

# CUDA

## 编程模型

NVIDIA亚太区技术市场经理 邓培智

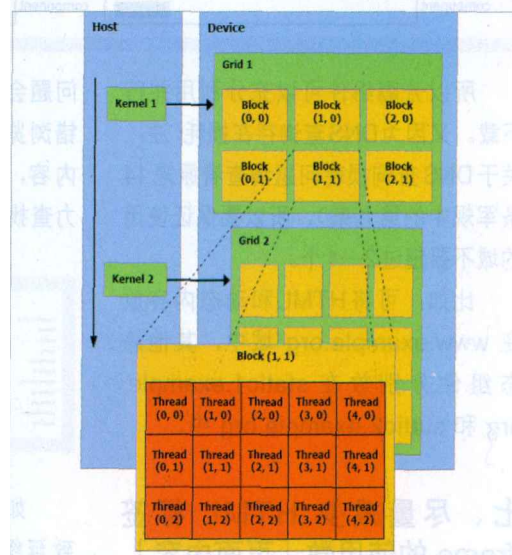
在CUDA的架构中，GPU可视为一个计算设备，是主机或者CPU的协处理器，用于处理高度并行的计算。GPU（或者称为“设备，device”）均具备自己的存储器（device memory，设备内存），可以并行地运行许多线程。在CUDA程序中，并行计算的部分可以被分离到一个被称为kernel（内核）的函数。在设备上许多线程执行同一个kernel。

CUDA的线程和CPU的线程有很大的不同。GPU线程是超级轻量的，创建线程的开销非常小，而且切换几乎是立即的。现在多核的CPU之需要有限的几个线程就可以使得CPU满负荷工作，但是GPU需要上千的线程才能够使之获得足够高的效率。

在设备中，执行内核的许多线程组织为线程块网格（grid）的方式。其组织如图所示：

线程块包含有一批线程SIMD（单指令多数数据流）方式执行的线程，每个线程都有自己的ID，同时每个线程块也有自己的ID。在每个块内的线程可以相互协作：可以共享结果以节省计算；可以同步；线程块可以对共享内存访问，从而大大节省设备内存带宽的消耗。可以说线程协作是CUDA强大性能之一。不过，不同线程块内的线程不能进行协作。

将线程组织为线程块的方式会带来线程协作的降低，因为不同块之间的线程不能协作，但是



好处是可以带来透明的硬件伸缩性。因为不同的设备有不同的并行能力。在并行度很高的设备上执行时很多块可以并行地执行；反之在较低并行度的设备上并行执行的块可以相对较少。

线程块和线程都有ID，其中线程ID可以是1D-3D的，3D的ID如Thread (x,y,z)；而线程块则可以是1D或者2D的。在图像处理等多维数据的场合，可以简化内存寻址。

在设备上执行的线程可以访问多种内存，包括设备上的DRAM和GPU芯片内置的寄存器和共享内存。下图是内存模型的示意图。线程、线程块和grid分别可以访问的内存如下：

# CUDA: 大规模并行计算的利器

线程可以读写寄存器

线程可以读写本地内存

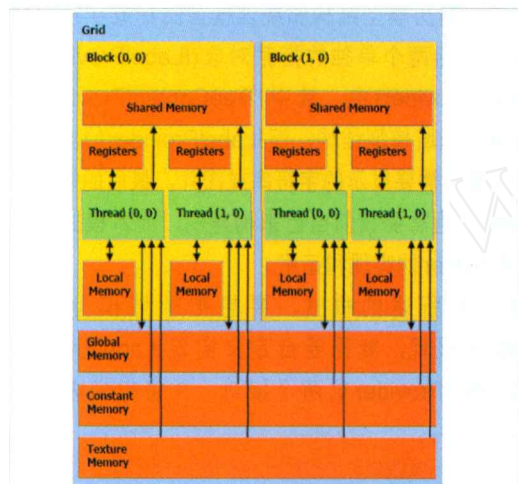
线程块可以读写共享内存

Grid 可以读写全局内存

Grid 对常量内存是只读的

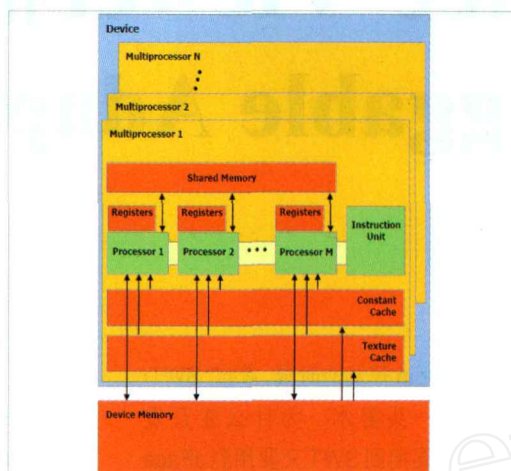
Grid 对纹理内存是只读的

主机可以读写存储在设备 DRAM 中的全局内存、常量内存和纹理内存。



共享内存位于芯片内部，其速度等同于寄存器。其速度比没有高速缓存的全局内存快得多。将线程按照块的方式进行组织，在块内进行线程间的协作以及多次的计算时使用共享内存，可以大大提高计算和存储的效率，并且极大地减少所需设备 DRAM 的带宽。

在硬件层面上，设备或者说 CUDA GPU 包含有数量不一的 Multiprocessor。每一个 Multiprocessor 里面有一组单指令多数据流架构的 32 位处理器。在每一个周期内，一个 multiprocessor 在被称为 warp 的一组线程上执行同样的指令，每一个 warp 包含的线程数量被称为 warp size。程序员所面对的本地、全局、常量和纹理存储器实际在物理上都位于设备存储器（或者 GPU DRAM）中，而在芯片内每个 multiprocessor 还包含有 32 位的寄存器（每个处理器均有一组）、一个芯片内的共享内存、一个只读的常量高速缓存以及一个只读的纹理高速缓存。常量高速缓存和纹理高速缓存可以加速常量存储器和纹理存储器的访问。CUDA GPU 硬件结构的示意图如下：



在 GPU 程序执行的时候，每个 grid 的线程块被拆分到 warp 以内，每个线程块只能被一个 multiprocessor 执行，因此共享内存空间就固定在芯片内的共享内存内了。寄存器被分配给线程。如果一个 kernel 要求过多的寄存器则会出现启动失败。一个 multiprocessor 可以同时执行多个线程块，片内的存储资源则会被分配给这些并发的线程块中。如果线程块和线程减少对共享内存和寄存器的使用，就会使得并发的线程块数量增加，从而获得更高的并行度。因此在编程的时候需要注意节约这些资源的使用。

以上是 CUDA 软硬件模型的一些描述。如果我们对编程模型做个简洁的总结的话可以用下面这些公式：

设备 (device) = GPU = 一组 multiprocessor

Multiprocessor = 一组处理器以及共享内存

Kernel = GPU 程序

Grid = 执行 kernel 的线程块的阵列

线程块 = 执行一个 kernel 并能够通过共享内存进行通信的一组 SIMD 线程

NVIDIA 目前推出了一系列的支持 CUDA 的 GPU，这些 GPU 具备不同的并行计算能力，每种 GPU 都有不同数量的 multiprocessor。这些 GPU 的架构是非常类似的，但是每种 GPU 的 multiprocessor 具备的内部资源以及 warp size 等参数可能会有区别。详细的情况可以进入 [http://www.nvidia.com/object/cuda\\_home.html](http://www.nvidia.com/object/cuda_home.html) 去查阅。更多有关 CUDA 编程的信息请参阅 CUDA 编程指南。