Master Informatics Eng.

2016/17 A.J.Proença

Data Parallelism 3 (GPU & CUDA) (most slides are borrowed)

AJProença, Advanced Architectures, MiEl, UMinho, 2016/17

Beyond Vector/SIMD architectures

\sim

XX

Vector/SIMD-extended architectures are hybrid approaches

- mix (super)scalar + vector op capabilities on a single device
- highly pipelined approach to reduce memory access penalty
- tightly-closed access to shared memory: lower latency

• Evolution of Vector/SIMD-extended architectures

- CPU cores with wider vectors and/or SIMD cores:
 - <u>DSP</u> VLIW cores with vector capabilities: **Texas Instruments** (...?)
 - PPC cores coupled with SIMD cores: Cell (past...), IBM Power BQC...
 - <u>ARM64</u> cores coupled with SIMD cores: from Tegra to Parker (**NVidia**) (...?)
 - x86 many-core: Intel MIC / Xeon KNL, PEZY-SC, AMD FirePro...
 - other many-core: ShenWay 260, Adapteva Epiphany-V...

- coprocessors (require a host scalar processor): accelerator devices

- on disjoint physical memories (e.g., Xeon KNC with PCI-Expr, PEZY-SC)
- focus on SIMT/SIMD to hide memory latency: GPU-type approach
- ISA-free architectures, code compiled to silica: **FPGA**

What is an FPGA

\sim

Field-Programmable Gate Arrays (FPGA)

A fabric with 1000s of simple configurable logic cells with LUTs, on-chip SRAM, configurable routing and I/O cells



FPGA as a multiple configurable ISA



Beyond Vector/SIMD architectures

\sim

- Vector/SIMD-extended architectures are hybrid approaches
 - mix (super)scalar + vector op capabilities on a single device
 - highly pipelined approach to reduce memory access penalty
 - tightly-closed access to shared memory: lower latency

Evolution of Vector/SIMD-extended architectures

- CPU cores with wider vectors and/or SIMD cores:
 - <u>DSP</u> VLIW cores with vector capabilities: **Texas Instruments** (...?)
 - <u>PPC</u> cores coupled with SIMD cores: **Cell** (past...) , **IBM Power BQC...**
 - <u>ARM64</u> cores coupled with SIMD cores: from Tegra to Parker (NVidia) (...?)
 - x86 many-core: Intel MIC / Xeon KNL, PEZY-SC, AMD FirePro...
 - other many-core: **ShenWay** 260, **Adapteva** Epiphany-V...

- coprocessors (require a host scalar processor): accelerator devices

- on disjoint physical memories (e.g., Xeon KNC with PCI-Expr, PEZY-SC)
- focus on SIMT/SIMD to hide memory latency: **GPU**-type approach
- ISA-free architectures, code compiled to silica: FPGA

AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

Graphical Processing Units

- Question to GPU architects:
 - Given the hardware invested to do graphics well, how can we supplement it to improve the performance of a wider range of applications?

Key ideas:

- Heterogeneous execution model
 - CPU is the *host*, GPU is the *device*
- Develop a C-like programming language for GPU
- Unify all forms of GPU parallelism as CUDA_threads
- Programming model follows SIMT:
 "Single Instruction Multiple Thread "



5

Graphical Processing Units

cores/processing elements in several devices



Theoretical peak performance in several computing devices (DP)





NVIDIA GPU Architecture

- Similarities to vector machines:
 - Works well with data-level parallel problems
 - Scatter-gather transfers
 - Mask registers
 - Large register files

Differences:

- No scalar processor
- Uses multithreading to hide memory latency
- Has many functional units, as opposed to a few deeply pipelined units like a vector processor





The GPU as a compute device: the G80

NVidia GPU structure & scalability



AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

The NVidia Fermi architecture



AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

Fermi Architecture Innovations

- Each SIMD processor has
- Each SIMD processor has
 Two SIMD thread schedulers, two instruction dispatch units
 16 SIMD lanes (SIMD width=32, chime=2 cycles), 16 load-store units, 4 special function units
 Thus, two threads of SIMD instructions are scheduled every two clock cycles clock cycles
- Fast double precision
- Caches for GPU memory
- 64-bit addressing and unified address space
- Error correcting codes
- Faster context switching
- **Faster atomic instructions**



NVIDIA GPU Memory Structures

- Each SIMD Lane has private section of off-chip DRAM
 - "Private memory" (Local Memory)
 - Contains stack frame, spilling registers, and private variables
- Each multithreaded SIMD processor also has local memory (Shared Memory)
 - Shared by SIMD lanes / threads within a block
- Memory shared by SIMD processors is GPU Memory (Global Memory)
 - Host can read and write GPU memory







From Fermi into Kepler: The Memory Hierarchy



AJProença, Advanced Architectures, MiEI, UMinho, 2016/17



AJProença, Advanced Architectures, MiEI, UMinho, 2016/17



From Fermi into Kepler: Compute capabilities

	FERMI GF100	FERMI GF104	KEPLER GK104	KEPLER GK110				
Compute Capability	2.0	2.1	3.0	3.5				
Threads / Warp	32	32	32	32				
Max Warps / Multiprocessor	48	48	64	64				
Max Threads / Multiprocessor	1536	1536	2048	2048				
Max Thread Blocks / Multiprocessor	8	8	16	16				
32-bit Registers / Multiprocessor	32768	32768	65536	65536				
Max Registers / Thread	63	63	63	255				
Max Threads / Thread Block	1024	1024	1024	1024				
Shared Memory Size Configurations (bytes)	16K	16K	16K	16K				
	48K	48K	32K	32K				
			48K	48K				
Max X Grid Dimension	2^16-1	2^16-1	2^32-1	2^32-1				
Hyper-Q	No	No	No	Yes				
Dynamic Parallelism	No	No	No	Yes				







AJProença, Advanced Architectures, MiEl, UMinho, 2016/17

23



Families in NVidia Tesla GPUs



Families in NVidia Tesla GPUs



Families in NVidia GPU

Nvidia Tesla Workstation GPU Specifications Comparison									
	Full P2xx	P100	M40	K80	К40				
Family	Pascal (2nd-gen)	Pascal	Maxwell (2nd-gen)	Kepler (2nd-gen)	Kepler				
Architecture	N/A	GP100	GM200	GK210	GK110				
Cores	3840	3584	3072	2 x 2496	2880				
ROPs	N/A	N/A	96	96	48				
Texture Units	240	224	192	416	240				
Core Clock	N/A	1328MHz	948MHz	562MHz	745MHz				
Memory Clock	N/A	1400MHz	1500MHz	2500MHz	3004MHz				
Memory Bandwidth	N/A	N/A	288GB/s	480GB/s	288GB/s				
Memory Bus Width	4096-bit	4096-bit	384-bit	2x 384-bit	384-bit				
Memory Type	HBM2	HBM2	GDDR5	GDDR5	GDDR5				
Memory Size	16GB or higher	16GB	24GB	2 x 12GB	12GB				
Die Size	N/A	610mm ²	601mm ²	561mm ²	551mm ²				
Transistors	N/A	15.3 billion	8 billion	2 x 7.08 billion	7.08 billion				
Register File Size / SM	256KB	256KB	256KB	512KB	256KB				
L2 Cache	N/A	4MB	3MB	1.5MB	1.5MB				
TDP	N/A	300W	250W	300W	235W				
Manufacturing Process	TSMC 16nm	TSMC 16nm	TSMC 28nm	TSMC 28nm	TSMC 28nm				
Release Date	N/A	Jul-16	Nov-15	Nov-14	Nov-13				

AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

1

<form>

Top500: Accelerator distribution over all 500 systems



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

XX

- 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

AJProença, Advanced Architectures, MiEI, UMinho, 2016/17





AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009 ECE 498AL, University of Illinois, Urbana-Champaign

Programming Model: SPMD + SIMT/SIMD

XX

- 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



AJProença, Advanced Architectures, MiEI, UMinho, 2016/17



The Computational Grid: Block IDs and Thread IDs



³⁴

Example



- 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 strip-mined vector loop 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



Copyright © 2012, Elsevier Inc. All rights reserved.

Graphical Processing Units

Terminology (and in NVidia)

- 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

\sim

- 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

AJProença, Advanced Architectures, MiEI, UMinho, 2016/17





Parallel Memory Sharing

CUDA Memory Model Overview

\sim



Hardware Implementation: Memory Architecture



Device memory

Courtesy NVIDIA

NVIDIA GPU Memory Structures

- Each SIMD Lane has private section of off-chip DRAM
 - "Private memory" (Local Memory)
 - Contains stack frame, spilling registers, and private variables
- Each multithreaded SIMD processor also has local memory (Shared Memory)
 - Shared by SIMD lanes / threads within a block
- Memory shared by SIMD processors is GPU Memory (Global Memory)
 - Host can read and write GPU memory







Vector Processor versus CUDA core



Copyright $\ensuremath{\textcircled{O}}$ 2012, Elsevier Inc. All rights reserved.

Graphical Processing Units

Conditional Branching

- 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

