

# 1 Hands-on; GPU parallelization

```
$ lspci | grep -i nvidia
01:00.0 VGA compatible controller: NVIDIA Corporation GP108
[GeForce GT 1030] (rev a1)
01:00.1 Audio device: NVIDIA Corporation GP108 High Definition
Audio Controller (rev a1)
```

Installation Instructions:

```
$ wget https://developer.download.nvidia.com/compute/cuda/repos/
ubuntu2204/x86_64/cuda-keyring_1.0-1_all.deb
$ sudo dpkg -i cuda-keyring_1.0-1_all.deb
$ sudo apt-get update
$ sudo apt-get -y install cuda
$ sudo reboot
$ export PATH=/usr/local/cuda-12.0/bin${PATH:+:${PATH}}
$ export LD_LIBRARY_PATH=/usr/local/cuda-12.0/lib64${LD_LIBRARY_PATH:
+:${LD_LIBRARY_PATH}}
```

1. **Hello world I:** [code24.cu](#) Follow the steps below for executing Cuda code;

```
nvcc -o code24 code24.cu
./code24
```

```
1 #include <stdio.h>
2 #include <unistd.h>
3 __device__ const char *STR = "HELLO WORLD!";
4 const int STRLENGTH = 12;
5 __global__ void cuda_hello(){
6     // blockIdx.x: Block index within the grid in the x direction
7     // threadIdx.x: Thread index within the block
8     // blockDim.x,y,z # of threads in a block
9     printf("Hello World from GPU! (%d,%d) : %c ThreadID %d \n",
10         blockIdx.x, threadIdx.x, STR[threadIdx.x % STRLENGTH], (threadIdx
11         .x +blockIdx.x*blockDim.x));
12 }
13 /*
14 ./deviceQuery Starting...
15
16 CUDA Device Query (Runtime API) version (CUDA static linking)
17
18 Detected 3 CUDA Capable device(s)
19
20 Device 2: "NVIDIA GeForce GT 1030"
21 CUDA Driver Version / Runtime Version          12.0 / 11.8
22 CUDA Capability Major/Minor version number:    6.1
23 Total amount of global memory:                 1998 MBytes
24         (2095185920 bytes)
25 (003) Multiprocessors, (128) CUDA Cores/MP:    384 CUDA Cores
```

```

23 GPU Max Clock rate: 1468 MHz (1.47 GHz)
24 Memory Clock rate: 3004 Mhz
25 Memory Bus Width: 64-bit
26 L2 Cache Size: 524288 bytes
27 Maximum Texture Dimension Size (x,y,z) 1D=(131072), 2D
   =(131072, 65536), 3D=(16384, 16384, 16384)
28 Maximum Layered 1D Texture Size, (num) layers 1D=(32768), 2048
   layers
29 Maximum Layered 2D Texture Size, (num) layers 2D=(32768, 32768),
   2048 layers
30 Total amount of constant memory: 65536 bytes
31 Total amount of shared memory per block: 49152 bytes
32 Total shared memory per multiprocessor: 98304 bytes
33 Total number of registers available per block: 65536
34 Warp size: 32
35 Maximum number of threads per multiprocessor: 2048
36 Maximum number of threads per block: 1024
37 Max dimension size of a thread block (x,y,z): (1024, 1024, 64)
38 Max dimension size of a grid size (x,y,z): (2147483647, 65535,
   65535)
39 Maximum memory pitch: 2147483647 bytes
40 Texture alignment: 512 bytes
41 Concurrent copy and kernel execution: Yes with 2 copy
   engine(s)
42 Run time limit on kernels: Yes
43 Integrated GPU sharing Host Memory: No
44 Support host page-locked memory mapping: Yes
45 Alignment requirement for Surfaces: Yes
46 Device has ECC support: Disabled
47 Device supports Unified Addressing (UVA): Yes
48 Device supports Managed Memory: Yes
49 Device supports Compute Preemption: Yes
50 Supports Cooperative Kernel Launch: Yes
51 Supports MultiDevice Co-op Kernel Launch: Yes
52 Device PCI Domain ID / Bus ID / location ID: 0 / 4 / 0
53 Compute Mode:
54 < Default (multiple host threads can use ::cudaSetDevice() with
   device simultaneously) >
55 deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 12.0, CUDA
   Runtime Version = 11.8, NumDevs = 3
56 Result = PASS
57 */
58 int main() {
59     printf("Hello World from CPU!\n");
60     sleep(2);
61     int threads_per_block=12;
62     int blocks_per_grid=2;
63     cuda_hello <<< blocks_per_grid, threads_per_block >>> ();
64     cudaDeviceSynchronize(); /* Halt host thread execution on CPU until
   the device has finished processing all previously requested tasks
   */
65     return 0;
66 }

```

2. **Hello world II:** [code25.cu](http://code25.cu) Follow the steps below for executing Cuda code;

```
nvcc -o code25 code25.cu
./code25
```

```
1 #include <stdio.h>
2 #include <cuda.h>
3 #include <cuda_runtime.h>
4
5 #define N 720 // number of computations
6 #define GRID_D1 20 // constants for grid and block sizes
7 #define GRID_D2 3 // constants for grid and block sizes
8 #define BLOCK_D1 12 // constants for grid and block sizes
9 #define BLOCK_D2 1 // constants for grid and block sizes
10 #define BLOCK_D3 1 // constants for grid and block sizes
11
12 --global-- void hello(void) // this is the kernel function called for
    each thread
13 {
14     // we use the CUDA variables {threadIdx, blockIdx, blockDim, gridDim
    } to determine a unique ID for each thread
15     int myblock = blockIdx.x + blockIdx.y * blockDim.x; // id of the
    block
16     int blocksize = blockDim.x * blockDim.y * blockDim.z; // size of
    each block (within grid of blocks)
17     int subthread = threadIdx.z*(blockDim.x * blockDim.y) + threadIdx.y*
    blockDim.x + threadIdx.x; // id of thread in a given block
18     int idx = myblock * blocksize + subthread; // assign overall id/
    index of the thread
19     int nthreads=blocksize*gridDim.x*gridDim.y; // Total # of threads
20     int chunk=20; // Vary this value to see the changes at the output
21     if(idx < chunk || idx > nthreads-chunk) { // print buffer from
    within the kernel is limited so only print for first and last
    chunks of threads
22         if (idx < N){
23             printf("Hello world! My block index is (%d,%d) [Grid dims=(%d,%d)
    ], 3D-thread index within block=(%d,%d,%d) => thread index=%d \n"
    , blockIdx.x, blockIdx.y, blockDim.x, blockDim.y, threadIdx.x,
    threadIdx.y, threadIdx.z, idx);
24         }
25         else
26         {
27             printf("Hello world! My block index is (%d,%d) [Grid dims=(%d,%d)
    ], 3D-thread index within block=(%d,%d,%d) => thread index=%d [###
    this thread would not be used for N=%d ###]\n", blockIdx.x,
    blockIdx.y, blockDim.x, blockDim.y, threadIdx.x, threadIdx.y,
    threadIdx.z, idx, N);
28         }
29     }
30 }
31
32 int main(int argc, char **argv)
33 {
34     // objects containing the block and grid info
35     const dim3 blockSize(BLOCK_D1, BLOCK_D2, BLOCK_D3);
36     const dim3 gridSize(GRID_D1, GRID_D2, 1);
37     int nthreads = BLOCK_D1*BLOCK_D2*BLOCK_D3*GRID_D1*GRID_D2; // Total
    # of threads
```

```

38  if (nthreads < N){
39      printf("\n===== NOT ENOUGH THREADS TO COVER N=%d\n",N);
40  }
41  else
42      {
43      printf("Launching %d threads (N=%d)\n",nthreads ,N);
44      }
45  hello<<<gridSize, blockSize>>>(); // launch the kernel on the
    specified grid of thread blocks
46  // Need to flush prints, otherwise none of the prints from within
    the kernel will show up
47  // as program exit does not flush the print buffer.
48  cudaError_t cudaerr = cudaDeviceSynchronize();
49  if (cudaerr){
50      printf("kernel launch failed with error \"%s\".\n",
51          cudaGetErrorString(cudaerr));
52  }
53  else
54      {
55      printf("kernel launch success!\n");
56      }
57  printf("That's all!\n");
58  return 0;
59  }

```

3. **Vector Addition:** [code26.cu](#) Follow the steps below for executing Cuda code;

```
nvcc -o code26 code26.cu
./code26
```

```
1 // https://www.olcf.ornl.gov/tutorials/cuda-vector-addition/
2 // https://github.com/olcf-tutorials/vector_addition_cuda
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <math.h>
6
7 #define n 100000 // Size of array
8
9 /* CUDA KERNEL Compute the sum of two vectors
10 * Each thread takes care of one element of C
11 * C[i] = A[i] + B[i]
12 */
13 __global__ void add_vectors(double *a, double *b, double *c)
14 {
15     int id = blockIdx.x*blockDim.x+threadIdx.x; // Get our global thread
16     ID
17     if (id < n) // Make sure we do not go out of bounds
18         c[id] = a[id] + b[id]; /* Compute the element of C */
19 }
20 int main( int argc , char* argv [] )
21 {
22     double *h_a; // Host input vectors
23     double *h_b;
24     double *h_c;
25
26     double *d_a; // Device input vectors
27     double *d_b;
28     double *d_c;
29
30     size_t bytes = n*sizeof(double); // Size, in bytes, of each vector
31
32     h_a = (double*)malloc(bytes); // Allocate memory for each vector on
33     host
34     h_b = (double*)malloc(bytes);
35     h_c = (double*)malloc(bytes);
36
37     cudaMalloc(&d_a, bytes); // Allocate memory for each vector on GPU
38     cudaMalloc(&d_b, bytes);
39     cudaMalloc(&d_c, bytes);
40
41     int i;
42     for( i = 0; i < n; i++ ) { // Initialize vectors on host
43         h_a[i] = sin(i)*sin(i);
44         h_b[i] = cos(i)*cos(i);
45     }
46
47     // Copy data from host arrays h_a and h_b to device arrays d_a and
48     d_b
49     cudaMemcpy( d_a, h_a, bytes, cudaMemcpyHostToDevice);
50     cudaMemcpy( d_b, h_b, bytes, cudaMemcpyHostToDevice);
51
52     int thr_per_blk = 1024; // blockSize. Number of threads in each
53     thread block
```

```

51 int blk_in_grid = ceil( (float) n/thr_per_blk ); // gridSize. Number
    of thread blocks in grid
52
53 add_vectors<<<blk_in_grid, thr_per_blk>>>(d_a, d_b, d_c); // Execute
    the kernel
54
55 cudaMemcpy( h_c, d_c, bytes, cudaMemcpyDeviceToHost ); // Copy data
    from device array d_c to host array h_c
56
57 // Sum up vector c and print result divided by n, this should equal
    1 within error
58 double sum = 0;
59 for(i=0; i<n; i++)
60     sum += h_c[i];
61 printf("final result: %f\n", sum/n);
62
63 cudaFree(d_a); // Free device memory
64 cudaFree(d_b);
65 cudaFree(d_c);
66
67 free(h_a); // Free host memory
68 free(h_b);
69 free(h_c);
70
71 printf("-----\n");
72 printf(" _ENDED_ \n");
73 printf("-----\n");
74 printf("N           = %d\n", n);
75 printf("Threads Per Block = %d\n", thr_per_blk);
76 printf("Blocks In Grid   = %d\n", blk_in_grid);
77 printf("-----\n\n");
78
79 return 0;
80 }

```