<section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><text><text><text><text><image><image>

Parallel computing

Parallel computing is not new but:

- often it is now *massively* parallel
- we cannot wait until our old serial program runs faster (it never will)
- hence today programming *is* parallel programming
- we are probably restrained in first thinking about an algorithm in a serial way (implicit serialism in programming languages) **example**

Many tools are tried for parallel computing:

- very explicit ones like MPI (for cluster machines and supercomputers)
- lightweight language extensions such as OpenMP, OpenACC, Array Building Blocks (ArBB), Cilk Plus, ...
- domain-specific languages: CUDA, OpenCL, OpenGL, ...
- intelligent compilers, automatic parallelizers: PGI Compilers



GPU computing



traditional interpretation of GPU computing



GPU computing



traditional interpretation of GPU computing

M. Weigel (Coventry/Mainz)

GPU basics

08/10/2018 4/36

GPU computing

۲

0

The	oretical GFLOP/s						
	5750						
	5500						
	5250			GeForce 780	Ti 🦻 👘		
	5000	NVIDIA GPU Single Precisio	n				
	4750	NVIDIA GPU Double Precisi	n				
	4500	Intel CPU Double Precision	GeF	orce GTX TITAN			
	4250	Intel CPU Single Precision		/			
1. Sec. 1. Sec	4000			/			
	3750						
	3500						
	3250						
	3000		GeForce	GTX 680			
	2750						
	2500						
	2250						
	2000						
	1750		GeForce GTX 580	TI- KOOV	Tesla K40		
	1500	G	Force GTX 480	e Testa K20X			
	1250	GeForce	TX 280				
	1000	GoForce 8800 GTX	Tes	a M2090			
	750		Tesla C2050				
	500	GeForce 7800 GTX	rpertown	Sandy Bridge	lvy Bridge		
	250 GeForce FX	5800 Woodcrest					
	Apr-01 Sep-	02 Jan-04 May-05 Oct-06	Bloomfield We	stmere	ug-13 Dec-14		
	Арг-от зер-	02 Jun-04 May-05 Oct-00	eb-00 501-07 110	Apr-12 A	dg-15 Dec-14		
Core i7 Ivyl	Bridge i7-3	3870: 122 GFLOP	s				
				`			
NVIDIA Te	sla K10: 4	4580 GFLOP/s (sir	gle precisio	n)			
M. Weigel (Coventr	y/Mainz)	GPU E	asics			08/10/2018	4/36

GPUs in sicentific computing



GPU computing



- $\, \bullet \,$ Core i7 IvyBridge i7-3870: \approx 21 GB/s
- NVIDIA Tesla K10: 320 GB/s

Outline			
1 Latency vs. throughput			
② GPU architecture			
3 Execution model			
4 CUDA Programming			
5 Memory hierarchy			
6 The CUDA distribution			
M. Weigel (Coventry/Mainz)	GPU basics	08/10/2018	6 / 36

CPU vs. GPU hardware



Latency vs. throughput



- CPU must minimize latency of individual thread for responsiveness
- GPU hides latency through interleaved execution

General processing flow



GPU basics

General processing flow



Latency vs. throughput



General processing flow



Latency vs. throughput





GPGPU history

NVIDIA	Other architectures	Main memory	GigaT
 introduced CUDA in 2007 developed into a fully blown ecosystem series of computing cards academic support programs CUDA professorships CUDA research centers CUDA teaching centers 	 Intel MIC started as Larabee in 2006 prototype board Knight's Ferry (2010) 32 cores with 4 threads/core 2 GB DDR5 memory Knight's Corner, released in 2012 more than 50 cores per chip Knight's Landing, released in 2016 up to 72 cores per board used in a number of recent supercomputers 	 up to 16GB in current GPUs maximum bandwidth up to 900 GB/s on Volta (V100) accessible from CPU and GPU sides large latency (see below) optional error correction (ECC on/off, Fermi onwards) 	Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Scheduler Sche
 ATI ATI Stream introduced in 2007 less viral marketing higher peak performance somewhat less flexible architecture 	 (some) special-purpose machines ANTON by D. E. Shaw research, composed of purpose-built ASICs Janus, FPGA based machine for spin models QPace, based on the Cell processor known from Playstation 3 GRAPE, based on FPGAs and used for astrophysical N-body simyulations 	Several multiprocessors similar to a multi-core CPU each has its own set of registers, scheduler, caches etc. + scheduling units, PCIe logic etc.	
M. Weigel (Coventry/Mainz)	GPU basics 08/10/2018 12/30	6 M. Weigel (Coventry/Mainz) GPU basics	

GPU architecture

GPU architecture: streaming multiprocessor (Fermi)



Sc	hedul	ler	S	chedu	uler		
Di	spato	:h		Dispat	ch		
Core	Core	Core	Core	Register Register	05	Г)
Core	Core	Core	Core	Register Register	55		
Core	Core	Core	Core	Register			
Core	Core	Core	Core	Register Register Register	SF	Inter	
						°,	
Core	Core	Core	Core	Register Register	SE	Inec	
Core	Core	Core	Core	Register	01	×	
Core	Core	Core	Core	Register			
Core	Core	Core	Core	Register Register Register	SF		
						<u> </u>	,
	Ins	truct	ion	Cach	ne		
		-1.1.4					
5	are		emo	ry/Ca	ache		

GPU architecture **NVIDIA** microarchitectures



Latest generations: Volta (2017/18), Turing (2019/20) Disclaimer: this is promotional material!

GPU architecture: main components

GPU architecture

GigaThread scheduler							
SM	SM	SM					
Scheduler	Scheduler	Scheduler					
Shared Memory	Shared Memory	Shared Memory					
	Interconnect						
L2							
DRAM							
DIVIN							

08/10/2018

13/36

GPU architecture: computing core

Computing core

- Integer and floating-point units:
 - IEEE-2008 compliant floating point arithmetic (starting from Fermi)
 - Fused multiply-and-add instruction in hardware
- Logic unit
- Move, compare unit
- Branch unit



NVIDIA Geforce GTX Titan

- Kepler generation card
- 2688 streaming processor cores
- 837 MHz clock frequency (1.5 GHz memory clock)
- single precision peak performance 4.4 TFLOP/s
- 1.5 GB GDDR5 RAM
- memory bandwidth 288 GB/s



M. Weigel (Coventry/Mainz)

GPU basics

08/10/2018 16/36

M. Weigel (Coventry/Mainz)

GPU basics

Grid 1

Block (0.0

Execution model

Host (CPU) 08/10/2018 17/36

Device (GPU)

GPU architecture

GPU computation frameworks

 $\mathsf{GPGPU} = \mathsf{General} \; \mathsf{Purpose} \; \mathsf{Computation} \; \mathsf{on} \; \mathsf{Graphics} \; \mathsf{Processing} \; \mathsf{Unit}$

"Old" times: use original graphics primitives

- OpenGL
- DirectX

Vendor specific APIs for GPGPU:

- NVIDIA CUDA: library of functions performing computations on GPU (C, C++, Fortran), additional preprocessor with language extensions
- ATI/AMD Stream: similar functionality for ATI GPUs

Device independent schemes:

- superseded approaches: BrookGPU, AMD Stream, Sh
- OpenCL (Open Computing Language): open framework for parallel programming across a wide range of devices, ranging from CPUs, Cell processors and GPUs to handheld devices
- OpenACC: pragma based parallelization of code parts

Definitions

Kernel GPU program

Thread scalar execution

unit

Warp block of 32

that runs on a

grid of threads

threads executed

executed on the same SM

usually executed

on different SMs

in lockstep

Block a set of warps

Grid a set of blocks

Threading hierarchy

- one kernel is executed at a time (later on modified in concept of streams)
- each kernel executes in many threads, but on one device
- Compare CUDA threads to CPU threads
 - CUDA threads are lightweight
 - very little creation overhead
 - low cost of thread switching
 - CUDA needs thousands of threads for efficiency
- Kernels are executed by an array of threads
 - $\circ\,$ all threads run the same code
 - $\circ\,$ each thread has a unique threadid for control decisions and memory access
- Kernel launched are in grid of thread blocks
 - threads within block cooperate via shared memory
 - threads within a block can synchronize
 - threads within different blocks cannot cooperate (or only via global memory)
 - ${\scriptstyle \circ }$ block executes on one SM and does not migrate
- allows programs to transparently scale to GPUs with different numbers of cores

```
M. Weigel (Coventry/Mainz)
```

GPU basics

CUDA Programming

CUDA variables

Variable declaration	Memory	Scope	Lifetime	Penalty/Latency
int var;	register	thread	thread	1X
<pre>int array_var[10];</pre>	local	thread	thread	100X (pre-Fermi)
shared int shared_var;	shared	block	block	10X
device int global_var;	global	grid	application	100X
constant int constant_var;	constant	grid	application	1X

- automatic scalar variables reside in registers, compiler will spill into local memory in shortage of registers (use -ptxas-options=-v to check)
- automatic array variables (in the absence of qualifiers) reside in thread-local memory
- the type of memory used will be crucial for the performance of the application

08/10/2018 21/36

CUDA code: C with some extra reserved words

Consider a simple "SAXPY" computation, i.e., "Single-Precision $A\cdot X+Y$ ".

Standard C code

```
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];
}
// Invoke serial SAXPY kernel
saxpy_serial(n, 2.0, x, y);</pre>
```

CUDA C code

__global__ void saxpy_parallel(int n, float a, float *x, float *y)
{

int i = blockIdx.x*blockDim.x + threadIdx.x; if (i < n) y[i] = a*x[i] + y[i];</pre>

// Invoke parallel SAXPY kernel with 256 threads/block int nblocks = (n + 255) / 256; saxpy_parallel <<< nblocks, 256>>>(n, 2.0, x, y);

 $\mathsf{M}. \ \mathsf{Weigel} \ \ \mathsf{(Coventry/Mainz)}$

08/10/2018

22/36

CUDA Programming Elementary data transfers with CUDA

,

Memory allocation

Arrays in device global memory are typically allocated from CPU code. Functions:

GPU basics

- cudaMalloc(void ** pointer, size_t nbytes);
- o cudaMemset(void * pointer, int value, size_t count);
- cudaFree(void* pointer);

```
int n = 1024;
int nbytes = 1024*sizeof(int);
int *a_d = 0;
cudaMalloc( (void**)&a_d, nbytes );
cudaMemset( a_d, 0, nbytes);
cudaFree(a_d);
```

Data transfers

The elementary function for data transfers is

- - direction is one of cudaMemcpyHostToDevice, cudaMemcpyDeviceToHost or cudaMemcpyDeviceToDevice and specifies location of src and dst
 - blocks CPU thread (asynchronous transfers possible in streams)

GPU basics

Kernel execution

Function	qua	lifiers
----------	-----	---------

__global__ void f()

- function called from host, executed on device
- must return void

M. Weigel (Coventry/Mainz)

__device__ int f()

• function called from device, executed on device

__host__ int f()

- function called from host, executed on host
- __host__ and __device__ can be combined to generate CPU and GPU code

Built-in variables

All __global__ and __device__ functions have the following automatic variables:

- dim3 gridDim; dimension of the grid in blocks
- dim3 blockDim; dimension of the block in threads
- dim3 blockIdx; block index within grid
- dim3 threadIdx; thread index within block

The indices can be used to construct a global thread index, for instance for a block size of 5 threads,

thread_index = blockIdx.x*blockDim.
 x + threadIdx.x;

08/10/2018 26/36

CUDA Programming

Kernel execution

Execution configuration	Built-in variables
Modified C function call syntax:	Allglobal anddevice functions have the following automatic variables:
 Execution configuration: ag: dimension and size of grid in blocks two-dimensional, ag.x and ag.y 	 dim3 gridDim; — dimension of the grid in blocks dim3 blockDim; — dimension of the block in threads dim3 blockIdx; — block index within
 total number of blocks launched is dG.x × dG.y dB: dimension and size of each block two- or three-dimensional, dB.x, dB.y, and dB.z total number of threads per block is dB.x × dB.y × dB.z if not specified, dB.z = 1 is assumed 	 grid dim3 threadIdx; — thread index within block The indices can be used to construct a global thread index, for instance for a block size of 5 threads, thread_index = blockIdx.x*blockDim.x + threadIdx.x;
M. Weigel (Coventry/Mainz) GPU	J basics 08/10/2018 26 /
Memory hierarchy	
Per device: global memory	

Memory hierarchy Per thread • Registers (extra fast, no copy for ops) Local memory Thread blocks: shared memory Regs Regs Reas allocated by thread block, same lifetime as block allocate as Local Local Local Local Loca __shared__ int shared_array[DIM]; Iow latency (of the order of 10 **Shared Memory** cycles), bandwidth up to 1 TB/s • use for data sharing and user-managed cache

GPU basics

Memory hierarchy



Global Memory

Host Memory

PCIe Bus

Memory hierarchy

Memory hierarchy (summary)

More generally, the different types of memory have the following characteristics:

Memory	Location	Cached	Access	Scope	Lifetime
Register	On-chip	N/A	R/W	One thread	Thread
Local	Off-chip	No	R/W	One thread	Thread
Shared	On-chip	N/A	R/W	All threads in a block	Block
Global	Off-chip	(Yes)	R/W	All threads $+$ host	Application
Constant	Off-chip	Yes	R	All threads $+$ host	Application
Texture	Off-chip	Yes	R	All threads $+$ host	Application

M. Weigel (Coventry/Mainz)

GPU basics

08/10/2018 30/36

Memory hierarchy

Unified virtual addressing



Unified virtual addressing



Memory hierarchy

DirectGPU data transfers



Memory hierarchy

- GPUDirect/CUDA 5.0: direct communication between GPUs in different nodes
- MPI integration under development

CUDA distributrion

CUDA is available together with documentation and libraries on CUDA zone,

https://developer.nvidia.com/cuda-zone

current version is CUDA 10.

Important tools:

- CUDA compiler, nvcc
- CUDA visual profiler
- NVIDIA Nsight
- CUDA gdb
- OpenACC

Additional libraries:

- cuRAND
- \circ cuFFT
- nvGRAPH, ...
- M. Weigel (Coventry/Mainz)

OpenACC

CUDA GDE

OpenACC

The CUDA distribution

Summary and outlook

This lecture

This lecture has given a basic introduction into GPGPU and, in particular, the CUDA framework of GPU programming. Some basic examples hopefully provided a feel for how to go about in using these devices for scientific computing.

GPU basics

Next lecture

In lecture 2, we will start using GPUs for simulating spin models with local algorithms. In terms of GPU programming, a number of additional concepts such as thread synchronization, memory coalescence, and atomic operations will be introduced.

Reading

- Zillions of internet resources, e.g., N. Matloff, "Programming on Parallel Machines", http://heather.cs.ucdavis.edu/~matloff/158/PLN/ParProcBook.pdf
- D. B. Kirk, W.-m. W. Hwu, "Programming Massively Parallel Processors", 3rd edition (Morgan Kaufmann, Amsterdam, 2016).
- J. Sanders, E. Kandrot: "CUDA by example An Introduction to General-Purpose GPU Programming", (Addison Wesley, Upper Saddle River, 2011).

Visual Profile

CUDA MemCher

CLIDA Profiling Tools Interfac

08/10/2018

34 / 36

CUDA distributrion

- Programming Guide
- Best Practices Guide
- Compatibility and Tuning Guides for different microarchitectures (Tesla to Turing)

CUDA compiler invocation:

Simplest case

nvcc -o test test.cu

Some options

M. Weigel (Coventry/Mainz)

GPU basics

08/10/2018 35/36