# GPU Programming CS6023

Fall 2018

## Assignment 3

## Report

Author: Varun Sundar

EE16B068

21st September 2018

# 1 Implementation and Technical Notes

The code was run on a *GTX 1080 Ti* and not the cluster of *K'40s*. This could be a reason for variance in the optimal configuration found for the GPUs. However, the code does not make any hardware assumptions and can therefore be run on any suitable cluster.

CUDA version 9.0 was used to compile the code, with `nvcc` as the device compiler and `gcc` as the host compiler.

# 2 Analysis of Problem Statement

We are required to build N-count-grams of a given text (development on Shakespeare's *Procreation Sonnets* ; test elsewhere). This is associated with a *MAXWORD* limit.

Note that, unlike usual, we treat the whole text as a single sentence.

# 3 Pseudo Code

1. Read file, check if string is a word.

2. Save words into an array.

3. For a given N, generate the sliding window array, with count. *(Can this be optimised?)*

4. Bin this histogram. *(This can be optimised)*

We shall implement ideas from shared memory atomics, privatisation etc, and also consider recent paper developments for this. Wherever this has helped, we shall cite the relevant papers for credibility.

# 4 Implementing Binning

We largely follow the following method:

- Declare shared memory with private output histogram

- Cooperatively initialize the histogram to 0

- Synchronize

- Identify the index/indices of the input on which to operate. For each

    - Access each input item such that warps have coalesced access.
    - Use atomic add to update appropriate bin in the output histogram

- Cooperatively update the global output histogram with local one with

- atomic add

Now when N>5, we notice that shared memory can no longer fit our histograms, so we split the private histograms into private sub-histograms:

- Declare shared memory with private sub-histogram

- Cooperatively initialize the histogram to 0

- Synchronize

- Identify the index/indices of the input on which to operate. For each

  – Ensure the global index you want to update falls in this sub-histogram.
  – Access each input item such that warps have coalesced access.
  – Use atomic add to update appropriate bin in the output histogram

- Cooperatively update the global output histogram with local one with

- atomic add

# 5   Code and Timings

Histogram binning via thread-wise access to shared memory was used.

We implement two kernels: one for the case of $N < 5$ and one for the other (case of full privat. The results may be replicate histograms and the case of parallel ed by running,

```
./ee16b068.cu
```

## 5.1   Run Timing and Interpretation

**Figure-1** summarises the run-times (GPU kernel) for various $N$, tested on the provided Shake-spearean text.

This can be attributed to the fact that shared atomic collisions grow as the number of simultaneous additions increase, but a $log(n)$ dependency may be attributed to increased loops in each block.

For questions 2,3 , we shall consider column major timing only.

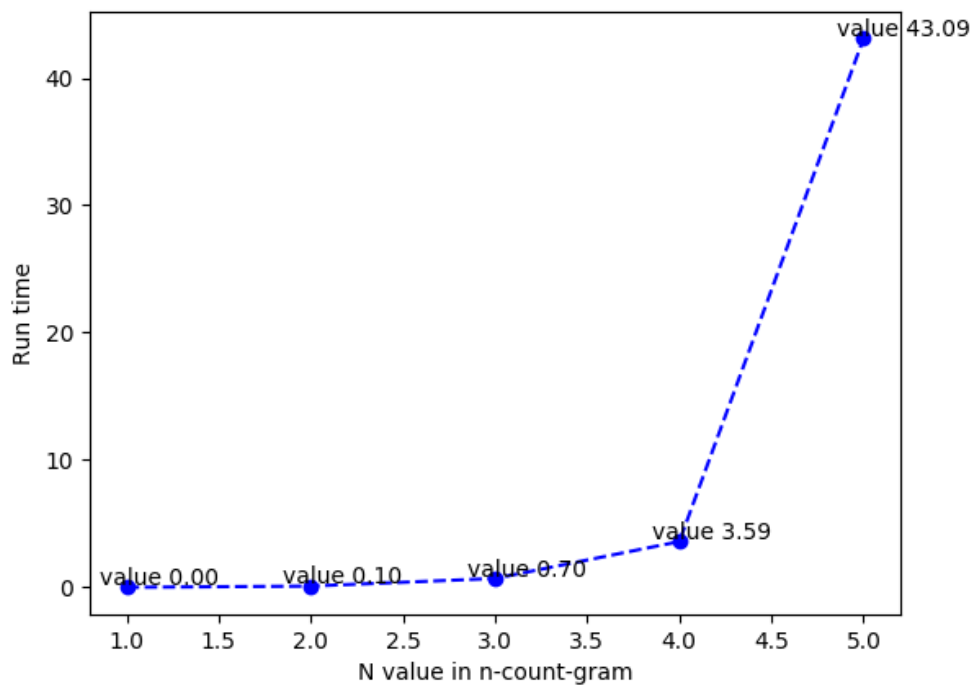## 5.2   Code Blocks for pre-processing, word-counts (pertinent only)

Figure 1: N-count-gram runtimes on**GTX1080 Ti**

```
2   int checkWord(char* word,char* words,int* count_array,int offset){
3       // Check if word meets, else pre-process
4       // Args:
5       // word >> word of consideration from fscanf
6       // words >> Array where, every 20 chars is a word
7       // offset >> Which entry to start writting at (modulo 20)
8       // Returns:
9       // new offset
10      // Modifies:
11      // words
12      int loop=0;
13      int count=0;
14
15      for (loop=0;loop<strlen(word)-1;loop++)
16      {
17
18          if (word[loop]=='-')
19          {
20              words[offset*20+loop]=0;
21              printf("Word %s \n",&words[offset*20]);
22              offset+=1;
23              count_array[offset]=count;
24              count=0;
25          }
```

```
26          else{
27              /* Copy character */
28              words[offset*20+loop]=word[loop];
29              count+=1;
30          }
31      }
32
33      if (ispunct((unsigned char)word[strlen(word)-1]))
34          {
35              /* Skip this character */
36              words[offset*20+strlen(word)-1]=0;
37              count_array[offset]=count;
38              offset+=1;
39          }
40      else{
41          words[offset*20+strlen(word)-1]=word[strlen(word)-1];
42          count+=1;
43          words[offset*20+strlen(word)]=0;
44          count_array[offset]=count;
45          offset+=1;
46      }
47      return offset;
48
49 }
```

## 5.3   Code Blocks for N=1,2,3,4 (pertinent only)

```
51 __global__ void nCountGram(int* d_count, int* d_hist, int N, int
      totalWordCount){
52      extern __shared__ int buffer[];
53      int *temp = &buffer[0];
54
55      //__shared__ int temp[1024];
56      // Helper var
57      int index, j, p;
58      int a, b;
59
60      a=1;
61      for (p=0;p<N;p++){
62          a*=20;
63      }
64
65      for (p=0; p<a/1024+1; p++){
66          if (threadIdx.x + p*1024< a){
67              temp[threadIdx.x + p*1024] = 0;
68          }
69      }
70
71      __syncthreads();
```

```
 72
 73      int i = threadIdx.x + blockIdx.x * blockDim.x;
 74      int offset = blockDim.x * gridDim.x;
 75
 76      while (i < totalWordCount - N + 1)
 77      {
 78          // Since 0,0 is invalid
 79          index=-1;
 80          b=a/20;
 81          for (j = 0;j < N; j++){
 82              index+=(d_count[i+j])*b;
 83              b/=20;
 84          }
 85          atomicAdd( &temp[index], 1);
 86          i += offset;
 87          //printf("Index %d",index);
 88      }
 89
 90      __syncthreads();
 91
 92      for (p=0;p<a/1024+1;p++){
 93          if (threadIdx.x + p*1024< a){
 94              atomicAdd( &(d_hist[threadIdx.x + p*1024]), temp[threadIdx.x
                      + p*1024] );
 95              if (temp[threadIdx.x+p*1024]>0){
 96                  //printf("Hist val at %d is %d \n",threadIdx.x+p*1024,
                          d_hist[threadIdx.x + p*1024]);
 97              }
 98          }
 99      }
100
101      __syncthreads();
102
103  }
```

## 5.4  Code Blocks for N>=5 (pertinent only)

```
105  __global__ void nCountGram_optimal(int* d_count, int* d_hist, int N, int
         totalWordCount, int sub_hist_size){
106      extern __shared__ int buffer[];
107      int *temp = &buffer[0];
108
109      //__shared__ int temp[1024];
110      // Helper var
111      int index, j, p;
112      int a, b;
113
114      a=1;
115      for (p=0;p<N;p++){
```

```
116          a*=20;
117      }
118
119      for (p=0; p<sub_hist_size/1024 +1; p++){
120          if (threadIdx.x + p*1024 < sub_hist_size){
121              temp[threadIdx.x + p*1024] = 0;
122          }
123      }
124
125      __syncthreads();
126
127      int i = threadIdx.x + blockIdx.x * blockDim.x ;//blockIdx.y*gridDim.
             y;
128      int offset = blockDim.x * gridDim.x*blockIdx.y*gridDim.y;
129
130      while (i < totalWordCount - N + 1)
131      {
132          // Since 0,0 is invalid
133          index=-1;
134          b=a/20;
135          for (j = 0;j < N; j++){
136              index+=(d_count[i+j])*b;
137              b/=20;
138          }
139          if ((index<sub_hist_size*(blockIdx.y+1)) && (index >
                 sub_hist_size*blockIdx.y)){
140              //printf("Index %d",index);
141          atomicAdd( &temp[index - blockIdx.y*sub_hist_size], 1);
142          }
143          i += offset;
144      }
145
146      __syncthreads();
147
148      for (p=0;p<sub_hist_size/1024+1;p++){
149          if (threadIdx.x + p*1024 < sub_hist_size){
150              atomicAdd( &(d_hist[threadIdx.x + sub_hist_size*blockIdx.y +
                     p*1024]), temp[threadIdx.x + p*1024] );
151              if (d_hist[threadIdx.x+ sub_hist_size*blockIdx.y + p
                     *1024]>0){
152                  printf("Hist val at %d is %d \n",threadIdx.x+
                         sub_hist_size*blockIdx.y+p*1024,d_hist[threadIdx.x +
                         sub_hist_size*blockIdx.y+ p*1024]);
153              }
154          }
155      }
156
157      __syncthreads();
158
```

```
159  }
```

# 6   References

- Classroom lectures and Slides.

- CUDA Dev Documentation.