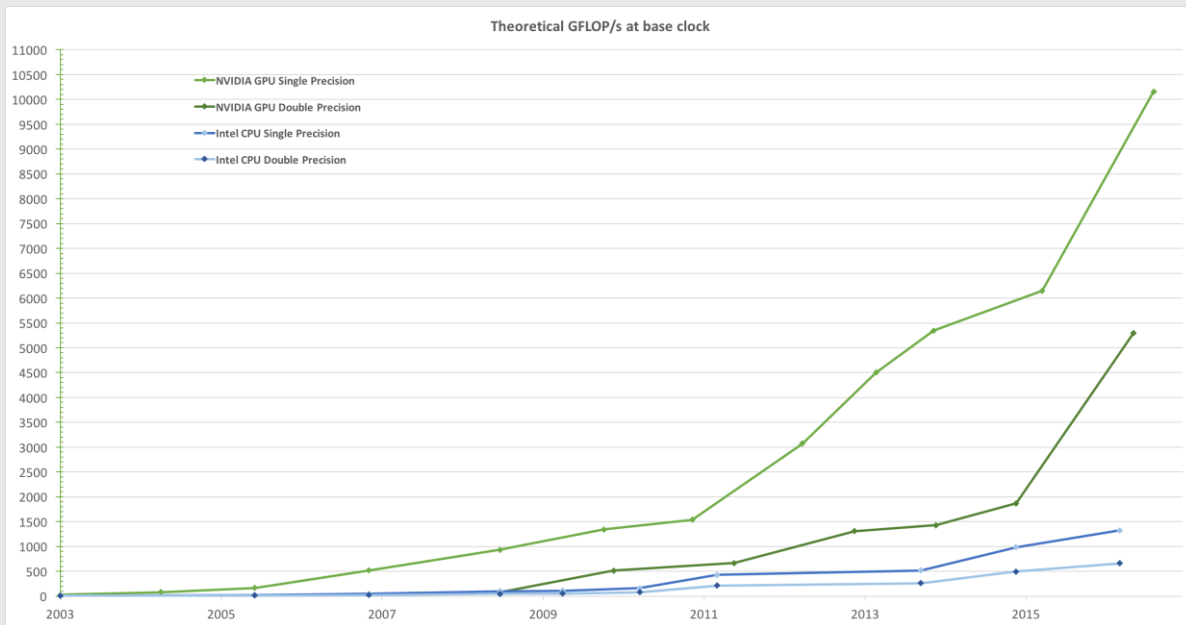


Cuda: NVIDIA GPU programming

Paulius Milmantas

Introduction - Purpose of GPU computing

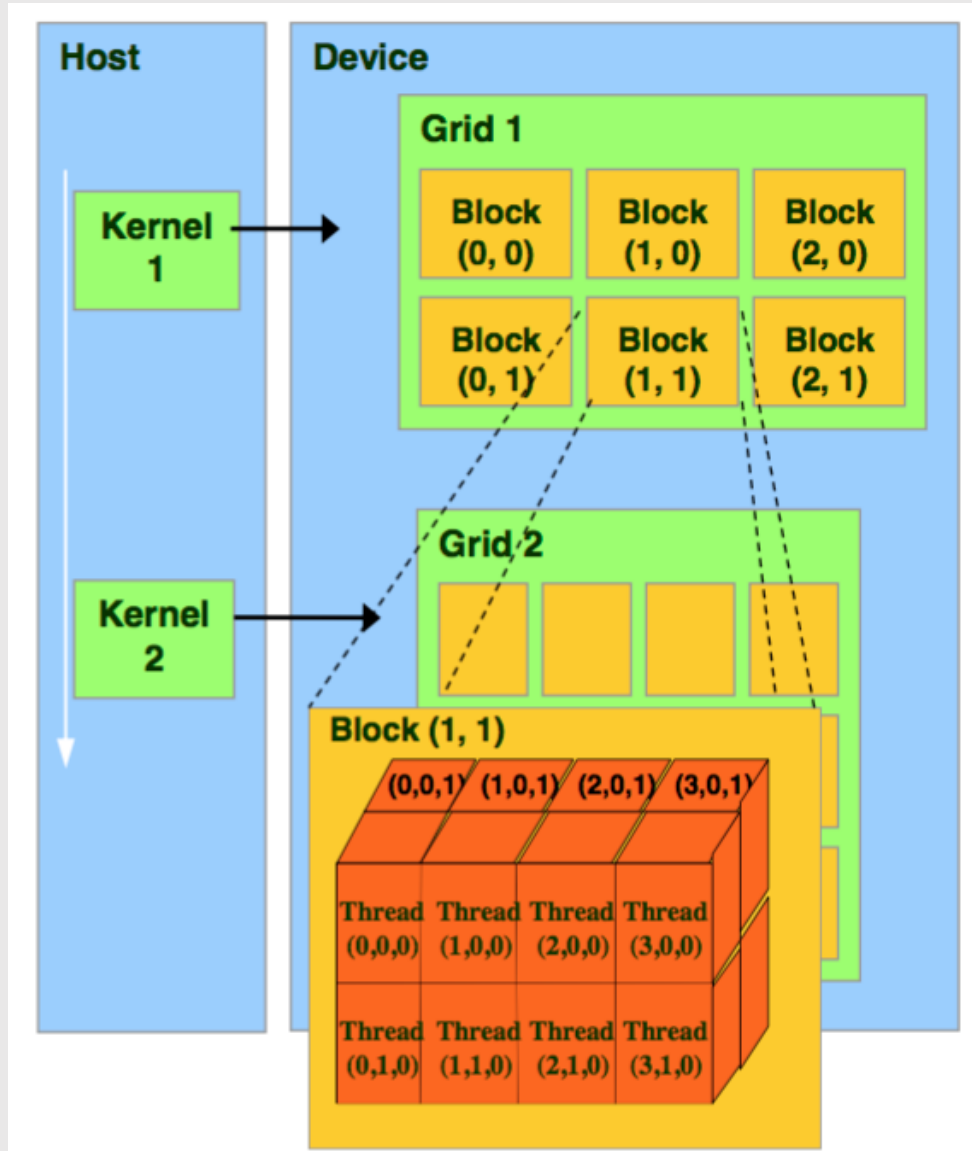
- **CPU** is designed to handle **complex tasks**
- **GPU** is designed to handle repetitive **low-level tasks**.



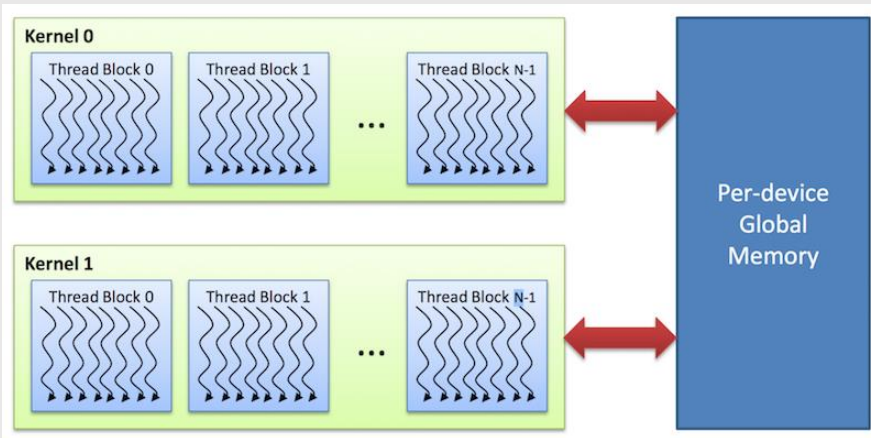
Memory structure

Kernel – Shared function executed by GPU
Kernel is executed many times

- **Thread** executes **kernel** code
- **Threads** are grouped into **threads blocks**
- **Blocks** are grouped into **grids**

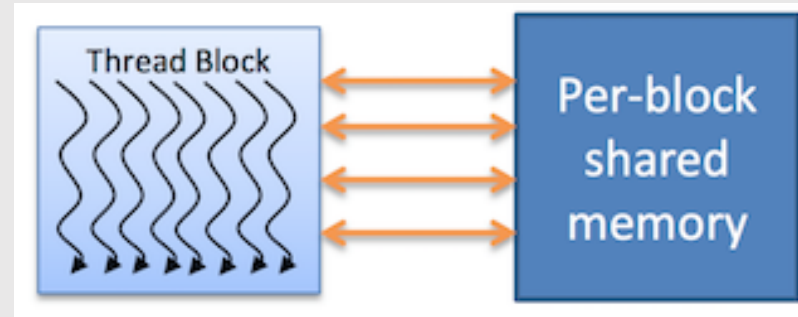


Global memory



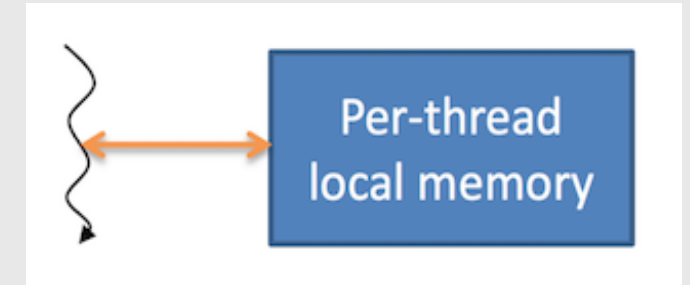
- Memory allocated by the host
- Use by all GPU processes

Shared memory



- Accessible only by the threads with a block
- Memory not shared across different blocks
- Much faster than local/global memory
- Only exists for the lifetime of the block

Local memory



- Only exists for the lifetime of the thread
- Handled automatically by the compiler

Typical execution process

1

Copy data from host memory to GPU memory

2

Execute program

3

Copy results from GPU memory to CPU memory.

Example #1 – Kernel

```
1  __global__ void Add(float *a, float *b, float *c)
2  {
3      int i = threadIdx.x;
4      c[i] = a[i] + b[i];
5  }
6
7  int main()
8  {
9      Add<<<1, N>>>(a, b, c);
10 }
```

For a 1-dimensional grid: For a 2-dimensional grid:

```
tx = cuda.threadIdx.x
bx = cuda.blockIdx.x
bw = cuda.blockDim.x
i = tx + bx * bw
array[i] = compute(i)
```

```
tx = cuda.threadIdx.x
ty = cuda.threadIdx.y
bx = cuda.blockIdx.x
by = cuda.blockIdx.y
bw = cuda.blockDim.x
bh = cuda.blockDim.y
x = tx + bx * bw
y = ty + by * bh
array[x, y] = compute(x, y)
```

- Declared using:

Keyword	Executed on	Call on
__device__	Device	Device
__global__	Device	Host
__host__	Host	Host

- Executed with <<< X, Y >>>
 - X – Number of blocks per grid
 - Y – Threads per block

Example #2 – 2D kernel

```
1  ✓ global void Add(float* a[N][N], float* b[N][N], float* c[N][N])
2  {
3      int x = threadIdx.x;
4      int y = threadIdx.y;
5      c[x][y] = a[x][y] + b[x][y];
6  }
7
8  ✓ int main()
9  {
10     dim3 threadsPerBlock(N, N);
11     Add<<<1, threadsPerBlock>>>(a, b, c);
12 }
```

Example #3 – Memory allocation & data transfer

```
1  ✓ int main()
2  {
3      const unsigned int X=1048576; //1 Megabyte
4      const unsigned int bytes = X*sizeof(int);
5      int *hostArray= (int*)malloc(bytes);
6      int *deviceArray;
7      cudaMalloc((int**)&deviceArray,bytes);
8      memset(hostArray,0,bytes);
9      cudaMemcpy(deviceArray,hostArray,bytes,cudaMemcpyHostToDevice);
10     cudaMemcpy(hostArray,deviceArray,bytes,cudaMemcpyDeviceToHost);
11
12     cudaFree(deviceArray);
13
14 }
```

Allocate memory on the device

`cudaMalloc (void** devPtr, size_t size)`

Copy data between host and device

`cudaMemcpy (void* dst, const void* src, size_t count, cudaMemcpyKind kind)`

kind = 0 => Host -> Host

kind = 1 => Host -> Device

kind = 2 => Device -> Host

kind = 3 => Device -> Device

kind = 4 => Direction of the transfer is inferred from the pointer values. Requires unified virtual addressing

Free memory on the device

`cudaFree (void* devPtr)`

Example #4 – Reading GPU properties

```
1  ✓ int main()
2  {
3      cudaDeviceProp prop;
4      cudaGetDeviceProperties(&prop, device);
5
6      printf("GPU clock rate: %d\n", prop.clockRate);
7      printf("Device can concurrently copy memory and execute a kernel: %d\n", prop.deviceOverlap);
8  }
```

Return information about the compute-device

`cudaGetDeviceProperties (cudaDeviceProp* prop, int device)`

Conclusions

- Simple low-level operations (e.g., addition, multiplication...) can be more effectively parallelized on GPU than on CPU.
- GPU is not meant for complex operations such as slicing array or complex data comparisons.
- GPU has more complex data management structure than CPU does.