Tam Phi
Hanh Do Phung

# Assignment 5: Jacobi parallelization with CUDA

### I. Introduction

Jacobi solver is a method of solving a system of n-linear equations of the form Ax = b in which x is a (n x 1) vector of all x_1 to x_n solutions and A is an n x n sized matrix representing the coefficients of variables x_1 to x_n.

### II. Design and Implementation

There are 2 steps in implementation.

### 1) Exec() function:

-Allocate memory on GPU for the output matrix U, call the kernel function to perform parallelization

Pseudo code{

  /*Allocate memory on device */

  /* Copy matrix to device */

  /* Initialize mutex lock */

}

Naïve version

```
129    while (! done){
130        jacobi_iteration_kernel_naive
131            <<< grid,thread_block >>>
132            (device_A.elements, device_B.elements, curr_X.elements, prev_X.elements, input_elements, ssd_device);
133
134        cudaDeviceSynchronize();
135        cudaMemcpy(ssd_host, ssd_device, elements * sizeof(double), cudaMemcpyDeviceToHost);
136
137        double ssd = 0;
138        int i = 0;
139        for (i = 0; i < elements; i++)
140            ssd += ssd_host[i];
141        // ping pong
142        float *X_temp = prev_X.elements;
143        prev_X.elements = curr_X.elements;
144        curr_X.elements = X_temp;
145
146        num_iter++;
147
148        double mse = sqrt (ssd);
149        if (mse <= THRESHOLD){
150            done = 1;
151        }
152    }
153
```

Optimized version

Tam Phi
Hanh Do Phung

```
206    while(!done)
207    {
208        (ping_pong)
209            ?
210            (jacobi_iteration_kernel_optimized<<< grid, threads >>>(device_A, opt, device_B, out_X, ssd_device, mutex_on_device))
211            : \
212            (jacobi_iteration_kernel_optimized<<< grid, threads >>>(device_A, out_X, device_B, opt, ssd_device, mutex_on_device));
213
214        cudaDeviceSynchronize();
215
216        num_iter++;
217
218        ping_pong = !ping_pong;
219        cudaMemcpy(&ssd, ssd_device, sizeof(double), cudaMemcpyDeviceToHost);
220        mse = sqrt (ssd);
221        if (mse <= THRESHOLD)
222            done = 1;
223    }
```

## 2) GPU kernel functions:

- **naive_kernel():** A function that performs division of current row by pivot element

```
14 __global__ void jacobi_iteration_kernel_naive(float *d_A, float * d_B, float *curr_X, float *prev_X, int *input_elements, double *ssd)
15 {
16
17    int tid = threadIdx.y;
18    double sum = 0;
19    int num_elements = *input_elements;
20
21    int i = (blockIdx.y * THREAD_BLOCK_SIZE) + tid;
22    for(int j = 0; j < num_elements; j++){
23        if (i != j)
24            sum += d_A[i * num_elements + j] * curr_X[j];
25    }
26
27    prev_X[i] = (d_B[i] - sum)/d_A[i * num_elements + i];
28
29    ssd[i] = (curr_X[i] - prev_X[i]) * (curr_X[i] - prev_X[i]);
30
31    return;
32 }
```

- **opt_kernel():** A function that performs reduction of under rows

```
34 __global__ void jacobi_iteration_kernel_optimized(float *d_A, float *opt, float *d_B, float *out_X, double *ssd, int *mutex)
35 {
36    __shared__ double ssd_device[THREAD_BLOCK_SIZE];
37
38    unsigned int tid = threadIdx.x;
39    unsigned int blkid = blockIdx.x * blockDim.x + threadIdx.x;
40
41    //reset ssd = 0
42    if (blkid == 0)
43        *ssd = 0.0;
44
45    // begin iteration
46    double sum = -d_A[blkid * MATRIX_SIZE + blkid] * opt[blkid];
47    for (int j = 0; j < MATRIX_SIZE; j++) {
48        sum += d_A[blkid + MATRIX_SIZE * j] * opt[j];
49    }
50    out_X[blkid] = (d_B[blkid] - sum)/d_A[blkid * MATRIX_SIZE + blkid];
51
52    if (blkid < MATRIX_SIZE)
53        ssd_device[tid] = (out_X[blkid] - opt[blkid]) * (out_X[blkid] - opt[blkid]);
54    else
55        ssd_device[tid] = 0.0;
56
57    __syncthreads();
58
59    for (unsigned int stride = blockDim.x >> 1; stride > 0; stride = stride >> 1) {
60        if(tid < stride)
61            ssd_device[tid] += ssd_device[tid + stride];
62        __syncthreads();
63    }
64
65    if (tid == 0) {
66        lock(mutex);
67        *ssd += ssd_device[0];
68        unlock(mutex);
69    }
70
71    return;
72 }
```

Tam Phi
Hanh Do Phung

The optimized version calculate single value of ssd on GPU while the naïve version copy the ssd_array back to CPU, then calculate the ssd on host side.

**Result and discussion**

All accuracy test was passed for a variety of matrix size as suggested in the figure below.

```
Generating 1024 x 1024 system

Performing Jacobi iteration on the CPU

Convergence achieved after 16797 iterations
Execution time = 28.953085s
Average diff between LHS and RHS: 0.001917

Performing naive Jacobi iteration on device
Execution time = 10.669594s
Done in 16797 iterations
Average diff between LHS and RHS: 0.001917

Performing optimized Jacobi iteration on device
Execution time = 9.194625s
Done in 16797 iterations
Average diff between LHS and RHS: 0.001917
```

Execution time

| Matrix width | Num threads | gold_time | naïve_time | opt_time |
| --- | --- | --- | --- | --- |
| 512 | 128 | 3.656533 | 4.199220 | 4.108942 |
| 512 | 256 | 3.673595 | 4.445559 | 4.141872 |
| 1024 | 128 | 28.30203 | 9.337974 | 8.826369 |
| 1024 | 256 | 28.953085 | 10.669594 | 9.194625 |
| 2048 | 128 | 258.089325 | 24.00955 | 20.450293 |
| 2048 | 256 | 236.729752 | 26.709866 | 20.975603 |

`

**Conclusion**

As the data moving step between the CPU and GPU is not included in the execution time measurements, the trend of more speed up is observed as the matrix size increases. This is because of the one-time overhead of cudaMalloc() and cudaMemcopy() executed is more significant overall compared to the execution time of the algorithm itself on the CPU/GPU. However, size of the thread block does not affect speed up for matrix of size 2048x2048 or smaller.