

GPU Programming Assignment - 3

AE13B064

Question - 1

Since there is no `__syncthreads()` operation before the `printf()` 3 cases are possible. 2 warps are possible one with 32 active threads and the other with 2 active threads

Case 1: Both warps finishes the job of increment simultaneously, in this case all there will be 34 - 34's in the output screen.

Case 2: Warp 1 finishes the job (first 32 threads) before warp 2, in this case there will be 32 - 32's and 2 - 34's in the output screen.

Case 3: Warp 2 finishes the job (2 active threads) before warp 1, in this case there will be 2 - 2's and 32 - 34's in the output screen.

Question - 2

Since coalescing is access of memory which is used by the warp threads at single time step, the following code snippet achieves various degrees of coalescing by accessing useful memory for warp threads in different time steps (depending on the degree).

So, the idea for coalesced access of an array (of size 32 for example) is to find the number of time steps required to access the array elements, such that when parameter passed is 1 elements are accessed in 32 steps. Similarly if parameter 2 is passed elements are accessed in 31 time steps. Finally if parameter 32 is passed elements are accessed in a single time step.

```
1
2
3 __global__ void coalescing(int * arr, int param) {
4     int id = threadIdx.x;
5     int timeSteps = 32 - param + 1; // Offset which is required to
        calculate how many access at a time
6
7     for(int i = 0; i < timeSteps; i++){
8         if(id >= i and id < param + i){ // Checks if the thread id is
            with in the range for accessing memory from array at that
            particular time Step
9             arr[id]++; // Process something with the fetched memory
10        }
```

```

11 }
12 }

```

Listing 1: Degree of coalescing Kernel

Testing the Code

The code is tested with following `main.cu` where the kernel call is run for 10000 times to find the average time it took to finish the kernel for various degrees of coalescing.

```

1 #include <cuda.h>
2 #include <stdio.h>
3 #include "timer.h"
4 #include "kernel.h"
5
6 int main() {
7     int N = 32; // Array of size 32
8     int *counter;
9     int *hcounter = (int*)malloc(N*sizeof(int));
10    int i = 0;
11    for(i = 0; i < N; i++){
12        hcounter[i] = i;
13    }
14
15    cudaMalloc(&counter, N*sizeof(int));
16    cudaMemcpy(counter, &hcounter, sizeof(int),
17               cudaMemcpyHostToDevice);
18    GPUSync timer;
19    float t = 0;
20    for(i = 0; i < 10000; i++){
21        timer.Start();
22        coalescing<<<1, N>>>(counter, 16);
23        cudaDeviceSynchronize();
24        timer.Stop();
25        t += 1000*timer.Elapsed();
26    }
27    printf("GPU elapsed %f ms \n", t/10000);
28    return 0;
29 }

```

Listing 2: Degree of coalescing Kernel

Experiment is performed for various degrees of coalescing for increment operation of array values. Following are the results:

Degree of Coalescing	Average time for 10000 calls (ms)
1	20.67
16	17.51
32	14.11

Question - 3

Following is the code snippet for finding the saddle point in the kernel. Logic is to assign single row to each thread. Then, find the minimum in that corresponding thread; store the col index of minimum; find the maximum in that particular column where minimum is found; store the row index where maximum is found in that particular column; if the row index matches with thread id then we found our saddle point

```
1  __global__ void saddleFinder(const int *arr, int N){
2      int id = threadIdx.x;
3      int max = INT_MIN, min = INT_MAX; // Using limits.h
4      int i1 = 0, i2 = 0 ;
5
6
7      __syncthreads(); // It is not required
8
9      // Finding minimum in a row
10     for(int i = 0; i < N; i++){
11         if(min > arr[id*N+i]){
12             min = arr[id*N + i]; // Storing the minimum for checking in
13             // next iteration
14             i1 = i; // Storing the coloumn where minimum is
15             found;
16             //printf("threadIdx: %d, min: %d, i1: %d\n", id, min, i1);
17         }
18     }
19
20     __syncthreads(); // Not required
21
22     // Finding Maximum in the coloumn corresponding to found minimum
23     (in i1)
24     for(int i = 0; i < N; i++){
25         if(max < arr[i*N + i1]){
26             max = arr[i*N + i1]; // Storing the Maximum
27             i2 = i; // Storing the row where Maximum is
28             found
29             //printf("threadIdx: %d, max: %d, i2: %d\n", id, max, i2);
30         }
31     }
32
33     // Check if the row matches with corresponding thread id (which
34     // determines the saddle condition)
35     if (i2 == id){
36         printf("Saddle Found: %d, at (%d,%d)\n", arr[id*N + i1], id,
37             i1 );
38     }
39 }
```

Listing 3: Saddle Finder Kernel

Testing the code

The code is tested with following matrix `main.cu`

```

1 int main(){
2     // specify the dimensions of the input matrix
3     const int M = 3; // square matrix dimensionss
4     unsigned numbytes = M * M * sizeof(int);
5
6     int arr[M][M] = {{1,2,4},{-1,-2,2},{4,5,6}}; // 3x3 matrix;
7         Saddle at (2,0) == 4
8
9     int *out;
10    cudaMalloc(&out, numbytes);
11    cudaMemcpy(out, arr, numbytes, cudaMemcpyHostToDevice);
12
13    saddleFinder<<<1,M>>>(out, M);
14    cudaDeviceSynchronize();
15    return 0;
16 }

```

Listing 4: Saddle Finder Main