

CUDA Threads

1. Complete and run the following `vecAdd4.cu` program, and see if it works for any $n > T$, for example, when $n = 1234$ and $T = 64$

```

1  #define n 1024 // size of vectors
2  #define T 240 // number of threads per block
3
4  global__ void vecAdd(int *A, int *B, int *C){
5      int i = blockIdx.x * blockDim.x + threadIdx.x;
6
7      if (i < n){ // allows for more threads than vector elements
8          C[i] = A[i] + B[i]; // some unused
9      }
10 }
11
12 int main(int argc, char *argv[]){
13     int blocks = (n + T - 1) / T; // efficient way of rounding to the next integer
14
15     ... // see the old code
16     vecAdd<<<blocks, T>>>(devA, devB, devC);
17     ...
18 }

```

Figure 1: `vecAdd4.cu`

2. Create a CUDA program, `vecFill.cu`, to fill in the array `A[256]` using 4 thread blocks. Each block has 64 threads. Each element `A[i] = i` as shown below.

1	2	3	4	5	6	...	253	254	255
---	---	---	---	---	---	-----	-----	-----	-----

3. Create and run the `matmul2.cu` program which multiplies two square matrices. Width is strictly a multiple of `TILE_WIDTH`.

```

1  #include <stdio.h>
2  #define Width 32 // size of Width x Width matrix
3  #define TILE_WIDTH 16
4
5  global__ void MatrixMulKernel(float *Md, float *Nd, float *Pd, int ncols){
6      int row = blockIdx.y * blockDim.y + threadIdx.y;
7      int col = blockIdx.x * blockDim.x + threadIdx.x;
8
9      // Pvalue is used to store the element of the output matrix
10     // that is computed by the thread
11
12     float Pvalue = 0;
13     for(int k = 0; k < ncols; k++){
14         float Melement = Md[row * ncols + k];
15         float Nelement = Nd[k * ncols + col];
16         Pvalue += Melement * Nelement;
17     }
18
19     Pd[row * ncols + col] = Pvalue;
20 }
21

```

Figure 2: `matmul2.cu` (part 1/2)

```
22 int main(int argc, char **argv){
23     int i,j;
24     int size = Width * Width * sizeof(float);
25     float M[Width][Width], N[Width][Width], P[Width][Width];
26     float *Md, *Nd, *Pd;
27
28     for(i = 0; i < Width; i++){
29         for(j = 0; j < Width; j++){
30             M[i][j] = 1; N[i][j] = 2;
31         }
32     }
33
34     cudaMalloc((void**)&Md, size);
35     cudaMalloc((void**)&Nd, size);
36     cudaMalloc((void**)&Pd, size);
37
38     cudaMemcpy( Md, M, size, cudaMemcpyHostToDevice);
39     cudaMemcpy( Nd, N, size, cudaMemcpyHostToDevice);
40
41     // setup the execution configuration
42     dim3 dimBlock(TILE_WIDTH, TILE_WIDTH);
43     dim3 dimGrid(Width/TILE_WIDTH, Width/TILE_WIDTH);
44
45     // launch the device computation thread!
46     MatrixMulKernel<<<dimGrid, dimBlock>>>(Md, Nd, Pd, Width);
47
48     // read P from the device
49     cudaMemcpy(P, Pd, size, cudaMemcpyDeviceToHost);
50
51     // free device matrices
52     cudaFree(Md); cudaFree(Nd); cudaFree(Pd);
53
54     for(i = 0; i < Width; i++){
55         for(j = 0; j < Width; j++){
56             printf("%.2f ", P[i][j]);
57         }
58         printf("\n");
59     }
60 }
```

Figure 3: mutmal2.cu (part 2/2)

4. Modify the program matmul2.cu into matmul3.cu to work with square matrices of arbitrary size. (Width not necessary a multiple of TILE_WIDTH)
5. Create a MS Word document (u5xxxxxx.docx), and put your source code from step 1. to step 4., along with the screenshots of their result.

- vecAdd4.cu
- vecFill1.cu
- matmul2.cu
- matmul3.cu