



BE THE DIFFERENCE.

Credits: Slides adapted from presentations of Sudeep Pasricha and others: Kubiatowicz, Patterson, Mutlu, Elsevier 💦

1

Flynn's Classification (1966)

Broad classification of parallel computing systems

- SISD: Single Instruction, Single Data
 - ° conventional uniprocessor
- SIMD: Single Instruction, Multiple Data
 - $^{\circ}$ one instruction stream, multiple data paths
 - $^\circ~$ distributed memory SIMD
 - ° shared memory SIMD
- MIMD: Multiple Instruction, Multiple Data
 - ° conventional multiprocessors
 - ° message passing machines
 - ° non-cache-coherent shared memory machines
 - ° cache-coherent shared memory machines
- MISD: Multiple Instruction, Single Data
 - ° Not a practical configuration





- •Convergence of application demands and technology constraints drives architecture choice
- •New applications, such as graphics, machine vision, speech recognition, machine learning, etc. - all require large numerical computations that are often trivially data parallel
- •SIMD-based architectures (vector-SIMD, subword-SIMD, SIMT/GPUs) are most efficient way to execute these algorithms







A Shift in the GPU Landscape

- •Some users noticed they could do general-purpose computation by mapping input and output data to images, and computation to vertex and pixel shading computations
- •Referred to as general-purpose computing on graphics processing units (GP-GPU)
- •Incredibly difficult programming model it had to use graphics pipeline model for general computation
 - ° A programming revolution was needed!

General-Purpose GPUs (GP-GPUs)

- In 2006, Nvidia introduced GeForce 8800 GPU supporting a new programming language:
 - ° CUDA "Compute Unified Device Architecture"
 - $^\circ\,$ Subsequently, broader industry pushing for OpenCL, a vendor-neutral version of same ideas.
- Idea: Take advantage of GPU computational performance and memory bandwidth to accelerate some kernels for general-purpose computing
- Attached processor model: host CPU issues data-parallel kernels to GP-GPU for execution







CUDA Revolution!
 CUDA Community Showcase <u>http://www.nvidia.com/object/gpu-applications.html</u> Computational fluid dynamics, EDA, finance, life sciences, signal processing, Speed-up's of >300x for some applications
 GPU Technology Conference <u>http://www.gputechconf.com/page/home.html</u> Include archive of previous editions
 Download CUDA <u>https://developer.nvidia.com/cuda-downloads</u> And start using it!
 NVIDIA YouTube Videos: <u>https://www.youtube.com/user/nvidia/videos</u> Many universities have already courses dedicated to teaching and using CUDA for research









```
#include <stdio.h>
#define SIZEOFARRAY 64
extern void fillArray(int *a, int size); Simple.c

/* The main program */
int main(int argc, char *argv[])
{
    /* Declare the array that will be modified by the GPU */
    int a[SIZEOFARRAY];
    int i;
    /* Initialize the array */
    for(i=0; i < SIZEOFARRAY; i++) {
        a[i]=i;
    }
    /* Print the initial array */
    printf("Initial state of the array:\n");
    for(i = 0; i < SIZEOFARRAY; i++) {
        printf("%d ",a[i]);
    }
    printf("\n");

    /* Call the function that will in turn call the function in
        the GPU that will fill the array */
    fillArray(a,SIZEOFARRAY);

    /* Now print the array after calling fillArray */
    printf("Final state of the array:\n");
    for(i = 0; i < SIZEOFARRAY);

    /* Now print the array after calling fillArray */
    printf("%d ",a[i]);
    }
    printf("\n");
    return 0;
}
</pre>
```

simple.cu

- simple.cu contains two functions
 - 1. fillArray(): A function that will be executed on the host and which takes care of:
 - Allocating variables in the global GPU memory
 - Copying the array from the host to the GPU memory
 - Setting the grid and block sizes
 - Invoking the kernel that is executed on the GPU
 - Copying the values back to the host memory
 - Freeing the GPU memory



```
/* execution configuration... */
fillArray
                /* Indicate the dimension of the block */
                dim3 dimblock(BLOCK_SIZE);
(part 2)
                /* Indicate the dimension of the grid in blocks */
                dim3 dimgrid(arraySize/BLOCK SIZE);
                /* actual computation: Call the kernel, the
                   function that is executed by each and every
                   processing element on the GPU card */
                cu fillArray<<<dimgrid, dimblock>>>(array d);
                /* read results back: */
                /* Copy results from GPU back to memory on the host */
                result = cudaMemcpy(array, array_d, sizeof(int)*arraySize,
                         cudaMemcpyDeviceToHost);
                /* Release the memory on the GPU card */
                cudaFree(array_d);
              }
```

simple.cu (cont.)

- •The other function in simple.cu is
 - 2. cu_fillArray():
 - This is the **kernel** that will be executed in every Stream Processor (SP) in the GPU
 - It is identified as a kernel by the use of the keyword: __global___
 - This function uses the built-in variables
 - blockldx.x
 - threadIdx.x
 - to identify a particular position in the array

```
cu_fillArray
__global__ void cu_fillArray(int *array_d)
{
    int x;
    /* blockIdx.x is a built-in variable in CUDA
    that returns the blockId in the x axis
    of the block that is executing this block of code
    threadIdx.x is another built-in variable in CUDA
    that returns the threadId in the x axis
    of the thread that is being executed by this
    stream processor in this particular block
    */
    x = blockIdx.x*BLOCK_SIZE + threadIdx.x;
    array_d[x] += array_d[x];
}
```



OpenCL – Open Compute Language

CUDA alternative

- Developed by Khronos
 - $^\circ\,$ Industry Consortium that includes: AMD, ARM, Intel, and NVIDIA
- Designed as an open standard for cross-platform parallel programming
- Allows for more general programming across multiple GPUs/CPUs
- Not as popular as CUDA, at least initially...

	OpenCL	CUDA
Programming Language	С	C/C++
Supported GPUs	AMD, NVIDIA	NVIDIA
Supported CPUs	AMD, Intel, ARM	None
Method of Creating GPU Work	Kernel	Kernel
Run-time compilation of kernels	Yes	No
Multiple Kernel Execution	Yes (in certain hardware)	Yes (in certain hardware)
Execution Across Multiple Components	Yes	Yes – only GPUs
Need to Optimize for Best Performance	High	High
Coding Complexity	High	Medium

23

Quick Guide to GPU Terms

A major obstacle to understanding GPUs has been the jargon, with some terms even having misleading names. This obstacle has been surprisingly difficult to overcome.

Туре	More descrip- tive name	Closest old term outside of GPUs	Official CUDA/ NVIDIA GPU term	Book definition
suo	Vectorizable Loop	Vectorizable Loop	Grid	A vectorizable loop, executed on the GPU, made up of one or more Thread Blocks (bodies of vectorized loop) that can execute in parallel.
im abstract	Body of Vectorized Loop	Body of a (Strip-Mined) Vectorized Loop	Thread Block	A vectorized loop executed on a multithreaded SIMD Processor, made up of one or more threads of SIMD instructions. They can communicate via Local Memory.
Progra	Sequence of SIMD Lane Operations	One iteration of a Scalar Loop	CUDA Thread	A vertical cut of a thread of SIMD instructions corresponding to one element executed by one SIMD Lane. Result is stored depending on mask and predicate register.
ine object	A Thread of SIMD Instructions	Thread of Vector Instructions	Warp	A traditional thread, but it contains just SIMD instructions that are executed on a multithreaded SIMD Processor. Results stored depending on a per-element mask.
Mach	SIMD Instruction	Vector Instruction	PTX Instruction	A single SIMD instruction executed across SIMD Lanes.
Í	Multithreaded SIMD Processor	(Multithreaded) Vector Processor	Streaming Multiprocessor	A multithreaded SIMD Processor executes threads of SIMD instructions, independent of other SIMD Processors.
hardware	Thread Block Scheduler	Scalar Processor	Giga Thread Engine	Assigns multiple Thread Blocks (bodies of vectorized loop) to multithreaded SIMD Processors.
rocessing	SIMD Thread Scheduler	Thread scheduler in a Multithreaded CPU	Warp Scheduler	Hardware unit that schedules and issues threads of SIMD instructions when they are ready to execute; includes a scoreboard to track SIMD Thread execution.
-	SIMD Lane	Vector Lane	Thread Processor	A SIMD Lane executes the operations in a thread of SIMD instructions on a single element. Results stored depending on mask.
g	GPU Memory	Main Memory	Global Memory	DRAM memory accessible by all multithreaded SIMD Processors in a GPU.
ardwa	Private Memory	Stack or Thread Local Storage (OS)	Local Memory	Portion of DRAM memory private to each SIMD Lane.
mory	Local Memory	Local Memory	Shared Memory	Fast local SRAM for one multithreaded SIMD Processor, unavailable to other SIMD Processors.
Me	SIMD Lane Registers	Vector Lane Registers	Thread Processor Registers	Registers in a single SIMD Lane allocated across a full thread block (body of vectorized loop).



Thread Batching: Grids and Blocks

- Kernel executed as a grid of thread blocks
 - All threads share data memory space
- Thread Block is a batch of threads, can cooperate with each other by:
 - Synchronizing their execution: For hazard-free shared memory accesses
 - Efficiently sharing data through a low latency shared memory
- Two threads from two different blocks cannot cooperate
 - (Unless thru slow global memory)
- Threads and Blocks have IDs









CUDA Thread Scheduling

•GPU hardware has two levels of hardware schedulers:

- 1) Thread Block Scheduler (top level) that assigns Thread Blocks to multithreaded SIMD processors, which ensures that thread blocks are assigned to the processors whose local memories have corresponding data
- 2) SIMD Thread Scheduler (lower level) (warp scheduler) within a SIMD Processor, which schedules when threads of SIMD instructions should run







CUDA Device Memory Space Overview

- Each Thread can:
 - 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
- Host can R/W global, constant, and texture memories



33

Global, Constant, and Texture Memories (Long Latency Accesses)





	TESLA K10 ^a	TESLA K20	TESLA K20X	
Peak double precision floating point performance (board)	0.19 teraflops	1.17 teraflops	1.31 teraflops	Evampla
Peak single precision floating point performance (board)	4.58 teraflops	3.52 teraflops	3.95 teraflops	Example:
Number of GPUs	2 x GK104s	1 x GK110		
Number of CUDA cores	2 x 1536	2496	2688	Invidia lesia k20
Memory size per board (GDDR5)	8 GB	5 GB	6 GB	
Memory bandwidth for board (ECC off)*	320 GBytes/sec	208 GBytes/sec	250 GBytes/sec	
SPU computing applications	Seismic, image, signal processing, video analytics	CFD, CAE, financial computin and physics, data analytics, s modeling	g, computational chemistry atellite imaging, weather]
Architecture features	SMX	SMX, Dynamic Parallelism, H	Graphics Pro	cessor Graphics Card Clock Speeds
System	Servers only	Servers and Workstations	GPU Name: GK	110 Released: Nov 12th, 2012 GPU Clock: 706 MHz
			Process Size: 28 Transistors: 7,0 Die Size: 561	Imm Production Status Active Memory Clock: 1300 MHz 5200 MHz 1300 MHz effect 30 million Launch Price: 3,199 USD Memory Memory mm* Bus Interface: PCIe 2.0 x16 Memory Memory Status
	ESIA		Render Co Shading Units: 249	Reference Board Memory Type: GDDR5 6 Clawlitht, Doddate Memory Type: 200 bit
- arman	0		TMUs: 208	Siot width: Dual-siot memory Bus: 320 bit
			ROPs: 40	Length: 267 mm
A DECEMBER OF THE OWNER	a and		SMX Count: 13	TDP: 225 W Graphics Features
16			Pixel Rate: 36.	7 GPixel/s VGA BIOS DirectX: 11.0
	com/content/	m/content/PDF/kepler/te	Tauture Data: 147	OpenGL: 4.4
https://www.nvidia.	comy conterney		Texture rate: 147	Glexel/s Eind graphics card BIOS for this card















	Unified Memory
Þ	Simpler Programming and Memory Model
•	Performance Through Data Locality
	<pre>- cudaMallocManaged();</pre>
	// Allocate input
I	malloc(input,);
•	cudaMallocManaged(d_input,);
•	amopy (a_input, input,,, // copy to managet memory
1	// Allocate output
•	cudaMallocManaged(d_output,);
	// Launch GPU kernel
Ģ	<pre>gpu_kernel<<<blocks, threads="">>> (d_output, d_input,);</blocks,></pre>
	// Synchronize
	cudaDeviceSynchronize();

Collaborative Computing Algorithms

- Case studies using CPU and GPU
- Kernel launches are asynchronous
 - ° CPU can work while waits for GPU to finish
 - Traditionally, this is the most efficient way to exploit heterogeneity
- Fine-grain heterogeneity becomes possible with Pascal/Volta architecture
- Pascal/Volta Unified Memory
 - ° CPU-GPU memory coherence
 - ° System-wide atomic operations
- Benefits of Collaboration Example: Bézier Surfaces
 - ° [1] J. Gomez-Luna et. al, Chai: Collaborative Heterogeneous Applications for Integrated-architectures, ISPASS 2017.
 - ° Data partitioning improves performance
 - AMD Kaveri (4 CPU cores + 8 GPU CUs)



Bézier Surfaces (up to 47% improvement over GPU only)

GPUs for Mobile Clients and Servers

Goal is for the graphics quality of a movie such as *Avatar* to be achieved in real time on a server GPU in 2015 and on your mobile GPU in 2020

	NVIDIA Tegra 2	NVIDIA Fermi GTX 480
Market	Mobile client	Desktop, server
System processor	Dual-Core ARM Cortex-A9	Not applicable
System interface	Not applicable	PCI Express 2.0 × 16
System interface bandwidth	Not applicable	6 GBytes/sec (each direction), 12 GBytes/sec (total)
Clock rate	Up to 1 GHz	1.4 GHz
SIMD multiprocessors	Unavailable	15
SIMD lanes/SIMD multiprocessor	Unavailable	32
Memory interface	32-bit LP-DDR2/DDR2	384-bit GDDR5
Memory bandwidth	2.7 GBytes/sec	177 GBytes/sec
Memory capacity	1 GByte	1.5 GBytes
Transistors	242 M	3030 M
Process	40 nm TSMC process G	40 nm TSMC process G
Die area	57 mm ²	520 mm ²
Power	1.5 watts	167 watts

Comparison of a GPU and a MIMD with Multimedia SIMD

Purpose is *not* to determine how much faster one product is than another, but to try to understand the relative value of features of these two contrasting architecture styles

	Core i7- 960	GTX 280	GTX 480	Ratio 280/i7	Ratio 480/i7
Number of processing elements (cores or SMs)	4	30	15	7.5	3.8
Clock frequency (GHz)	3.2	1.3	1.4	0.41	0.44
Die size	263	576	520	2.2	2.0
Technology	Intel 45 nm	TSMC 65 nm	TSMC 40 nm	1.6	1.0
Power (chip, not module)	130	130	167	1.0	1.3
Transistors	700 M	1400 M	3030 M	2.0	4.4
Memory bandwidth (GBytes/sec)	32	141	177	4.4	5.5
Single-precision SIMD width	4	8	32	2.0	8.0
Double-precision SIMD width	2	1	16	0.5	8.0
Peak single-precision scalar FLOPS (GFLOP/Sec)	26	117	63	4.6	2.5
Peak single-precision SIMD FLOPS (GFLOP/Sec)	102	311 to 933	515 or 1344	3.0-9.1	6.6-13.1
(SP 1 add or multiply)	N.A.	(311)	(515)	(3.0)	(6.6)
(SP 1 instruction fused multiply-adds)	N.A.	(622)	(1344)	(6.1)	(13.1)
(Rare SP dual issue fused multiply-add and multiply)	N.A.	(933)	N.A.	(9.1)	
Peak double-precision SIMD FLOPS (GFLOP/sec)	51	78	515	1.5	10.1

Relative Performance

Kernel	Units	Core i7-960	GTX 280	GTX 280/ i7-960
SGEMM	GFLOP/sec	94	364	3.9
MC	Billion paths/sec	0.8	1.4	1.8
Conv	Million pixels/sec	1250	3500	2.8
FFT	GFLOP/sec	71.4	213	3.0
SAXPY	GBytes/sec	16.8	88.8	5.3
LBM	Million lookups/sec	85	426	5.0
Solv	Frames/sec	103	52	0.5
SpMV	GFLOP/sec	4.9	9.1	1.9
GJK	Frames/sec	67	1020	15.2
Sort	Million elements/sec	250	198	0.8
RC	Frames/sec	5	8.1	1.6
Search	Million queries/sec	50	90	1.8
Hist	Million pixels/sec	1517	2583	1.7
Bilat	Million pixels/sec	83	475	5.7

Reasons for Differences from Intel
• GPU has 4.4× the memory bandwidth ° Explains why LBM and SAXPY run 5.0 and 5.3× faster; their working sets are hundreds of megabytes and hence don't fit into the Core i7 cache
 Five of the remaining Kernels are compute-bound: SGEMM, Conv, FFT, MC, and Bilat GTX 280 single precision is 3 to 6× faster; DP performance is only 1.5× faster; has direct support for transcendental functions lacking in i7
• Cache blocking optimizations benefit i7 ° Convert RC, Search, Sort, SGEMM, FFT, and SpMV from memory-bound to compute-bound
 Multimedia SIMD extensions are of little help if the data are scattered throughout main memory Reinforces the importance of gather-scatter to vector and GPU architectures that is missing from SIMD extensions
47

