# Assignment 6- Hybrid MPI + MP + CUDA Programming
## CSCI 596: Scientific Computing & Visualization

Anup V Kanale

November 1, 2017

The goal of this assignment is understand the basics of cuda through two examples– for calculating the pair distribution function, and to calculate $\pi$.

## 1 Task I– Pair-Distribution Computation with CUDA

The pair distribution function (PDF) is the probability of finding any two atoms, which are a given distance $r$ apart, in a given volume. It is popularly used to measure the order, especially crystalline materials.

To compute the pdf, we fist calculate the histogram of atomic pair distances, i.e., calculate the distance between all pairs of atoms $i$ and $j$, say $r_{ij}$, put them in the appropriate '*bins*' $\Delta r$ to generate a histogram. Note that the maximum distance between pairs under periodic boundary condition is the diagonal of half the simulation box. The algorithm is as follows:

```
for all histogram bins i
    nhist[i] = 0
for all atomic pairs (i,j)
    ++nhist[r_ij/r]
```

Once this histogram is obtained, the pdf may then be computed using the following formula:

$$g(r_k) = \frac{\text{nhist}(k)}{2\pi r_k^2 \Delta r \rho N} \tag{1}$$

where $r_k$ is the distance between two atoms, $\rho$ is the number density of atoms, $N$ is the total number of atoms. Notice that the obtaining the histogram involves computing distance between all pairs of atoms, which can get expensive with a large number of atoms. The parallel processing power of the GPUs was used to carry out these computations.
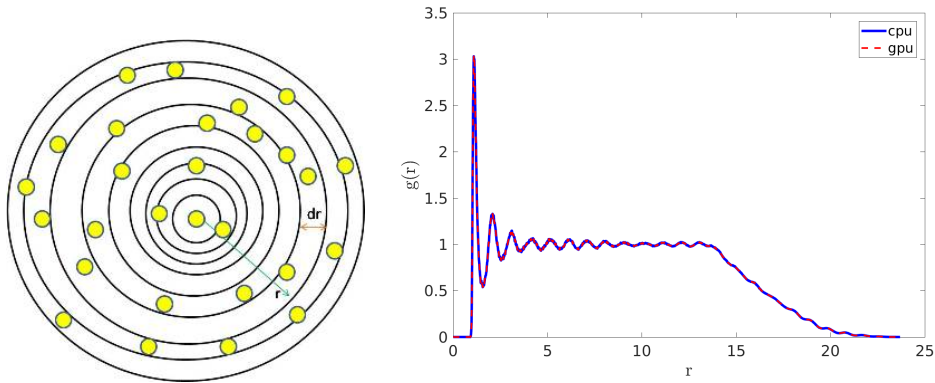


Figure 1: (a)Schematic for obtaining the histogram (b)PDF computation– CPU vs GPU

The code is attached in Appendix A. The pair distribution function computed using both the CPU and the GPUs shown for comparison in Fig 1 (the trailing off after $r = 14$ is explained in the notes.). The cpu implentation ran in 2.62 seconds whereas the gpu-parallel implementation ran in 0.28 seconds (almost 10 times faster!).

# 2 Task II– Parallel Computation of $\pi$ using MPI+OMP+CUDA

In this part, a triple-decker MPI+OpenMP+CUDA program was written to compute the value of $\pi$, by modifying the double-decker MPI+CUDA program provided in class, described in the lecture note on "Hybrid MPI+OpenMP+CUDA Programming". The formula used is

$$\pi = \int_0^1 \frac{4}{1+x^2} dx \approx \Delta \sum_{i=0}^{N-1} \frac{4}{1+x_i^2} \tag{2}$$
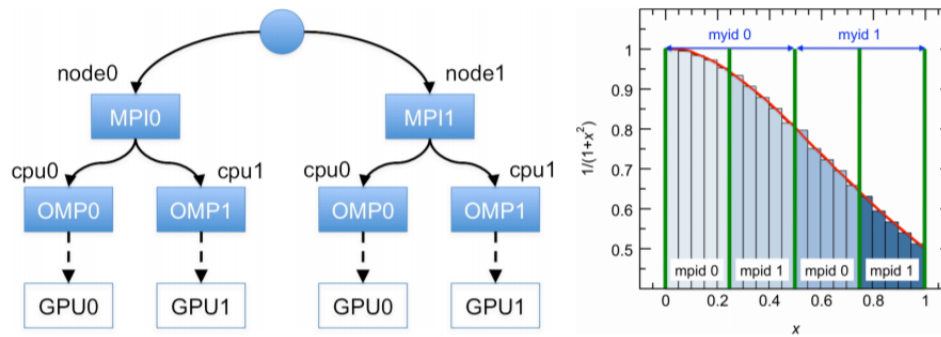


Figure 2: Output of program for computation of $\pi$

This program runs on two CPU cores, and 2 GPU devices on each compute node. This is achieved by launching one MPI rank per node, where each rank spawns two OpenMP threads, and each thread uses a different GPU device. This is illustrated in Fig 2. The above OpenMP multithreading introduces a race condition for variable pi. This can be circumvented by data privatization using `reduction(+:pi)`, which creates a `pi` variable for each thread. The '+:' means that the 'addition reduction' is to be performed on this variable.

The program can be found in appendix B, and the output of the program is as shown below in Fig 2. `myid` refers to the process rank, `mpid` is the thread number within the processor, and `device used` refers to the gpu.

```
[kanale@hpc3026 assgn6]$ mpirun -np 2 -machinefile nodefile ./pi --bind-to none
[hpc3026:44118] mca: base: component_find: unable to open /usr/usc/openmpi/1.6.4/lib/openmpi/mca_btl_mx: libmyriexpress.
so: cannot open shared object file: No such file or directory (ignored)
[hpc3027:07438] mca: base: component_find: unable to open /usr/usc/openmpi/1.6.4/lib/openmpi/mca_btl_mx: libmyriexpress.
so: cannot open shared object file: No such file or directory (ignored)
myid = 0; mpid = 0: device used = 0; partial pi = 0.979926
myid = 0; mpid = 1: device used = 1; partial pi = 0.874671
myid = 1; mpid = 1: device used = 1; partial pi = 0.567582
myid = 1; mpid = 0: device used = 0; partial pi = 0.719409
PI = 3.141588
```

Figure 3: Output of program for computation of $\pi$

### Important points

1. Finding compatible version of mpi, openMP and cuda is not trivial. And it varies from system to system.

2. pay attention to how the sequential thread index is calculated across blocks.

# A    Appendix– Calculating the PDF

Note that there were no major changes in the `main()` function, so it is not shown.

```c
 1 /*————————————————————————————
 2 Program pdf0.c computes a pair distribution function for n atoms
 3 given the 3D coordinates of the atoms.
 4 ————————————————————————————*/
 5 #include <stdio.h>
 6 #include <math.h>
 7 #include <time.h>
 8 #include <stdlib.h>
 9
10 #define NHBIN 2000   // Histogram size
11
12 float al[3];          // Simulation box lengths
13 int n;                // Number of atoms
14 float *r;             // Atomic position array
15 FILE *fp;
16
17 __constant__ float DALTH[3];
18 __constant__ int DN;
19 __constant__ float DDRH;
20
21 __device__ float d_SignR(float v, float x) {if (x > 0) return v; else
     return -v;}
22
23
24 __global__ void gpu_histogram_kernel(float *r, float *nhis) {
25   int i,j,a,ih;
26   float rij,dr;
27
28   int iBlockBegin = (DN/gridDim.x)*blockIdx.x;
29   int iBlockEnd = min((DN/gridDim.x)*(blockIdx.x+1),DN);
30   int jBlockBegin = (DN/gridDim.y)*blockIdx.y;
31   int jBlockEnd = min((DN/gridDim.y)*(blockIdx.y+1),DN);
32   for (i=iBlockBegin+threadIdx.x; i<iBlockEnd; i+=blockDim.x) {
33     for (j=jBlockBegin+threadIdx.y; j<jBlockEnd; j+=blockDim.y) {
34       if (i<j) {
35       // Process (i,j) atom pair
36       rij = 0.0;
37           for (a=0; a<3; a++) {
38             dr = r[3*i+a]-r[3*j+a];
39             /* Periodic boundary condition */
40             dr =
                 dr-d_SignR(DALTH[a],dr-DALTH[a])-d_SignR(DALTH[a],dr+DALTH[a]);
41             rij += dr*dr;
42           }
```

```
43              rij = sqrt(rij); /* Pair distance */
44              ih = rij/DDRH;
45              // nhis[ih] += 1.0; /* Entry to the histogram */
46              atomicAdd(&nhis[ih],1.0);
47         } // end if i<j
48       } // end for j
49    } // end for i
50 }
51
52
53 /*————————————————————————————————————————————————————*/
54 void histogram() {
55 /*————————————————————————————————————————————————————
56 Constructs a histogram NHIS for atomic−pair distribution.
57 ————————————————————————————————————————————————————*/
58    float alth[3];
59    float* nhis;   // Histogram array
60    float rhmax,drh,density,gr;
61    int a,ih;
62
63    float* dev_r; // Atomic positions
64    float* dev_nhis; // Histogram
65
66    /* Half the simulation box size */
67    for (a=0; a<3; a++) alth[a] = 0.5*al[a];
68    /* Max. pair distance RHMAX & histogram bin size DRH */
69    rhmax = sqrt(alth[0]*alth[0]+alth[1]*alth[1]+alth[2]*alth[2]);
70    drh = rhmax/NHBIN;   // Histogram bin size
71
72
73    nhis = (float*)malloc(sizeof(float)*NHBIN);
74    // for (ih=0; ih<NHBIN; ih++) nhis[ih] = 0.0; // Reset the histogram
75
76    cudaMalloc((void**)&dev_r,sizeof(float)*3*n);
77    cudaMalloc((void**)&dev_nhis,sizeof(float)*NHBIN);
78    cudaMemcpy(dev_r,r,3*n*sizeof(float),cudaMemcpyHostToDevice);
79    cudaMemset(dev_nhis,0.0,NHBIN*sizeof(float));
80    cudaMemcpyToSymbol(DALTH,alth,sizeof(float)*3,0,cudaMemcpyHostToDevice);
81    cudaMemcpyToSymbol(DN,&n,sizeof(int),0,cudaMemcpyHostToDevice);
82    cudaMemcpyToSymbol(DDRH,&drh,sizeof(float),0,cudaMemcpyHostToDevice);
83
84    dim3 numBlocks(8,8,1);
85    dim3 threads_per_block(16,16,1);
86    gpu_histogram_kernel<<<numBlocks,threads_per_block>>>(dev_r,dev_nhis);
87
88
89    // Compute dev_nhis on GPU: dev_r[] Âő dev_nhis[]
```

```
90     cudaMemcpy( nhis , dev_nhis ,NHBIN*sizeof(float) ,cudaMemcpyDeviceToHost) ;
91
92     density = n/( al[0]*al[1]*al[2]) ;
93     /* Print out the histogram */
94     fp = fopen("pdf_gpu.d" ,"w") ;
95     for (ih=0; ih<NHBIN; ih++) {
96       gr = nhis[ih]/(2*M_PI*pow((ih+0.5)*drh ,2)*drh*density*n) ;
97       fprintf(fp ,"%e %e\n" ,(ih+0.5)*drh , gr) ;
98     }
99     fclose(fp) ;
100    free(nhis) ;
101 }
```

<div align="center">pdf1writeup.cu</div>

# B   Appendix– Computation of $\pi$

```
 1 // Hybrid MPI+CUDA computation of Pi
 2 #include <stdio.h>
 3 #include <mpi.h>
 4 #include <cuda.h>
 5 #include <omp.h>
 6
 7 #define NUM_DEVICE 2  // # of GPU devices = # of OpenMP threads
 8 #define NBIN   10000000   // Number of bins
 9 #define NUM_BLOCK   13   // Number of thread blocks
10 #define NUM_THREAD 192   // Number of threads per block
11
12 // Kernel that executes on the CUDA device
13 __global__ void cal_pi(float *sum,int nbin,float step ,float offset ,int
       nthreads ,int nblocks) {
14   int i;
15   float x;
16   int idx = blockIdx.x*blockDim.x+threadIdx.x;  // Sequential thread
         index across the blocks
17   for (i=idx; i<nbin; i+=nthreads*nblocks) {  // Interleaved bin
         assignment to threads
18     x = offset+(i+0.5)*step;
19     sum[idx] += 4.0/(1.0+x*x);
20   }
21 }
22
23 int main(int argc ,char **argv) {
24   int myid ,nproc ,nbin ,tid ,dev_used ,mpid;
25   float step ,offset ,pi=0.0,pig;
26   dim3 dimGrid(NUM_BLOCK,1 ,1) ;  // Grid dimensions (only use 1D)
27   dim3 dimBlock(NUM_THREAD,1 ,1) ;  // Block dimensions (only use 1D)
```

<div align="center">5</div>

```
28    float *sumHost,*sumDev;   // Pointers to host & device arrays
29
30    MPI_Init(&argc,&argv);
31    MPI_Comm_rank(MPI_COMM_WORLD,&myid);   // My MPI rank
32    MPI_Comm_size(MPI_COMM_WORLD,&nproc);   // Number of MPI processes
33    omp_set_num_threads(NUM_DEVICE);
34    nbin = NBIN/(nproc*NUM_DEVICE);   // Number of bins per MPI process
35    step = 1.0/(float)(nbin*nproc*NUM_DEVICE);   // Step size with
          redefined number of bins
36
37    #pragma omp parallel private(mpid,offset,sumHost,tid) reduction(+:pi)
38    {
39        mpid = omp_get_thread_num();
40        offset = (NUM_DEVICE*myid+mpid)*step*nbin;   // Quadrature-point
              offset
41
42        cudaSetDevice(mpid%2);
43        size_t size = NUM_BLOCK*NUM_THREAD*sizeof(float);   //Array
              memory size
44        sumHost = (float *)malloc(size);   // Allocate array on host
45        cudaMalloc((void **) &sumDev,size);   // Allocate array on device
46        cudaMemset(sumDev,0,size);   // Reset array in device to 0
47        // Calculate on device (call CUDA kernel)
48        cal_pi <<<dimGrid,dimBlock>>>
              (sumDev,nbin,step,offset,NUM_THREAD,NUM_BLOCK);
49        // Retrieve result from device and store it in host array
50        cudaMemcpy(sumHost,sumDev,size,cudaMemcpyDeviceToHost);
51        // Reduction over CUDA threads
52        for(tid=0; tid<NUM_THREAD*NUM_BLOCK; tid++)
53          pi += sumHost[tid];
54        pi *= step;
55        // CUDA cleanup
56        free(sumHost);
57        cudaFree(sumDev);
58
59        // Reduction over MPI processes
60        cudaGetDevice(&dev_used);
61        printf("myid = %d; mpid = %d: device used = %d; partial pi =
              %f\n", myid,mpid,dev_used,pi);
62    } // End omp parallel
63    MPI_Allreduce(&pi,&pig,1,MPI_FLOAT,MPI_SUM,MPI_COMM_WORLD);
64    if (myid==0) printf("PI = %f\n",pig);
65    MPI_Finalize();
66    return 0;
67 }
```

piCudaMpiOmpwriteup.cu