

CUDA/GPU 编程模型

周 斌 @NVIDIA & USTC 2014年3月

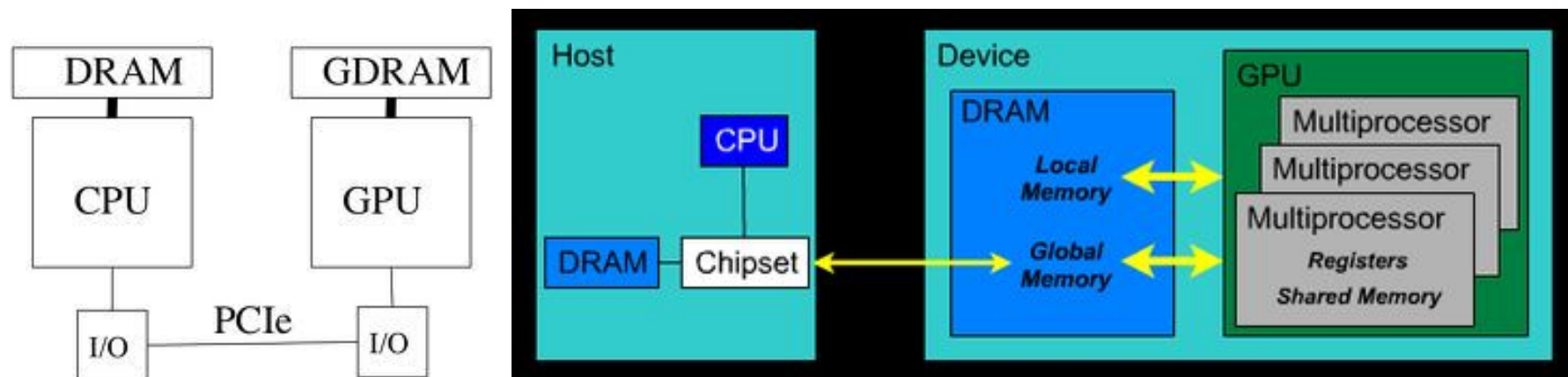
内容

- ▶ CPU和GPU互动模式
- ▶ GPU线程组织模型（不停强化）
- ▶ GPU存储模型
- ▶ 基本的编程问题

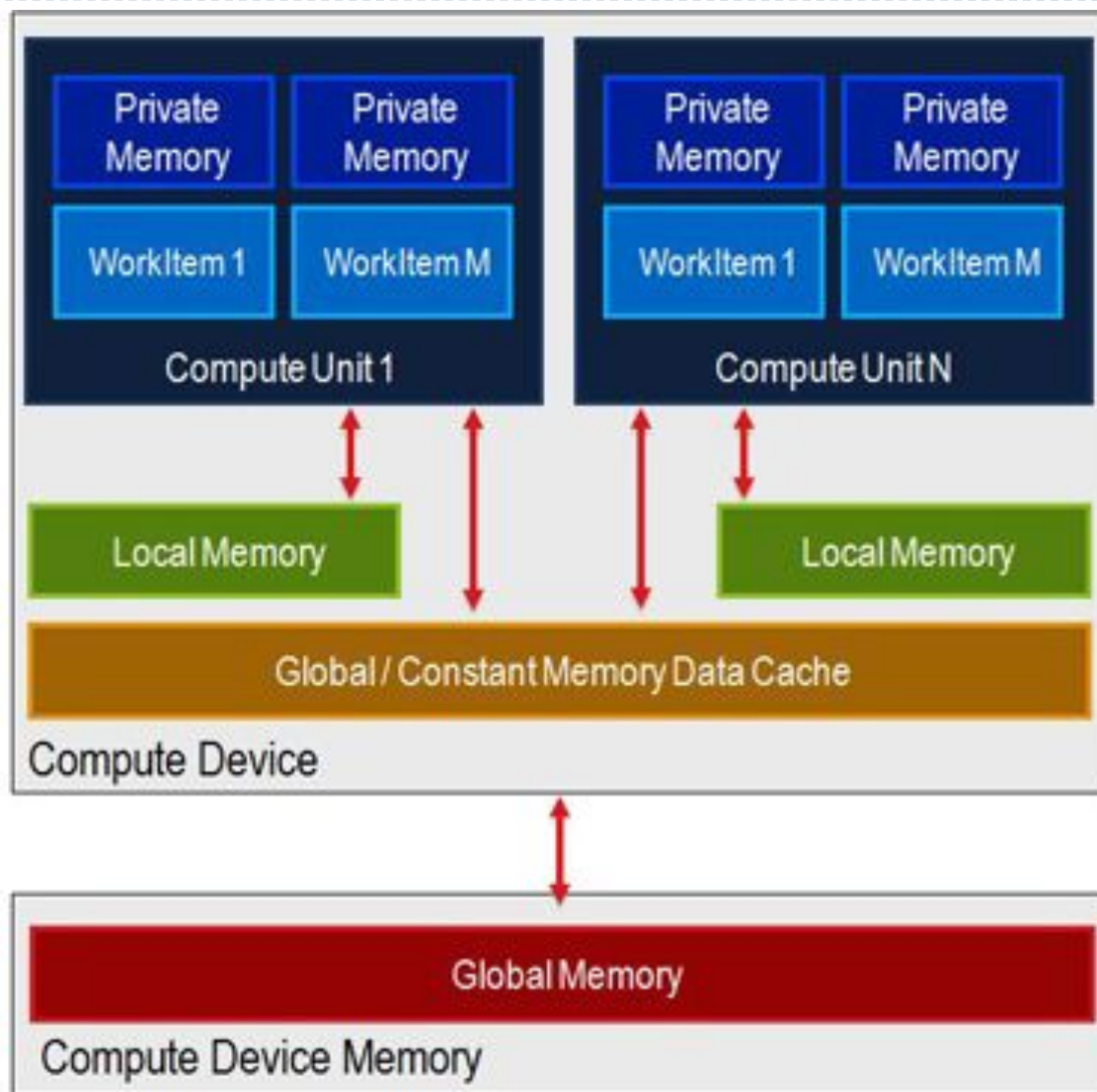


CPU-GPU 交互

- ▶ 各自的物理内存空间
- ▶ 通过PCIE总线互连(8GB/s~16GB/s)
- ▶ 交互开销较大



GPU存储器层次架构（硬件）

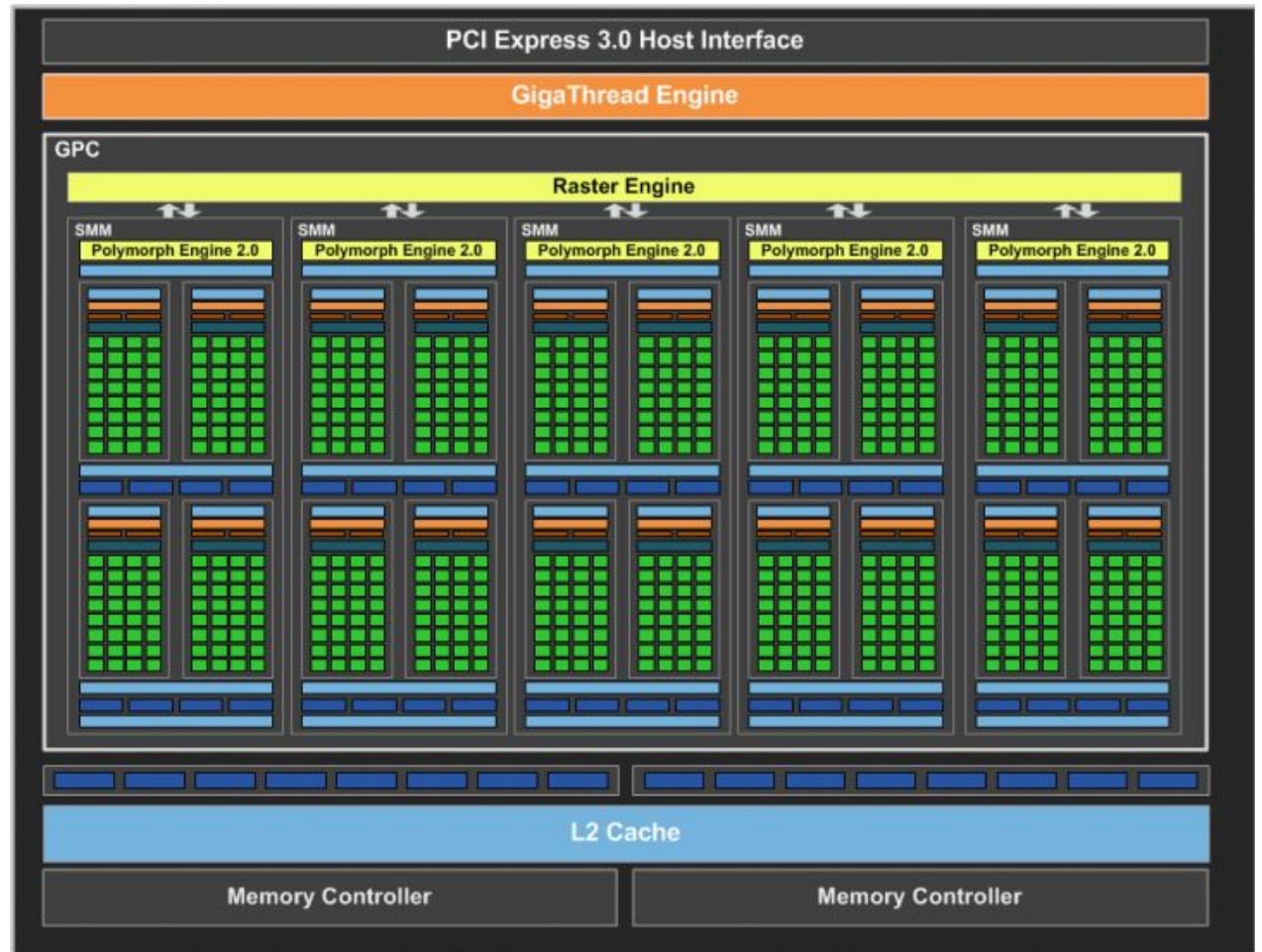
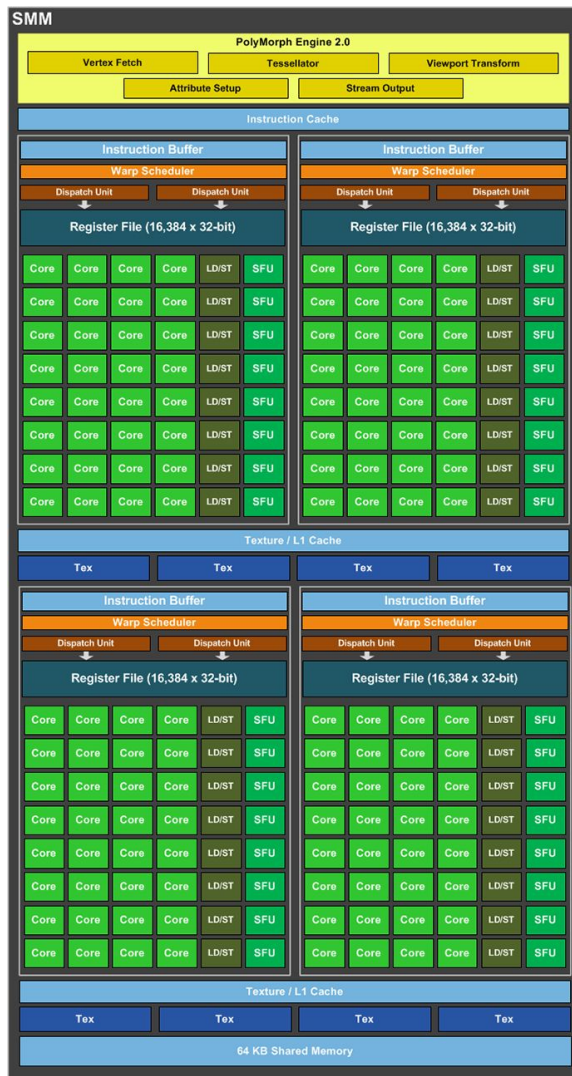


访存速度

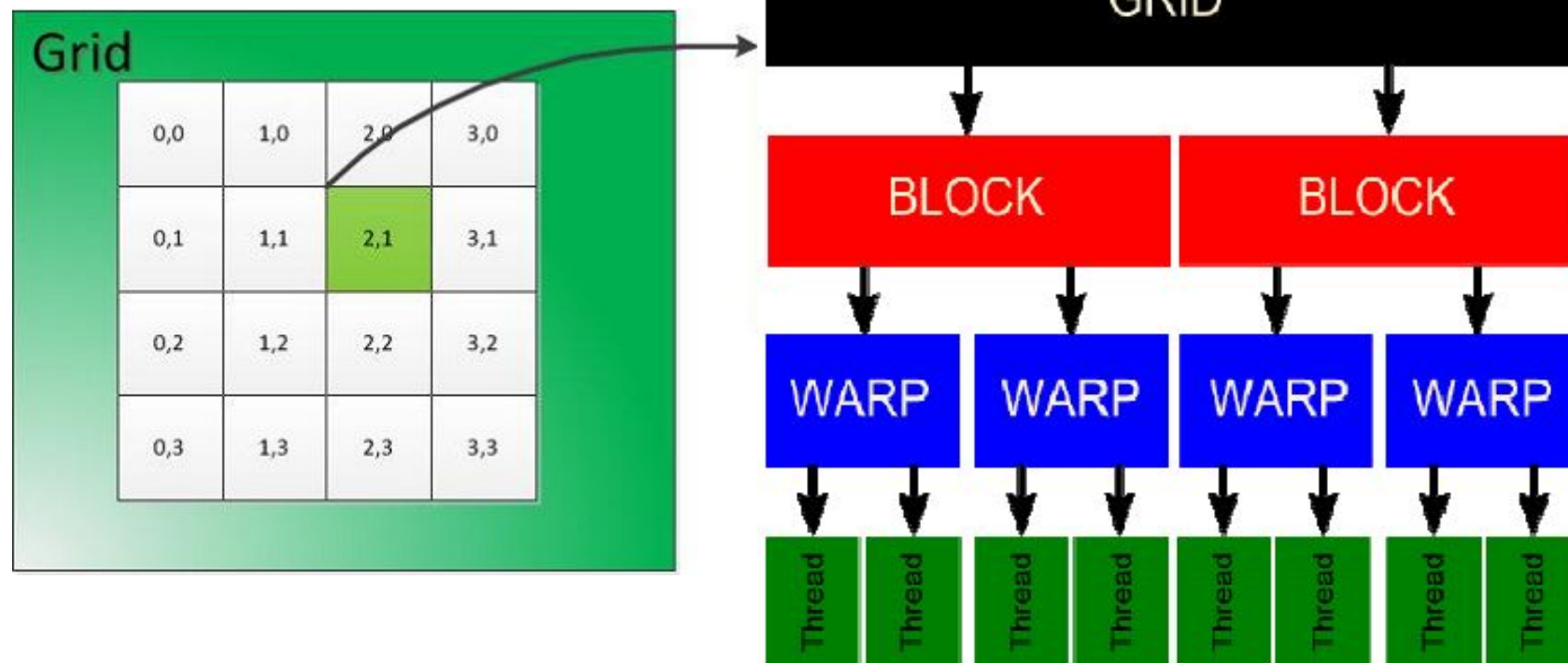
- ▶ Register – dedicated HW - single cycle
- ▶ Shared Memory – dedicated HW - single cycle
- ▶ Local Memory – DRAM, no cache - *slow*
- ▶ Global Memory – DRAM, no cache - *slow*
- ▶ Constant Memory – DRAM, cached, $1 \cdots 10s \cdots 100s$ of cycles, depending on cache locality
- ▶ Texture Memory – DRAM, cached, $1 \cdots 10s \cdots 100s$ of cycles, depending on cache locality
- ▶ Instruction Memory (invisible) – DRAM, cached



GPU架构回顾



GPU线程组织模型

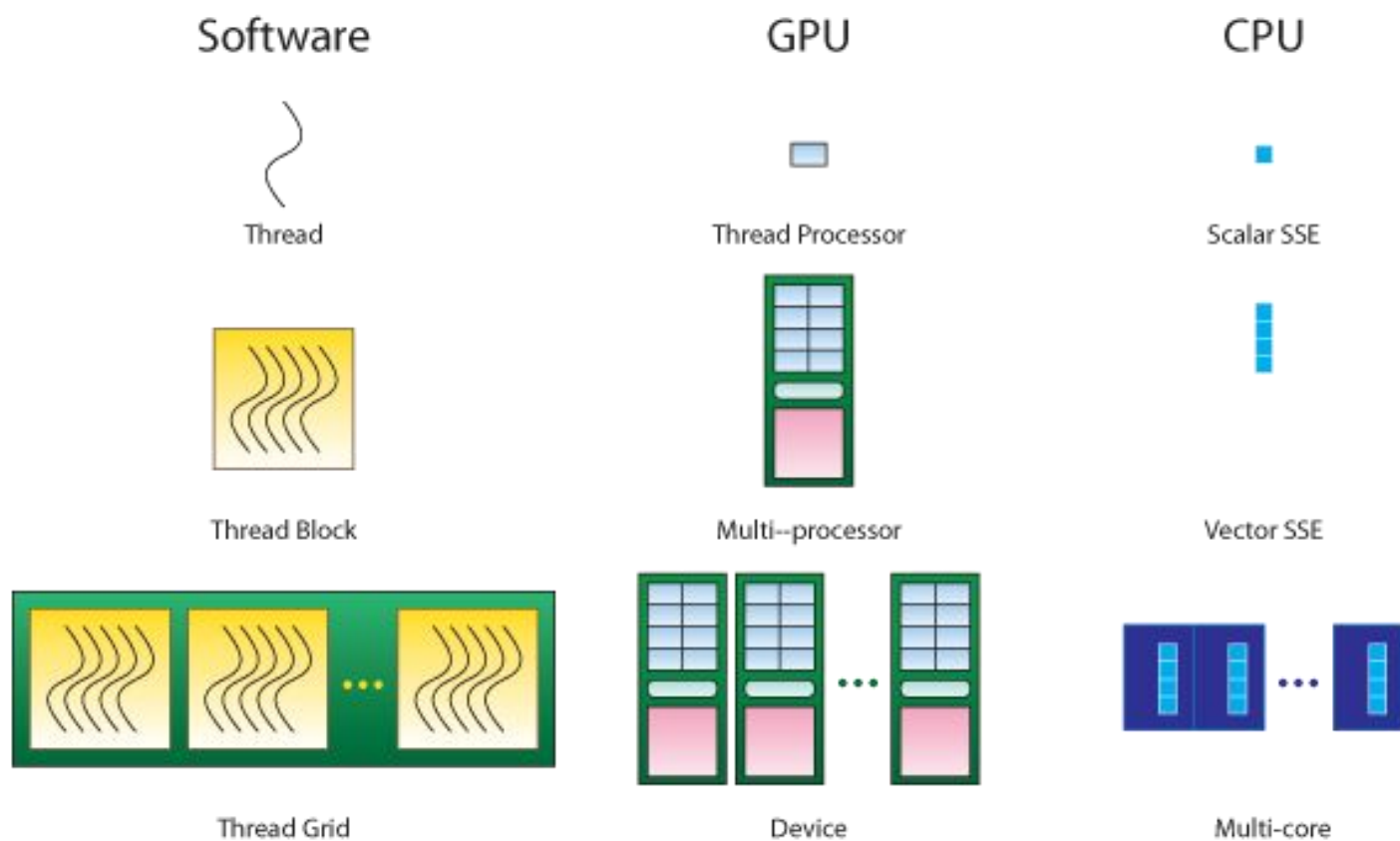


线程组织架构说明

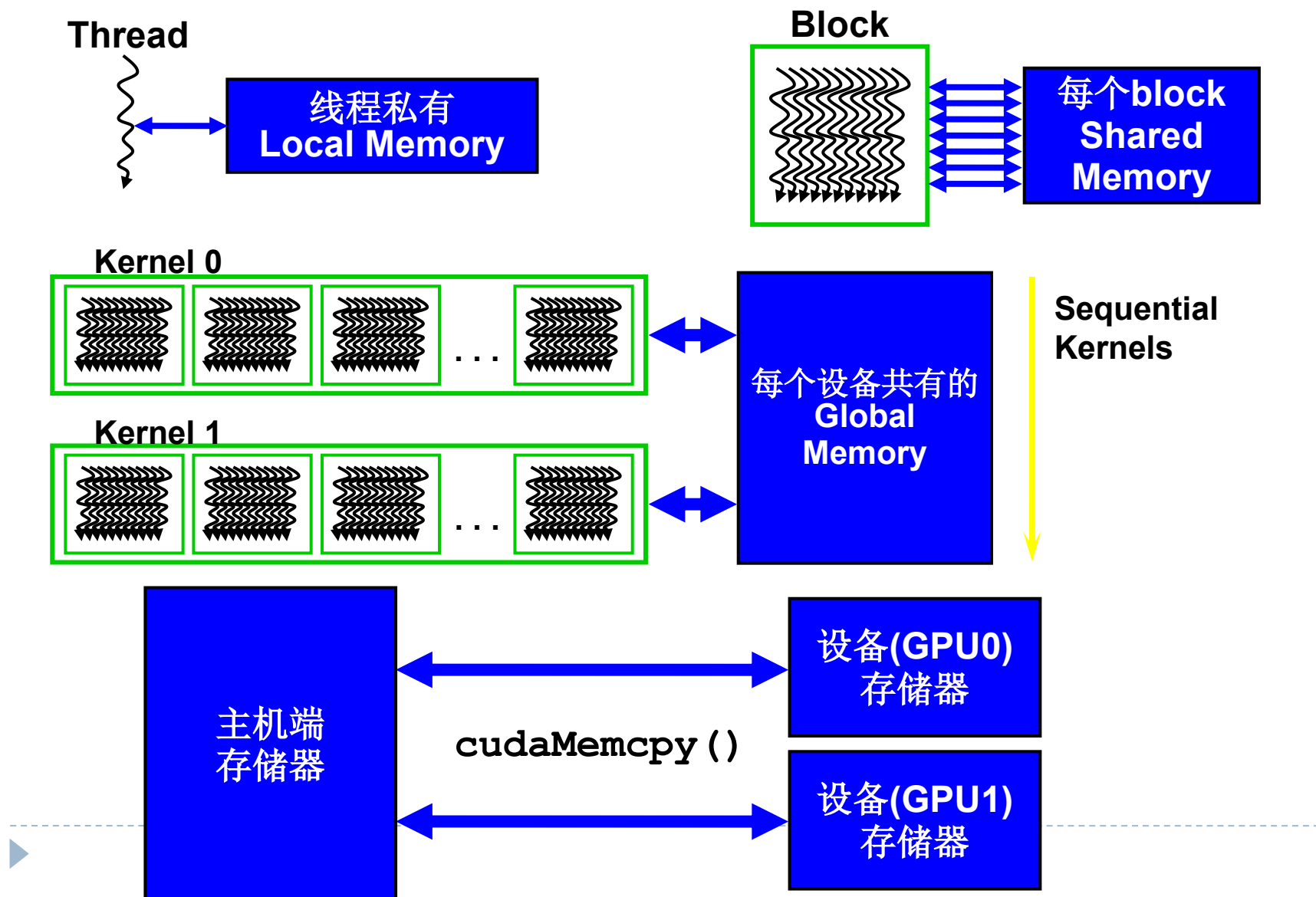
- ▶ 一个Kernel具有大量线程
- ▶ 线程被划分成线程块 ‘blocks’
 - ▶ 一个block内部的线程共享 ‘Shared Memory’
 - ▶ 可以同步 ‘_syncthreads()’
- ▶ Kernel启动一个 ‘grid’，包含若干线程块
 - ▶ 用户设定
- ▶ 线程和线程块具有唯一的标识

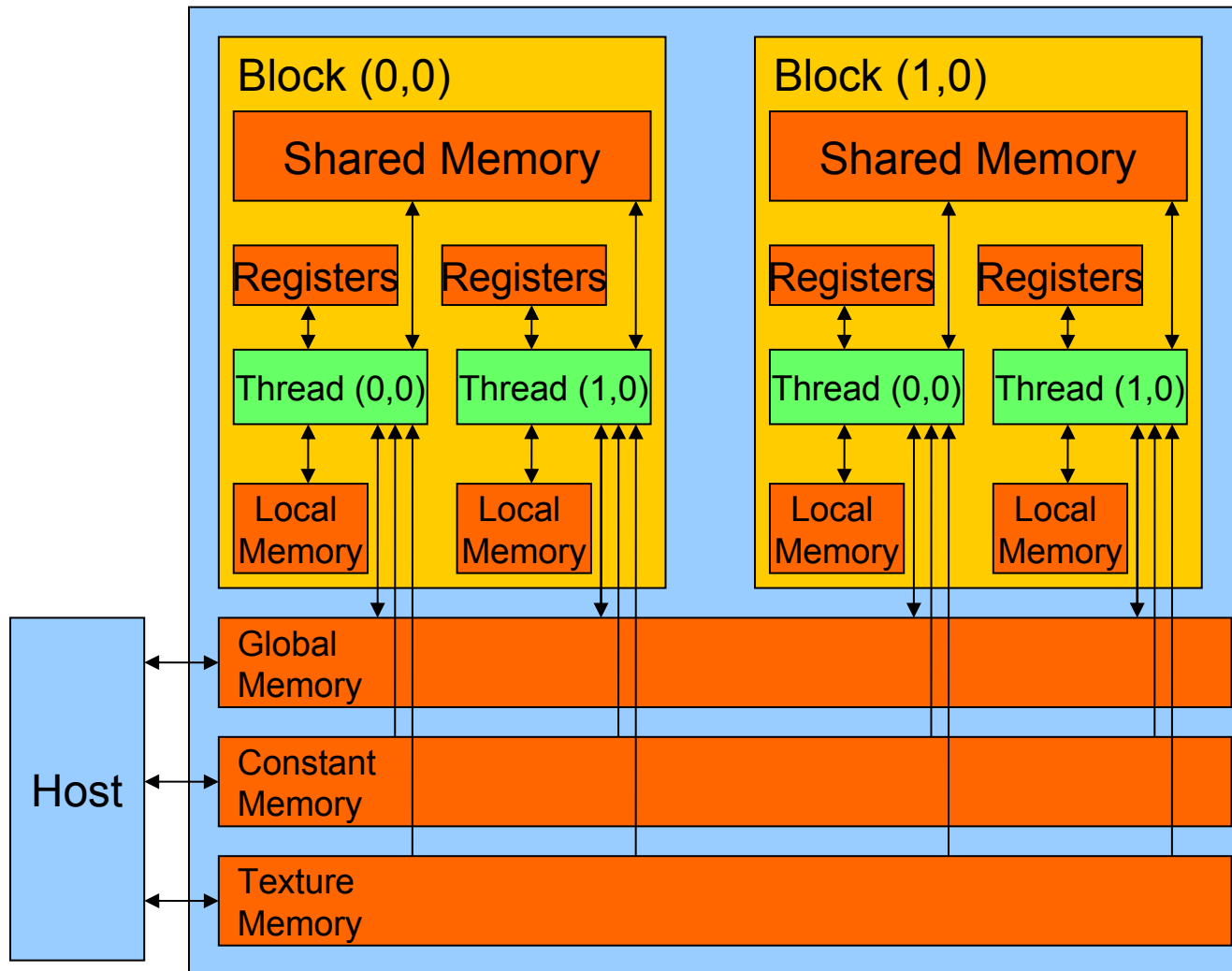


GPU线程映射关系



GPU内存和线程等关系





编程模型

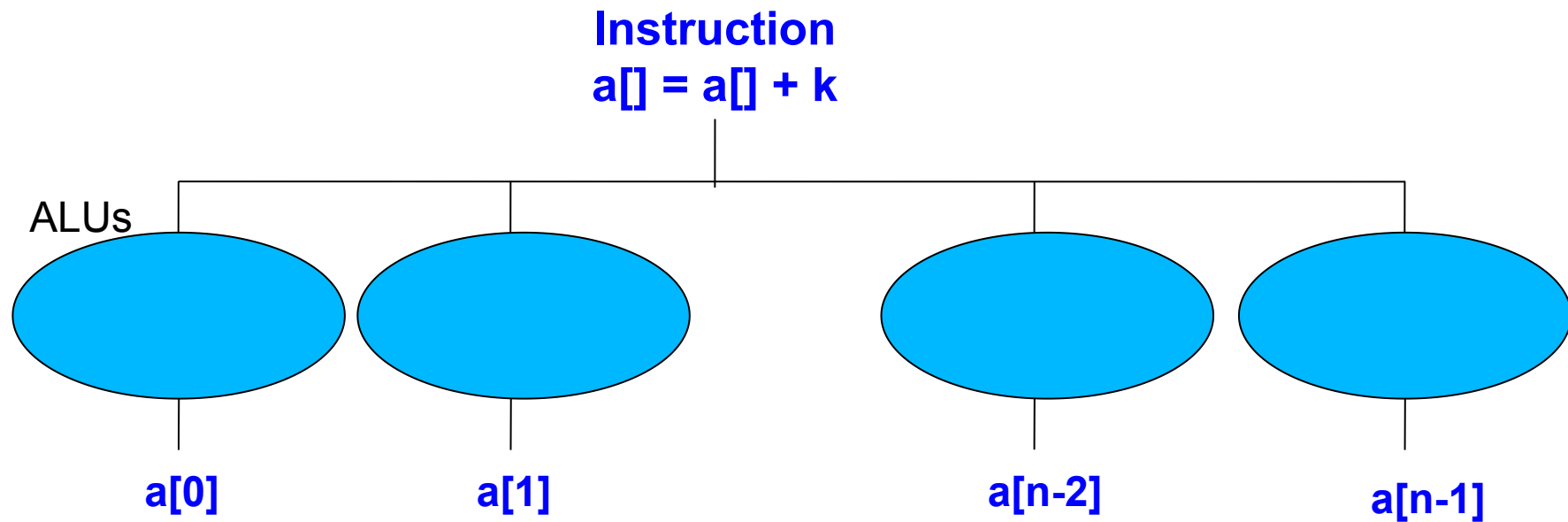
常规意义的GPU用于处理图形图像

操作于像素，每个像素的操作都类似

可以应用SIMD (single instruction multiple data)

SIMD (Single Instruction Multiple Data)

也可以认为是数据并行分割



Single Instruction Multiple Thread (SIMT)

GPU版本的 SIMD

大量线程模型获得高度并行

线程切换获得延迟掩藏

多个线程执行相同指令流

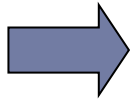
GPU上大量线程承载和调度



CUDA编程模式: Extended C

- ▶ **Declspecs**

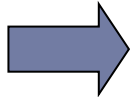
- ▶ global, device, shared, local, constant



```
__device__ float filter[N];  
  
__global__ void convolve (float *image) {  
    __shared__ float region[M];  
    ...
```

- ▶ **关键词**

- ▶ threadIdx, blockIdx



```
    region[threadIdx] = image[i];
```

- ▶ **Intrinsics**

- ▶ __syncthreads

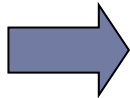
```
    __syncthreads()  
    ...
```

```
    image[j] = result;
```

```
}
```

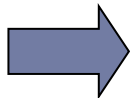
- ▶ **运行期API**

- ▶ Memory, symbol, execution management



```
// Allocate GPU memory  
void *myimage = cudaMalloc(bytes)
```

- ▶ **函数调用**



```
// 100 blocks, 10 threads per block  
convolve<<<100, 10>>> (myimage);
```



CUDA 函数声明

	执行 位置	调用 位置
<code>__device__ float DeviceFunc()</code>	device	device
<code>__global__ void KernelFunc()</code>	device	host
<code>__host__ float HostFunc()</code>	host	host

- ▶ `__global__` 定义一个 kernel 函数
 - ▶ 入口函数，CPU上调用，GPU上执行
 - ▶ 必须返回void
- ▶ `__device__` and `__host__` 可以同时使用

