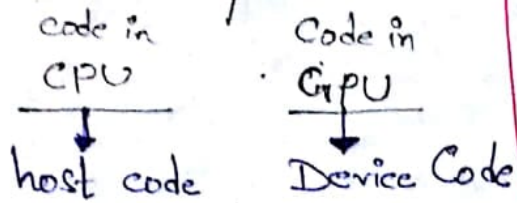


- 1) • Nvidia-smi (System Management Interface)
- cudaMallocManaged()
 - cudaDeviceSynchronize()
- cudaFree()



A.cu

↳ CUDA - accelerated Program

HelloCPU <<< 1, 1 >>> ();

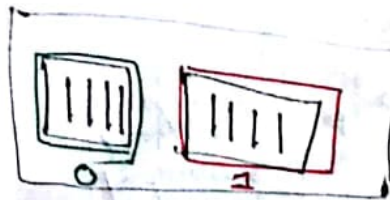
collection of thread is called block

1 block of thread

each block contains 1 thread

Memory allocated in this way is called UM (Unified Memory) which is automatically migrated between CPU & GPU

- 2) • nvcc = nvidia cuda Compiler.



→ A collection of block associated with a given kernel launch is a grid

• HelloGPU <<< 1, 1, >>> ();

↳ GPU functions are called kernel

3) grid Dim.x

↳ ②

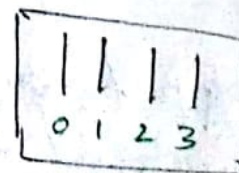
block Dim.x

↳ ④

block Id x.x

↳ 0

↳ 1



thread Id x.x

↳ 0

↳ 1, 2, 3

4) Each thread : has access to the size of its block via \Rightarrow blockDim.x
also \Rightarrow blockIdx.x

$$\text{threadIdx.x} + \text{blockIdx.x} * \text{blockDim.x}$$

max (0 - 1023)

0
1
2
3
4
5
6
7

These are CUDA variables

5) $\text{int } N = 2 \ll 20;$

~~size_t~~ size_t size = N * sizeof(int);

int *a;

cudaMallocManaged(&a, size); which can be used by both CPU & GPU.

cudaFree(a);

⑧ cudaError_t arrayError, asyncError;

\uparrow
= cudaGetLastError();

\downarrow
= cudaDeviceSynchronize();

if (arrayError != cudaSuccess)

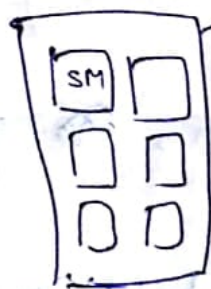
printf("Error: %s\n", cudaGetErrorString(arrayError));

• Managing Accelerated Application Memory with CUDA Unified Memory and nsys

CudaDe	— 1	— 2376 366692	8 31 36 245
cuda Ma	— 3	— 202862158	20 38 38 270
cuda Fr	— 3	— 18467746	18 06 07 21
cuda La	— 1	— 67525	48 899

cudaDeviceSynchronize	54 07 28 599
cudaMallocManaged	20 52 70 300
cudaFree	18 33 91 85
cudaLaunchKernel	50 67 9

• Streaming Multiprocessor (SM)



SM is functional unit of NVIDIA GPU

No of SM = No of Multiprocessor

• int deviceId;

• cudaGetDevice (&deviceId);

cudaDeviceProp prop;

~~Get~~

cudaGetDeviceProperties (&prop, deviceId);

• Grid Stride loop

int index = ~~thread~~

↳ Multi GPU program
↳ cuda Programming

int index = threadIdx.x + blockIdx.x *
blockDim.x;

int stride = blockDim.x * blockDim.x;

for (i = index; i < N; i = i + stride)

{

}

Run

» ! nvcc -o output 1st.cu -run

» ! nsys profile --stats=true ./output

↳ nsight system command line

■ There is a concept of fallow SM
or Wid SM

Matrix Multiply

Size = $N \times N \times \text{size of (int)}$;
 $\frac{64}{64} \frac{64}{64}$

```
for (int row = 0; row < N; row++)
```

```
{
    for (int col = 0; col < N; col++)
```

```
{
    a [row * N + col] = row;
```

```
    b [row * N + col] = col + 2;
```

```
    c [row * N + col] = 0;
```

```
}
```

```
}
```

```
int row = blockId.x * blockDim.x + threadIdx.x;
```

```
int col = blockId.y * blockDim.y + threadIdx.y;
```

```
if (row < N and col < N)
```

```
{
```

```
    for (int k = 0; k < N; k++)
```

```
{
```

```
        val += a [row * N + k] * b [k * N + col]
```

```
    }
```

```
    c [row * N + col] = val;
```

```
}
```

dim3 thread_per_block (16, 16, 1); // A 16*16 block threads.

dim3 number_of_blocks ((N / thread_per_block.x) + 1,

(N / thread_per_block.y) + 1, 1);

(~~4~~ 4+1, 4+1, 1)

matrix Mul GPU \lll no. of blocks, thread-per-block \ggg
 (a, b, c \ll gpu);

for (iStep = 0 : iStep <= nStep : iStep++)
 {

Kernel

2 cuda Function

Scrap

}

• %f = 93000.000.0 • ("%g", 3.14159)

%g = 9.3e+07 \rightarrow 3.14

• ("%g", 3.14159)

\rightarrow 3.141
3

• void foo() \rightarrow velocity

int i = threadIdx.x + blockDim.x

that stride = gridDim.x * blockDim.x

while (i < n)
 {

p[i].x += p[i].vx * dt

y

vy

z

vz

i += stride

}

• void foo2 (dt) \rightarrow body force

while (i < n) {

float Fx = 0.0f, Fy = 0.0f, Fz = 0.0f;

for (int j = 0; j < n; j++) {

float dx = P[j].x - P[i].x;

dy =

dz =

float distSqr = dx * dx + dy * dy + dz * dz +

SOFTENING;

float invDist = 1.0f / sqrtf (distSqr);

float invDist3 = invDist * invDist * invDist;

Fx += dx * invDist3;

Fy += dy * invDist3;

}

P[i].vx += dt * Fx;

i += stride;

}

}

cudaMemcpy

- Allocate memory on the device
- Copy data from host to device
- Perform some calculations
- Copy data from device to Host.
- Free Allocated device memory

cudaMemcpy (device Array, Host Array, bytes, cudaMemcpy Host To Device):

cudaMemcpy (Host Array, device Array, bytes, cudaMemcpy Device To Host):