |    261.2888 ms
-----------------------------------------------------------------
 Execution COMPLETE
```

We can see a majority of processing time is spent on Hashing.

# Task Flow:

In this lab, we'll input 100000 documents equivalent to about 350M words (1.4GB) and bloom filter coefficients to the kernel to compute hashes and flags, and then transmit the flags back to the host, and finally score every document.
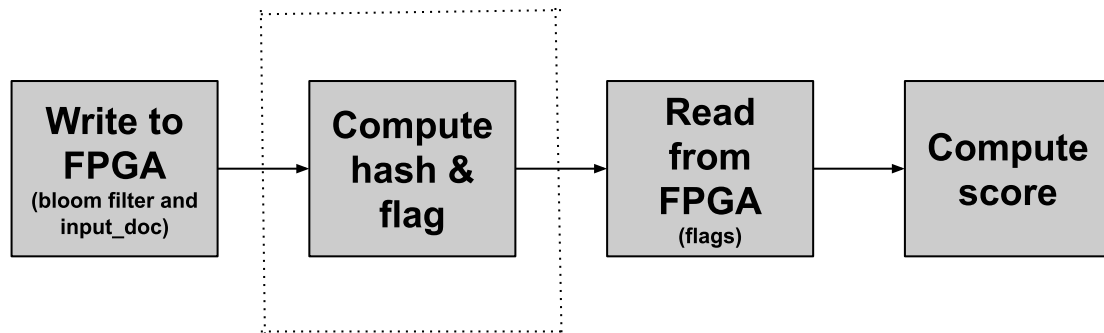
```
┌──────────────┐   ┌──────────────┐   ┌──────────────┐   ┌──────────────┐
│  Write to    │   │   Compute    │   │    Read      │   │              │
│    FPGA      │──▶│   hash &     │──▶│    from      │──▶│   Compute    │
│ (bloom filter│   │    flag      │   │    FPGA      │   │    score     │
│ and input_doc)│  │              │   │   (flags)    │   │              │
└──────────────┘   └──────────────┘   └──────────────┘   └──────────────┘
```

fig.1 The flow of the lab, the block surrounded by the dotted line is implemented on FPGA

It is obvious that if we implement the lab like fig 1, namely sequentially running the task blocks like running on CPU, we'll end up yielding bad performance, so we'll need to take advantage of concurrent processing and overlapping of the FPGA.

We'll start to discuss how to further optimize the above task flow in the following section.

# Optimizations of the Task Flow:

Process multiple word in the kernel (Kernel Parallelization) :

In this lab, we use the following interface requirements to create kernel:
- Read multiple words stored in the DDR as a 512-bit DDR access, equivalent of reading 16 words per DDR access.
- Write multiple flags to the DDR as a 512-bit DDR access, equivalent to writing 32 flags per DDR access.
- Compute 4 words in parallel with each word requiring two MurmurHash2 functions
- Compute the hash (two MurmurHash2 functions) functions for 4 words every cycle.

The kernel has 5 arguments:
1. 512-bit "output_flags" (through buffer)
2. 512-bit "input_words" (through buffer)
3. bloom_filter: Pointer of array with Bloom coefficients. (through buffer)
4. Total number of words to be computed
5. load_filter: Enable or disable loading bloom filter coefficients.

Following this requirements, we can expand the kernel (dotted lined block in fig.1) into fig.2



fig.2 The flow of the kernel

source:
https://github.com/Xilinx/Vitis-Tutorials/blob/2022.1/Hardware_Acceleration/Design_Tutorials/02-bloom/4_implement-kernel.md

The function compute_hash_flags_dataflow is responsible for this flow.

```
void compute_hash_flags_dataflow(
    ap_uint<512>*   output_flags,
    ap_uint<512>*   input_words,
    unsigned int    bloom_filter[PARALLELIZATION][bloom_filter_size],
    unsigned int    total_size)
{
#pragma HLS DATAFLOW

    hls::stream<ap_uint<512> >    data_from_gmem;
    hls::stream<parallel_words_t> word_stream;
    hls::stream<parallel_flags_t> flag_stream;
    hls::stream<ap_uint<512> >    data_to_gmem;
    . . . .
}
```

The kernel file includes:

1. Murmurhash2() — Hashing

```
1 #include <iostream>
2 #include <ctime>
3 #include <utility>
4 #include <cstdio>
5 #include <cstdlib>
6 #include <cstring>
7
8 #include "hls_stream_utils.h"
9 #include "sizes.h"
10
11 #define TOTAL_SIZE 3500000
12
13 //TRIPCOUNT identifiers
14 const unsigned int t_size = TOTAL_SIZE;
15 const unsigned int pf = PARALLELISATION;
16
17
18 #ifndef PARALLELISATION
19 #define PARALLELISATION 8
20 #endif
21
22 typedef ap_uint<sizeof(int )*8*PARALLELISATION> parallel_words_t;
23 typedef ap_uint<sizeof(char)*8*PARALLELISATION> parallel_flags_t;
24
25 const unsigned int bloom_filter_size = 1<<bloom_size;
26
27 unsigned int MurmurHash2(unsigned int key, int len, unsigned int seed)
28 {
29   const unsigned char* data = (const unsigned char *)&key;
30   const unsigned int m = 0x5bd1e995;
31   unsigned int h = seed ^ len;
32   switch(len) {
33     case 3: h ^= data[2] << 16;
34     case 2: h ^= data[1] << 8;
35     case 1: h ^= data[0];
36             h *= m;
37   };
38   h ^= h >> 13;
39   h *= m;
40   h ^= h >> 15;
41   return h;
42 }
```

2. compute_hash_flags — Compute the flag by the generated two hashes, the corresponding task block in fig 2. is the "compute" block

```
44 void compute_hash_flags (
45         hls::stream<parallel_flags_t>& flag_stream,
46         hls::stream<parallel_words_t>& word_stream,
47         unsigned int              bloom_filter_local[PARALLELISATION][bloom_filter_size],
48         unsigned int              total_size)
49 {
50   compute_flags: for(int i=0; i<total_size/PARALLELISATION; i++)
51   {
52     #pragma HLS LOOP_TRIPCOUNT min=1 max=t_size/pf
53     parallel_words_t parallel_entries = word_stream.read();
54     parallel_flags_t inh_flags = 0;
55
56     compute_hash_flags_loop: for (unsigned int j=0; j<PARALLELISATION; j++)
57     {
58 #pragma HLS UNROLL
59
60       unsigned int curr_entry = parallel_entries(31+j*32, j*32);
61       unsigned int frequency = curr_entry & 0x00ff;
62       unsigned int word_id = curr_entry >> 8;
63       unsigned hash_pu = MurmurHash2(word_id, 3, 1);
64       unsigned hash_lu = MurmurHash2(word_id, 3, 5);
65       bool doc_end= (word_id==docTag);
66       unsigned hash1 = hash_pu&hash_bloom;
67       bool inh1 = (!doc_end) && (bloom_filter_local[j][ hash1 >> 5 ] & ( 1 << (hash1 & 0x1f)));
68       unsigned hash2=(hash_pu+hash_lu)&hash_bloom;
69       bool inh2 = (!doc_end) && (bloom_filter_local[j][ hash2 >> 5 ] & ( 1 << (hash2 & 0x1f)));
70
71       inh_flags(7+j*8, j*8) = (inh1 && inh2) ? 1 : 0;
72     }
73
74     flag_stream.write(inh_flags);
75   }
76 }
77
```

3. compute_hash_flags_dataflow — implementing the flow in fig 2., this function will resize the stream data to desired number of words (e.g. Parallelization=4 words) and then feed into the function "compute_hash_flags":

```
77
78 void compute_hash_flags_dataflow(
79         ap_uint<512>*   output_flags,
80         ap_uint<512>*   input_words,
81         unsigned int    bloom_filter[PARALLELISATION][bloom_filter_size],
82         unsigned int    total_size)
83 {
84 #pragma HLS DATAFLOW
85
86     hls :: stream<ap_uint<512> >    data_from_gmem;
87     hls :: stream<parallel_words_t> word_stream;
88     hls :: stream<parallel_flags_t> flag_stream;
89     hls :: stream<ap_uint<512> >    data_to_gmem;
90
91     // Burst read 512-bit values from global memory over AXI interface
92     hls_stream :: buffer(data_from_gmem, input_words, total_size/(512/32));
93
94     // Form a stream of parallel words from stream of 512-bit values
95     // Going from Wi=512 to Wo= 256
96     hls_stream :: resize(word_stream, data_from_gmem, total_size/(512/32));
97
98     // Process stream of parallel word : word_stream is of 2k (32*64)
99     compute_hash_flags(flag_stream, word_stream, bloom_filter, total_size);
100
101     // Form a stream of 512-bit values from stream of parallel flags
102     // Going from Wi=64 to Wo=512
103     hls_stream :: resize(data_to_gmem, flag_stream, total_size/(512/8));
104
105     // Burst write 512-bit values to global memory over AXI interface
106     hls_stream :: buffer(output_flags, data_to_gmem, total_size/(512/8));
107 }
108
```

4. runOnfpga — The kernel top function, it'll load the bloom filter's coefficients first, and then perform "compute_hash_flags_dataflow" where the function "computer_hash_flags" is called:

```
108
109 extern "C"
110 {
111
112   void runOnfpga (
113         ap_uint<512>*  output_flags,
114         ap_uint<512>*  input_words,
115         unsigned int*  bloom_filter,
116         unsigned int   total_size,
117         bool           load_filter)
118   {
119   #pragma HLS INTERFACE ap_ctrl_chain port=return          bundle=control
120   #pragma HLS INTERFACE m_axi         port=output_flags    bundle=maxiport0   offset=slave
121   #pragma HLS INTERFACE m_axi         port=input_words     bundle=maxiport0   offset=slave
122   #pragma HLS INTERFACE m_axi         port=bloom_filter    bundle=maxiport1   offset=slave
123
124     static unsigned int bloom_filter_local[PARALLELISATION][bloom_filter_size];
125   #pragma HLS ARRAY_PARTITION variable=bloom_filter_local complete dim=1
126 printf("From runOnfpga : Total_size = %d\n", total_size);
127
128     if(load_filter==true)
129     {
130       read_bloom_filter: for(int index=0; index<bloom_filter_size; index++) {
131   #pragma HLS PIPELINE II=1
132         unsigned int tmp = bloom_filter[index];
133         for (int j=0; j<PARALLELISATION; j++) {
134           bloom_filter_local[j][index] = tmp;
135         }
136       }
137     }
138
139     compute_hash_flags_dataflow(
140       output_flags,
141       input_words,
142       bloom_filter_local,
143       total_size);
144   }
145 }
146
```

The host function includes two functions runOnFPGA & runOnCPU, the former will instantiate the kernel "runOnfpga" and perform data transmission between host and kernel:

```
159    std::cout << "Initializing data"<< endl;
160    block_size = num_iter*64;
161    setupData();
162
163    //printf ("Sending data on FPGA  Doc_sizes = %lu\n ", doc_sizes.size());
164    //std::cout << "Sending data on FPGA  Doc_sizes " <<  doc_sizes.size() << endl;
165    //std::cout << "Input Doc Words " <<  input_doc_words.size() << endl;
166    runOnFPGA(
167        doc_sizes.data(),
168        input_doc_words.data(),
169        bloom_filter.data(),
170        profile_weights.data(),
171        fpga_profileScore.data(),
172        total_num_docs,
173        size,
174        num_iter) ;
175
176    runOnCPU(
177        doc_sizes.data(),
178        input_doc_words.data(),
179        bloom_filter.data(),
180        profile_weights.data(),
181        cpu_profileScore.data(),
182        total_num_docs,
183        size) ;
184
185    printf("-------------------------------------------------------------\n");
186
187    for (unsigned doci = 0; doci < total_num_docs; doci++)
188    {
189        if (cpu_profileScore[doci] ≠ fpga_profileScore[doci]) {
190            std::cout << " Verification: FAILED "<< endl  << " : doc[" << doci << "]" << " score: CPU = " << cpu_profileScore[doci]<< ", FPGA = "<< fpga_profileScore[doci]  <<  endl;
191            return 0;
192        }
193    }
194    cout << " Verification: PASS" << endl;
195    cout << endl;
196
197    return 0;
198 }
199
200
```

////////////////////////////////////////////////////////////////////////////////////////////////////////////////

runOnFPGA creates buffer and setups kernel:

```
37    // Boilerplate code to load the FPGA binary, create the kernel and command queue
38    vector<cl::Device> devices = xcl::get_xil_devices();
39    cl::Device device = devices[0];
40    cl::Context context(device);
41    cl::CommandQueue q(context,device, CL_QUEUE_PROFILING_ENABLE|CL_QUEUE_OUT_OF_ORDER_EXEC_MODE_ENABLE);
42
43    string run_type = xcl::is_emulation()?(xcl::is_hw_emulation()?"hw_emu":"sw_emu"):"hw";
44    string binary_file = kernel_name + "_" + run_type + ".xclbin";
45    cl::Program::Binaries bins = xcl::import_binary_file(binary_file);
46    cl::Program program(context, devices, bins);
47    cl::Kernel kernel(program,kernel_name_charptr,NULL);
48
49    unsigned int total_size = total_doc_size;
50    unsigned char* output_inh_flags = (unsigned char*)aligned_alloc(4096, total_size*sizeof(unsigned char));
51    bool load_filter = true;
52
53    // Create buffers
54    cl::Buffer buffer_bloom_filter(context, CL_MEM_USE_HOST_PTR | CL_MEM_READ_ONLY, bloom_filter_size*sizeof(uint),bloom_filter);
55    cl::Buffer buffer_input_doc_words(context, CL_MEM_USE_HOST_PTR | CL_MEM_READ_ONLY, total_size*sizeof(uint),input_doc_words);
56    cl::Buffer buffer_output_inh_flags(context, CL_MEM_USE_HOST_PTR | CL_MEM_WRITE_ONLY, total_size*sizeof(unsigned char),output_inh_flags);
57
58    // Set buffer kernel arguments (needed to migrate the buffers in the correct memory)
59    kernel.setArg(0, buffer_output_inh_flags);
60    kernel.setArg(1, buffer_input_doc_words);
61    kernel.setArg(2, buffer_bloom_filter);
62
63    double mbytes_total  = (double)(total_doc_size * sizeof(int)) / (double)(1000*1000);
64    printf(" Processing %.3f MBytes of data\n", mbytes_total);
65    printf(" Single_Buffer: Running with a single buffer of %.3f MBytes for FPGA processing\n",mbytes_total);
66
67    // Create events for read,compute and write
68
69        vector<cl::Event> wordWait;
70        vector<cl::Event> krnlWait;
71        vector<cl::Event> flagWait;
72    cl::Event buffDone, krnlDone, flagDone;
73
74    printf("-------------------------------------------------------------\n");
75
76
77    chrono::high_resolution_clock::time_point t1, t2;
78    t1 = chrono::high_resolution_clock::now();
79
80
81    // Load the bloom filter and input document words buffers
82    q.enqueueMigrateMemObjects({buffer_bloom_filter, buffer_input_doc_words}, 0,NULL,&buffDone);
83        wordWait.push_back(buffDone);
84
85    // Start the FPGA compute
86    load_filter = true;
87    kernel.setArg(3, total_size);
88    kernel.setArg(4, load_filter);
89    q.enqueueTask(kernel,&wordWait,&krnlDone);
90        krnlWait.push_back(krnlDone);
91
```

\Users\USER\AppData\Local\Temp\Mxt221\RemoteFiles\7537 UNIX     C/C++           142 lines      Row #45       Col #35

////////////////////////////////////////////////////////////////////////////////////////////////////////////////

Let's use the following HLS synthesis report to see latencies of the functions
- The compute_hash_flags latency reported is 875,011 cycles, because there are in total 35,000,000 words, computed with 4 words in parallel, and with task-level parallelism pipelining we can achieve about 35M/4 cycles.
- The read_bloom_filter task takes about 16000 cycles because the loop is iterated over 16,000 cycles reading 32-bits data from the Bloom filter coefficients. This is also due to pipelining with the code #pragma HLS PIPELINE II=1 when loading the filter coefficients.



There are some issues when performing this lab with U50 (originally use U200), similar problem has been discussed in the thread https://github.com/Xilinx/Vitis-Tutorials/issues/95, and the solution was also provided.

Use the information provided on the lab website, after running on hardware:

1. Parallelization = 4 words simultaneously for computing their hashes:

```
Loading runOnfpga_hw.xclbin
Processing 1398.903 MBytes of data
  Running with a single buffer of 1398.903 MBytes for FPGA processing
-----------------------------------------------------------------
Executed FPGA accelerated version  |   838.5898 ms   ( FPGA 447.964 ms )
Executed Software-Only version     |  3187.0354 ms
-----------------------------------------------------------------
Verification: PASS
```

2. Parallelization = 8 words simultaneously for computing their hashes:

```
Single_Buffer: Running with a single buffer of 1398.903 MBytes for FPGA processing
-----------------------------------------------------------------
 Executed FPGA accelerated version  |   739.4475 ms   ( FPGA 315.475 ms )
 Executed Software-Only version     |  3053.9516 ms
-----------------------------------------------------------------
 Verification: PASS
```

3. Parallelization = 16 words simultaneously for computing their hashes:

```
Processing 1398.903 MBytes of data
 Single_Buffer: Running with a single buffer of 1398.903 MBytes for FPGA processing
-----------------------------------------------------------------
 Executed FPGA accelerated version  |   694.6162 ms   ( FPGA 270.275 ms )
 Executed Software-Only version     |  3052.5701 ms
-----------------------------------------------------------------
 Verification: PASS
```

we can see there is a significant improvement of execution time using FPGA compared to software-only version.

timeline trace of Parallelization = 8 words:



The figure also shows a similar task flow as in fig.1.

From this, we can see the tasks are still sequentially implemented, therefore if we divide the "write" into several parts, overlapping "write" and "runOnfpga", we can probably get better execution time. Next, we'll explore using multiple sub-buffers for better performance.

## Using multiple sub-buffers (split buffers):

The idea of splitting the buffer into several sub-buffers is shown in fig.3, where the input document buffer is split into 2 sub-buffers.
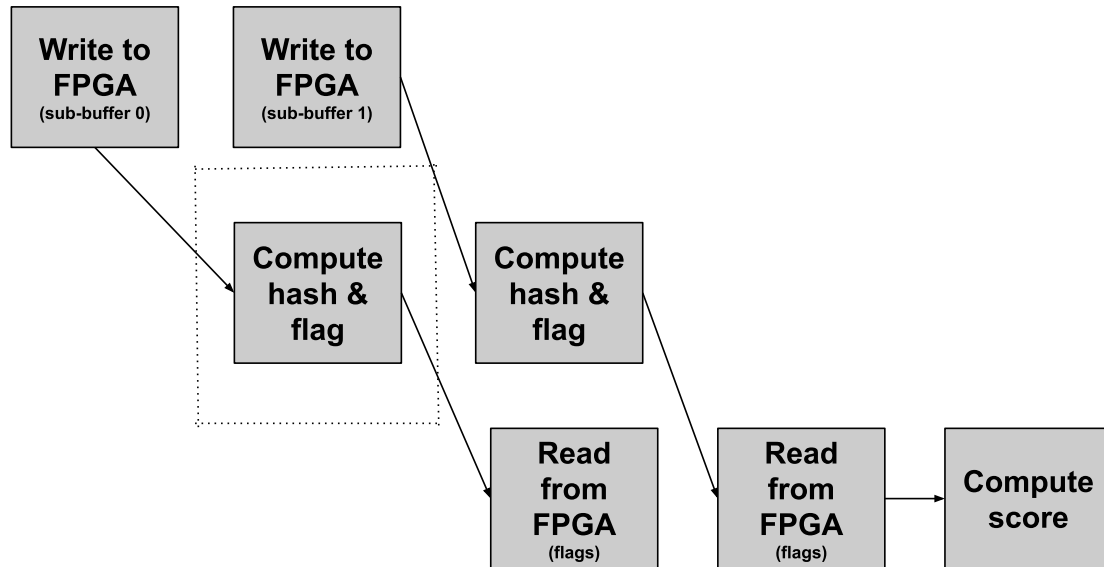


fig.3 implementation with multiple sub-buffer

The following code declares sub-buffer regions to specify offset and size of sub-buffer and creates two sub-buffers from buffers based on sub-buffer regions:

```cpp
// Make buffers resident in the device
q.enqueueMigrateMemObjects({buffer_bloom_filter, buffer_input_doc_words, buffer_output_inh_flags}, CL_MIGRATE_MEM_OBJECT_CONTENT_

// Specify size of sub-buffers, one for each transaction
unsigned subbuf_doc_sz = total_doc_size/2;
unsigned subbuf_inh_sz = total_doc_size/2;

// Declare sub-buffer regions to specify offset and size of sub-buffer
cl_buffer_region subbuf_inh_info[2];
cl_buffer_region subbuf_doc_info[2];

// Declare sub-buffers
cl::Buffer subbuf_inh_flags[2];
cl::Buffer subbuf_doc_words[2];

 // Specify offset and size of sub-buffers
subbuf_inh_info[0]={0, subbuf_inh_sz*sizeof(char)};
subbuf_inh_info[1]={subbuf_inh_sz*sizeof(char), subbuf_inh_sz*sizeof(char)};
subbuf_doc_info[0]={0, subbuf_doc_sz*sizeof(uint)};
subbuf_doc_info[1]={subbuf_doc_sz*sizeof(uint), subbuf_doc_sz*sizeof(uint)};

// Create sub-buffers from buffers based on sub-buffer regions
subbuf_inh_flags[0] = buffer_output_inh_flags.createSubBuffer(CL_MEM_WRITE_ONLY, CL_BUFFER_CREATE_TYPE_REGION, &subbuf_inh_info[0
subbuf_inh_flags[1] = buffer_output_inh_flags.createSubBuffer(CL_MEM_WRITE_ONLY, CL_BUFFER_CREATE_TYPE_REGION, &subbuf_inh_info[1
subbuf_doc_words[0] = buffer_input_doc_words.createSubBuffer (CL_MEM_READ_ONLY,  CL_BUFFER_CREATE_TYPE_REGION, &subbuf_doc_info[0
subbuf_doc_words[1] = buffer_input_doc_words.createSubBuffer (CL_MEM_READ_ONLY,  CL_BUFFER_CREATE_TYPE_REGION, &subbuf_doc_info[1

printf("\n");
 double mbytes_total  = (double)(total_doc_size * sizeof(int)) / (double)(1000*1000);
double mbytes_block  = mbytes_total / 2;
printf(" Processing %.3f MBytes of data\n", mbytes_total);
printf(" Splitting data in 2 sub-buffers of %.3f MBytes for FPGA processing\n", mbytes_block);
```

```
// Create Events to co-ordinate read,compute and write for each iteration
vector<cl::Event> wordWait;
vector<cl::Event> krnlWait;
vector<cl::Event> flagWait;
```

First iteration for setting kernel arguments, reading, enqueuing kernel and writing:

```
//  Set Kernel Arguments, Read, Enqueue Kernel and Write for first iteration
total_size = total_doc_size/2;
load_filter=false;
kernel.setArg(3, total_size);
kernel.setArg(4, load_filter);
kernel.setArg(0, subbuf_inh_flags[0]);
kernel.setArg(1, subbuf_doc_words[0]);
q.enqueueMigrateMemObjects({subbuf_doc_words[0]}, 0, &wordWait, &buffDone);
wordWait.push_back(buffDone);
q.enqueueTask(kernel, &wordWait, &krnlDone);
krnlWait.push_back(krnlDone);
q.enqueueMigrateMemObjects({subbuf_inh_flags[0]}, CL_MIGRATE_MEM_OBJECT_HOST, &krnlWait, &flagDone);
flagWait.push_back(flagDone);
```

Second iteration:

```
//  Set Kernel Arguments, Read, Enqueue Kernel and Write for second iteration
total_size = total_doc_size/2;
load_filter=false;
kernel.setArg(3, total_size);
kernel.setArg(4, load_filter);
kernel.setArg(0, subbuf_inh_flags[1]);
kernel.setArg(1, subbuf_doc_words[1]);
q.enqueueMigrateMemObjects({subbuf_doc_words[1]}, 0, &wordWait, &buffDone);
wordWait.push_back(buffDone);
q.enqueueTask(kernel, &wordWait, &krnlDone);
krnlWait.push_back(krnlDone);
q.enqueueMigrateMemObjects({subbuf_inh_flags[1]}, CL_MIGRATE_MEM_OBJECT_HOST, &krnlWait, &flagDone);
flagWait.push_back(flagDone);
```

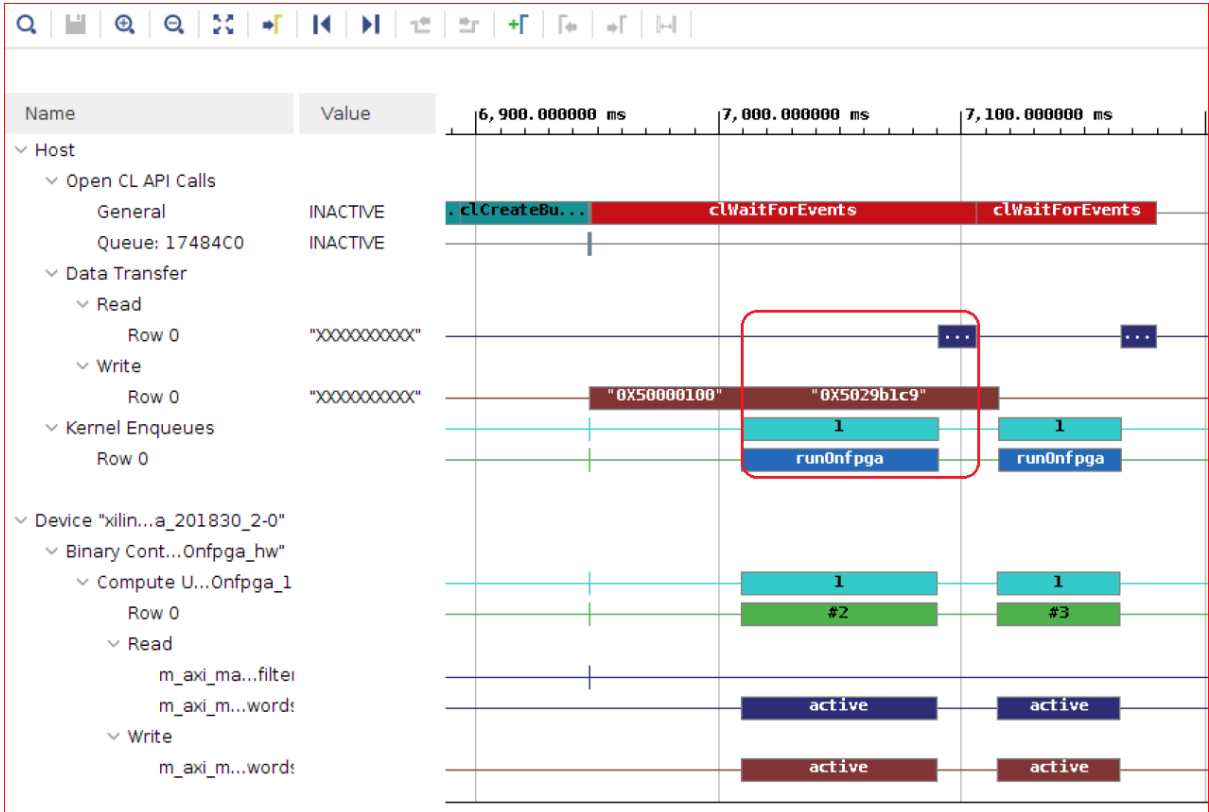Use wait() function to wait until results are transmitted to host:

```
// Wait until all results are copied back to the host before doing the post-processing
flagWait[0].wait();
flagWait[1].wait();
```

Running this optimization with Parallelization = 8 words yields:

```
Processing 1398.903 MBytes of data
Split_buffer : Splitting data in 2 sub-buffers of 699.452 MBytes for FPGA processing
----------------------------------------------------------------
Executed FPGA accelerated version |   734.0995 ms   ( FPGA 262.363 ms )
Executed Software-Only version    |   3246.2145 ms
----------------------------------------------------------------
Verification: PASS
```

This improves the FPGA execution time from 315 ms to 262 ms, while the overall performance does not improve much (739 ms vs 734 ms)

timeline trace:



The timeline trace is as expected as the idea in fig.3.

Now, we can change the code into more generic form, which allow us to split the buffer into desired number of sub-buffers:

```cpp
// Specify size of sub buffers for each iteration
unsigned subbuf_doc_sz = total_doc_size/num_iter;
unsigned subbuf_inh_sz = total_doc_size/num_iter;

// Declare sub buffer regions to specify offset and size for each iteration
cl_buffer_region subbuf_inh_info[num_iter];
cl_buffer_region subbuf_doc_info[num_iter];

// Declare sub buffers
cl::Buffer subbuf_inh_flags[num_iter];
cl::Buffer subbuf_doc_words[num_iter];

// Define sub buffers from buffers based on sub-buffer regions
for (int i=0; i<num_iter; i++)  {
  subbuf_inh_info[i]={i*subbuf_inh_sz*sizeof(char), subbuf_inh_sz*sizeof(char)};
  subbuf_doc_info[i]={i*subbuf_doc_sz*sizeof(uint), subbuf_doc_sz*sizeof(uint)};
  subbuf_inh_flags[i] = buffer_output_inh_flags.createSubBuffer(CL_MEM_WRITE_ONLY, CL_BUFFER_CREATE_TYPE_REGION, &subbuf_inh_in
  subbuf_doc_words[i] = buffer_input_doc_words.createSubBuffer (CL_MEM_READ_ONLY,  CL_BUFFER_CREATE_TYPE_REGION, &subbuf_doc_in
}

printf("\n");
double mbytes_total  = (double)(total_doc_size * sizeof(int)) / (double)(1000*1000);
double mbytes_block  = mbytes_total / num_iter;
printf(" Processing %.3f MBytes of data\n", mbytes_total);
if (num_iter>1) {
  printf(" Splitting data in %d sub-buffers of %.3f MBytes for FPGA processing\n", num_iter, mbytes_block);
}
```
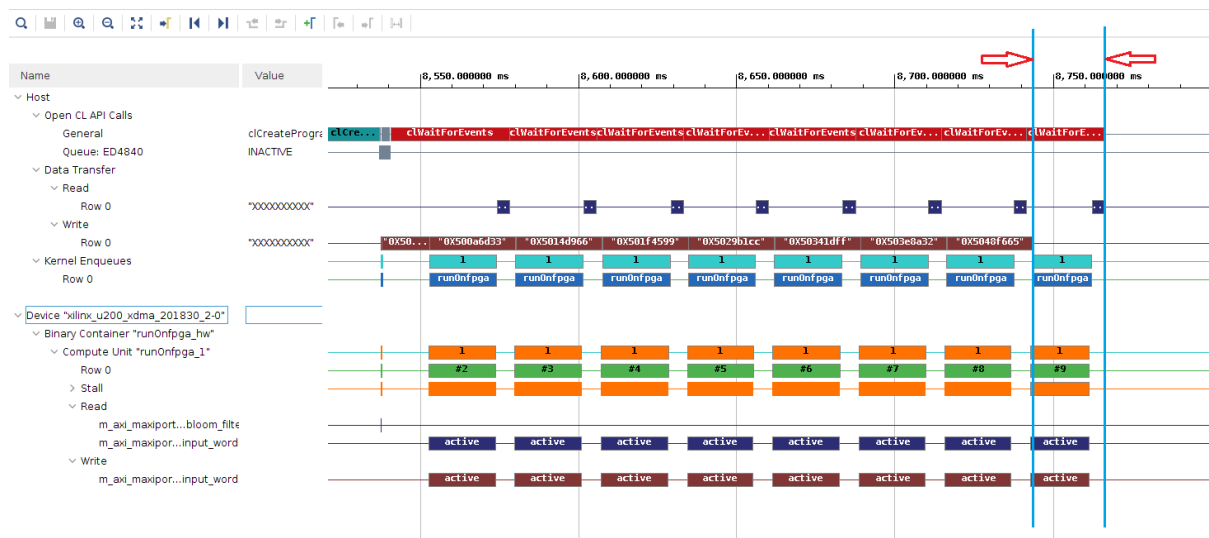
```cpp
// Set Kernel arguments. Read, Enqueue Kernel and Write for each iteration
for (int i=0; i<num_iter; i++)
{
  cl::Event buffDone, krnlDone, flagDone;
  total_size = subbuf_doc_info[i].size / sizeof(uint);
  load_filter = false;
  kernel.setArg(0, subbuf_inh_flags[i]);
  kernel.setArg(1, subbuf_doc_words[i]);
  kernel.setArg(3, total_size);
  kernel.setArg(4, load_filter);
  q.enqueueMigrateMemObjects({subbuf_doc_words[i]}, 0, &wordWait, &buffDone);
  wordWait.push_back(buffDone);
  q.enqueueTask(kernel, &wordWait, &krnlDone);
  krnlWait.push_back(krnlDone);
  q.enqueueMigrateMemObjects({subbuf_inh_flags[i]}, CL_MIGRATE_MEM_OBJECT_HOST, &krnlWait, &flagDone);
  flagWait.push_back(flagDone);
}
```

```cpp
// Wait until all results are copied back to the host before doing the post-processing
for (int i=0; i<num_iter; i++)
{
  flagWait[i].wait();
}
```
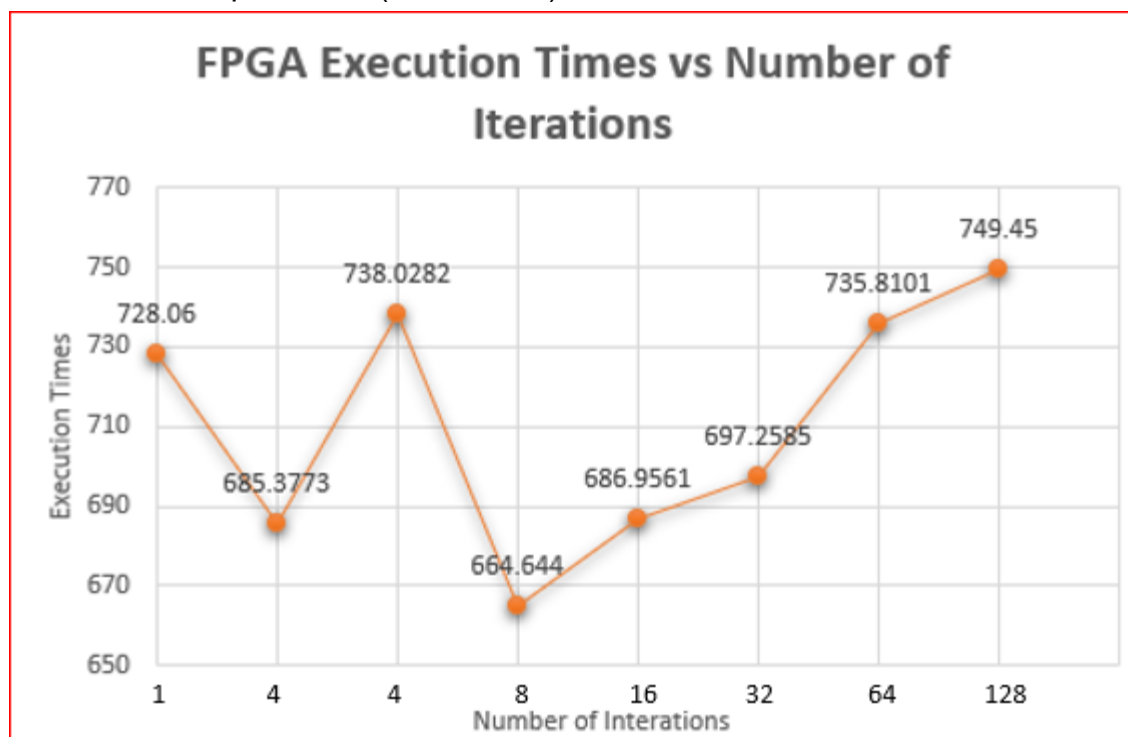
The timeline trace of splitting into 8 sub-buffers:



We can see the red flags clWaitForEvents last the whole kernel execution time, a further optimization idea is that we don't want to wait for the whole events, but compute score right after reading a subbuffer of flags.

We can run multiple sub-buffers cases and plot the following FPGA execution time v.s. number of split buffers(sub-buffers):



source:https://github.com/Xilinx/Vitis-Tutorials/blob/2022.1/Hardware_Acceleration/Design_Tutorials/02-bloom/5_data-movement.md

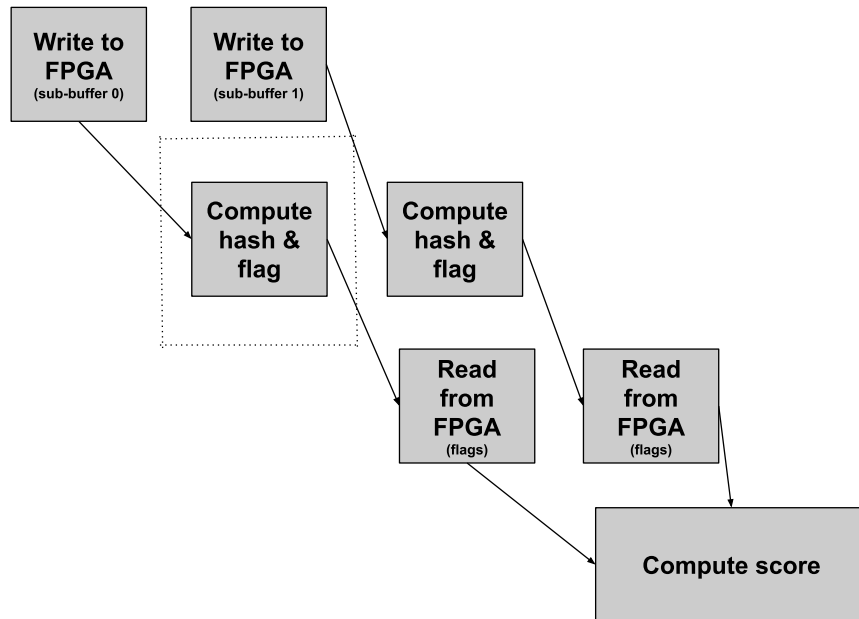We can see the case occurs when we split the buffer into 8 parts.

fig.4 implementation with multiple sub-buffer, overlapping compute score and read from FPGA

modify the compute score part with the following variables added:

```
// Create variables to keep track of number of words needed by CPU to compute score and number of words processed by FPGA such
unsigned int curr_entry;
unsigned char inh_flags;
unsigned int  available = 0;
unsigned int  needed = 0;
unsigned int  iter = 0;
```

```
for(unsigned int doc=0, n=0; doc<total_num_docs;doc++)
{
  unsigned long ans = 0;
  unsigned int size = doc_sizes[doc];

  // Calculate size by needed by CPU for processing next document score
  needed += size;

  // Check if flags processed by FPGA is greater than needed by CPU. Else, block CPU
  // Update the number of available words and sub-buffer count(iter)

  if (needed > available)
  {
    flagWait[iter].wait();
    available += subbuf_doc_info[iter].size / sizeof(uint);
    iter++;
  }

  for (unsigned i = 0; i < size ; i++, n++)
  {
    curr_entry = input_doc_words[n];
    inh_flags  = output_inh_flags[n];

    if (inh_flags)
    {
      unsigned frequency = curr_entry & 0x00ff;
      unsigned word_id = curr_entry >> 8;
      ans += profile_weights[word_id] * (unsigned long)frequency;
    }
  }
  profile_score[doc] = ans;
}
```
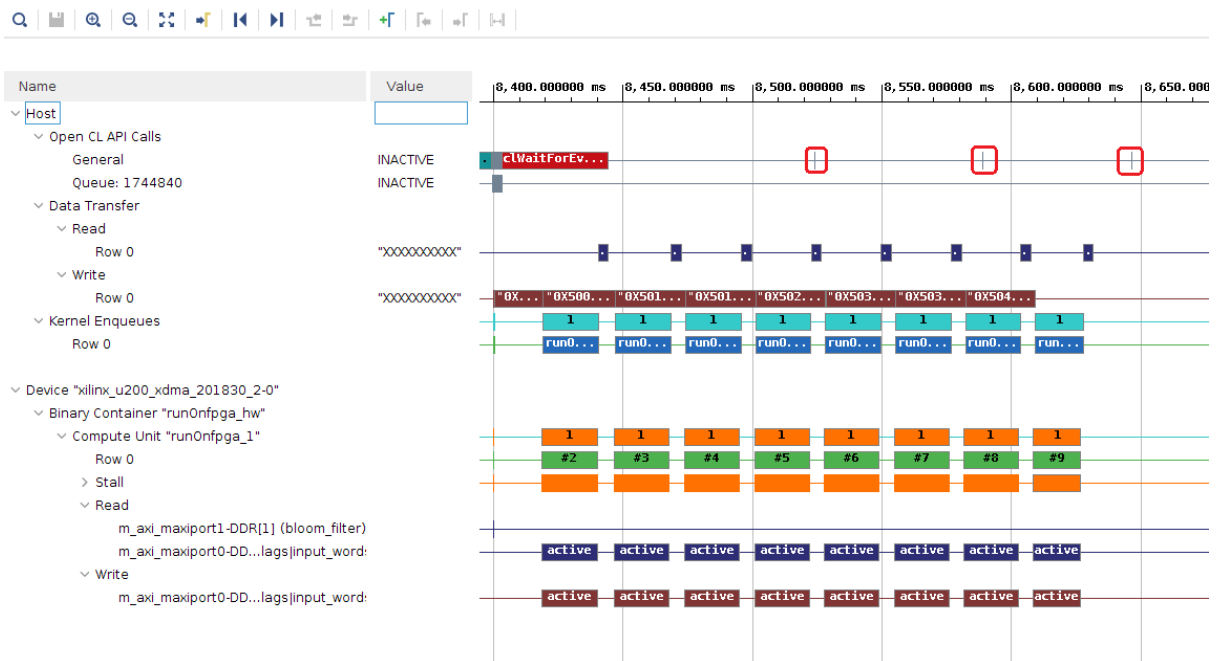
The needed variable indicates the number of flags needed (accumulated) to score the current document.

The available variable indicates the number of flags that has been generated.

needed > available means we should wait until more output flags are generated.

```
Processing 1398.905 MBytes of data
Splitting data in 8 sub-buffers of 174.863 MBytes for FPGA processing
----------------------------------------------------------------
Executed FPGA accelerated version  |   427.1341 ms   ( FPGA 230.345 ms )
Executed Software-Only version     |   3057.6307 ms
----------------------------------------------------------------
Verification: PASS
```

We can see the wait time in timeline trace is largely reduced, and it gives much better performance.

## Using multiple DDR banks:

Since there are memory contentions because the host and kernel both accessed the same bank at the same time, we can further optimize the performance by using multiple DDR banks in ping-pong fashion, meaning that:

- The host write words to DDR bank 1 and bank 2 alternatively
- when the host is writing DDR bank 1, the kernel is reading DDR bank 2
- when the host is writing DDR bank 2, the kernel is reading DDR bank 1

v++ linking option --connectivity.sp <arg> can be used to connect kernel arguments to specific memory resources. This option can be specified in a configuration file under the [connectivity] section head using the following format:

```
[connectivity]
sp=vadd_1.A:DDR[0:3]
sp=vadd_1.B:HBM[0:31]
sp=vadd_1.C:PLRAM[2]
```

NOTE: Any argument not explicitly mapped to a memory resource through the --connectivity.sp option is automatically connected to an available memory resource during the build process.

In this lab, this option is specified in connectivity.cfg:

```
1 [connectivity]
2 sp=runOnfpga_1.input_words:DDR[1:2]
3
```

The -sp option instructs the v++ linker that input_words is connected to both DDR banks 1 and 2. We should rebuild the kernel since the connectivity is changed.

There are two Xilinx extension pointer objects (cl_mem_ext_ptr_t) created. The corresponding flags will determine which DDR bank the buffer will be sent to, so that the kernel can access it.

```
cl_mem_ext_ptr_t buffer_words_ext[2];

buffer_words_ext[0].flags = 1 | XCL_MEM_TOPOLOGY; // DDR[1]
buffer_words_ext[0].param = 0;
buffer_words_ext[0].obj   = input_doc_words;
buffer_words_ext[1].flags = 2 | XCL_MEM_TOPOLOGY; // DDR[2]
buffer_words_ext[1].param = 0;
buffer_words_ext[1].obj   = input_doc_words;
```

modify the code as following:

```
buffer_doc_words[0] = cl::Buffer(context, CL_MEM_EXT_PTR_XILINX | CL_MEM_USE_HOST_PTR | CL_MEM_READ_ONLY, total_size*sizeof(uint), &buffer_words_ext[0]);
buffer_doc_words[1] = cl::Buffer(context, CL_MEM_EXT_PTR_XILINX | CL_MEM_USE_HOST_PTR | CL_MEM_READ_ONLY, total_size*sizeof(uint), &buffer_words_ext[1]);
buffer_inh_flags    = cl::Buffer(context, CL_MEM_USE_HOST_PTR | CL_MEM_WRITE_ONLY, total_size*sizeof(char),output_inh_flags);
buffer_bloom_filter = cl::Buffer(context, CL_MEM_USE_HOST_PTR | CL_MEM_READ_ONLY, bloom_filter_size*sizeof(uint),bloom_filter);

// Set buffer kernel arguments (needed to migrate the buffers in the correct memory)
kernel.setArg(0, buffer_inh_flags);
kernel.setArg(1, buffer_doc_words[0]);
kernel.setArg(2, buffer_bloom_filter);

// Make buffers resident in the device
q.enqueueMigrateMemObjects({buffer_bloom_filter, buffer_doc_words[0], buffer_doc_words[1], buffer_inh_flags}, CL_MIGRATE_MEM_OBJE

// Create sub-buffers, one for each transaction
unsigned subbuf_doc_sz = total_doc_size/num_iter;
unsigned subbuf_inh_sz = total_doc_size/num_iter;

cl_buffer_region subbuf_inh_info[num_iter];
cl_buffer_region subbuf_doc_info[num_iter];
cl::Buffer subbuf_inh_flags[num_iter];
cl::Buffer subbuf_doc_words[num_iter];

for (int i=0; i<num_iter; i++) {
    subbuf_inh_info[i]={i*subbuf_inh_sz*sizeof(char), subbuf_inh_sz*sizeof(char)};
    subbuf_doc_info[i]={i*subbuf_doc_sz*sizeof(uint), subbuf_doc_sz*sizeof(uint)};
    subbuf_inh_flags[i] = buffer_inh_flags.createSubBuffer(CL_MEM_WRITE_ONLY, CL_BUFFER_CREATE_TYPE_REGION, &subbuf_inh_info[i]);
    // The doc words sub-buffers will be alternating in DDR[1] and DDR[2]
    subbuf_doc_words[i] = buffer_doc_words[i%2].createSubBuffer (CL_MEM_READ_ONLY,  CL_BUFFER_CREATE_TYPE_REGION, &subbuf_doc_inf
}
```

The doc words sub-buffers will be alternating in DDR[1] and DDR[2].

```
Processing 1398.905 MBytes of data
MultiDDR- Splitting data in 8 sub-buffers of 174.863 MBytes for FPGA processing
----------------------------------------------------------------
Executed FPGA accelerated version  |   426.6388 ms   ( FPGA 175.113 ms )
Executed Software-Only version     |   3058.8499 ms
----------------------------------------------------------------
Verification: PASS
```
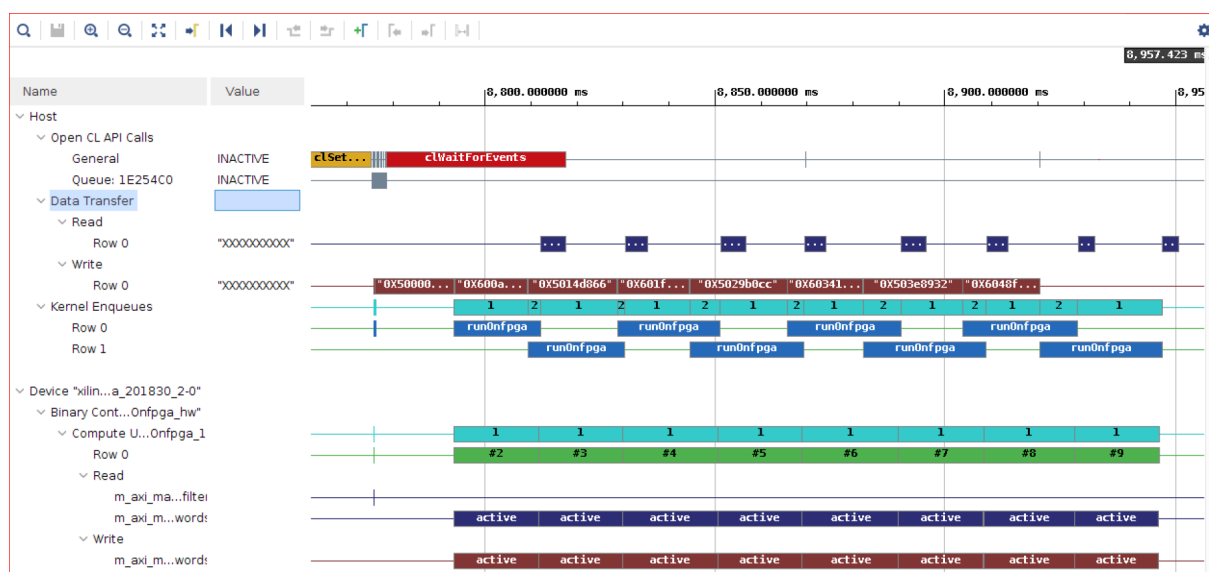
we can see the FPGA execution time has improved.



In WRITE transactions, we can observe that when the host is writing DDR bank 1 and bank 2 alternatively.

## What has been learned:

In this lab, I not only learned what is bloom filter but also learned a lot of optimization techniques, plus, Opencl API is reviewed thoroughly in the lab. I think this certainly gives me much intuition when design my own work.