Advanced Architectures



Master Informatics Eng.

2019/20 *A.J.Proença*

Data Parallelism 3 (GPU/CUDA)

(most slides are borrowed)

The CUDA programming model



- Compute Unified Device Architecture
- CUDA is a recent programming model, designed for
 - a multicore CPU *host* coupled to a many-core *device*, where
 - devices have wide SIMD/SIMT parallelism, and
 - the *host* and the *device* do not share memory
- CUDA provides:
 - a thread abstraction to deal with SIMD
 - synchr. & data sharing between small groups of threads
- CUDA programs are written in C with extensions
- OpenCL inspired by CUDA, but hw & sw vendor neutral
 - programming model essentially identical

CUDA Devices and Threads



- A compute device
 - is a coprocessor to the CPU or host
 - has its own DRAM (device memory)
 - runs many threads in parallel
 - is typically a GPU but can also be another type of parallel processing device
- Data-parallel portions of an application are expressed as device kernels which run on many threads - SIMT
- Differences between GPU and CPU threads
 - GPU threads are extremely lightweight
 - very little creation overhead, requires LARGE register bank
 - GPU needs 1000s of threads for full efficiency
 - multi-core CPU needs only a few

CUDA basic model: Single-Program Multiple-Data (SPMD)



Programming Model: SPMD + SIMT/SIMD

\sim

- Hierarchy
 - Device => Grids
 - Grid => Blocks
 - Block => Warps
 - Warp => Threads
- Single kernel runs on multiple blocks (SPMD)
- Threads within a warp are executed in a lock-step way called singleinstruction multiple-thread (SIMT)
- Single instruction are executed on multiple threads (SIMD)
 - Warp size defines SIMD granularity (32 threads)
- Synchronization within a block uses shared memory



The Computational Grid: Block IDs and Thread IDs



AJProença, Advanced Architectures, MiEI, UMinho, 2019/20

6

Example

公

- Multiply two vectors of length 8192
 - Code that works over all elements is the grid
 - Thread blocks break this down into manageable sizes
 - 512 threads per block
 - SIMD instruction executes 32 elements at a time
 - Thus grid size = 16 blocks
 - Block is analogous to a <u>strip-mined vector loop</u> with vector length of 32
 - Block is assigned to a *multithreaded SIMD processor* by the *thread block scheduler*
 - Current-generation GPUs (Fermi) have 7-16 multithreaded SIMD processors

Strip-mining and Cleanup

Strip-mining, also known as loop sectioning, is a loop transformation technique for enabling SIMD-encodings of loops, as well as a means of improving memory performance. By fragmenting a large loop into smaller segments or strips, this technique transforms the loop structure in two ways:

- It increases the temporal and spatial locality in the data cache if the data are reusable in different passes of an algorithm.
- It reduces the number of iterations of the loop by a factor of the length of each vector, or number of operations being performed per SIMD operation. In the case of Streaming SIMD Extensions, this vector or strip-length is reduced by 4 times: four floating-point data items per single Streaming SIMD Extensions single-precision floating-point SIMD operation are processed.

First introduced for vectorizers, this technique consists of the generation of code when each vector operation is done for a size less than or equal to the maximum vector length on a given vector machine.

The compiler automatically strip-mines your loop and generates a cleanup loop. The following examples demonstrate strip mining and cleaning up loops.

i = 1 do while (i<=n) a(i) = b(i) + c(i) ! Original loop code i = i + 1 end do</pre>

Example 2: After Vectorization

```
!The vectorizer generates the following two loops
i = 1
do while (i < (n - mod(n,4)))
! Vector strip-mined loop.
a(i:i+3) = b(i:i+3) + c(i:i+3)
i = i + 4
end do
do while (i <= n)
a(i) = b(i) + c(i)  !Scalar clean-up loop
i = i + 1
end do
```



\sim

C with CUDA Extensions: C with a few keywords

```
void saxpy_serial(int n, float a, float *x, float *y)
       for (int i = 0; i < n; ++i)
          y[i] = a*x[i] + y[i];
                                                    Standard C Code
   // Invoke serial SAXPY kernel
   saxpy_serial(n, 2.0, x, y);
    int i = blockIdx.x*blockDim.x + threadIdx.x:
       if (i < n) y[i] = a*x[i] + y[i];
                                                      Parallel C Code
   // Invoke parallel SAXPY kernel with 256 threads/block
   int nblocks = (n + 255) / 256;
   saxpy_parallel<<<nblocks, 256>>>(n, 2.0, x, y);
NVIDIA Confidential
```

Terminology (and in NVidia)

シン

- Threads of SIMD instructions (warps)
 - Each has its own IP (up to 48/64 per SIMD processor, Fermi/Kepler)
 - Thread scheduler uses scoreboard to dispatch
 - No data dependencies between threads!
 - Threads are organized into blocks & executed in groups of 32 threads (*thread block*)
 - Blocks are organized into a grid
- The <u>thread block scheduler</u> schedules blocks to SIMD processors (Streaming Multiprocessors)
- Within each SIMD processor:
 - 32 SIMD lanes (thread processors)
 - Wide and shallow compared to vector processors

CUDA Thread Block

公

- Programmer declares (Thread) Block:
 - Block size 1 to 512 concurrent threads
 - Block shape 1D, 2D, or 3D
 - Block dimensions in threads
- All threads in a Block execute the same thread program
- Threads share data and synchronize while doing their share of the work
- Threads have thread id numbers within Block
- Thread program uses thread id to select work and address shared data

CUDA Thread Block



Parallel Memory Sharing



CUDA Memory Model Overview



公

- R/W per-thread registers
- R/W per-thread local memory
- R/W per-block shared memory
- R/W per-grid global memory
- Read only per-grid constant memory
- Read only per-grid texture memory
- The host can R/W global, constant, and texture memories



Terminology: CUDA and OpenCL

 $\langle \rangle$

CUDA and OpenCL



Hardware Implementation: Memory Architecture

公

- Device memory (DRAM)
 - Slow (2~300 cycles)
 - <u>Local</u>, global, constant, and texture memory
- On-chip memory
 - Fast (1 cycle)
 - Registers, shared memory, constant/texture cache



Courtesy NVIDIA



Example

AJProença, Advanced Architectures, MiEI, UMinho, 2019/20

17

Vector Processor versus CUDA core



Copyright © 2012, Elsevier Inc. All rights reserved.

Conditional Branching

\sim

- Like vector architectures, GPU branch hardware uses internal masks
- Also uses
 - Branch synchronization stack
 - entries consist of masks for each SIMD lane
 - i.e. which threads commit their results (all threads execute)
 - Instruction markers to manage when a branch diverges into multiple execution paths
 - push on divergent branch
 - ...and when paths converge
 - act as barriers
 - pops stack
- Per-thread-lane 1-bit predicate register, specified by programmer