




GPU 101



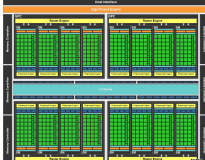
Oregon State University
Mike Bailey
mjb@cs.oregonstate.edu



This work is licensed under a [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nc-nd/4.0/).



Oregon State University
Computer Graphics

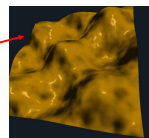
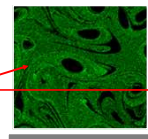



gp101.pdf mjb - March 28, 2023

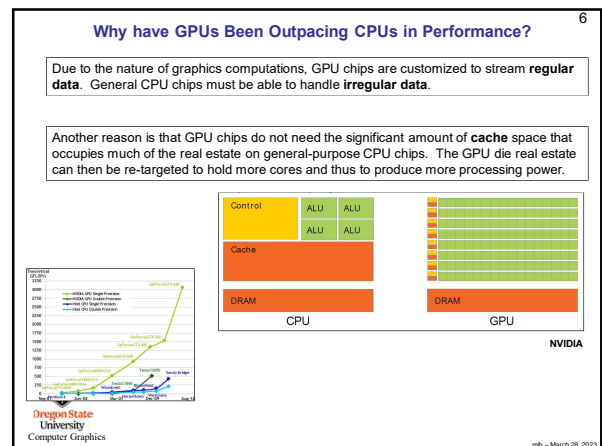
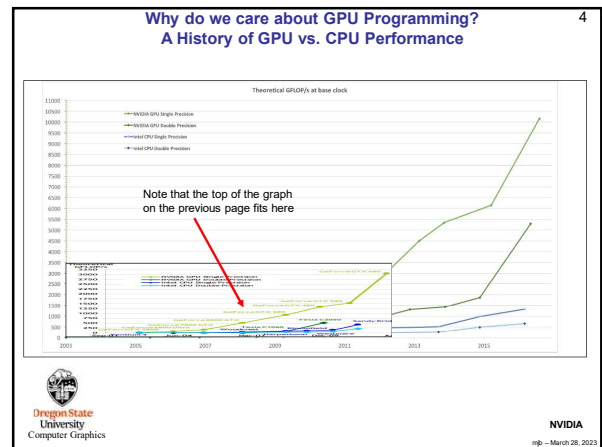
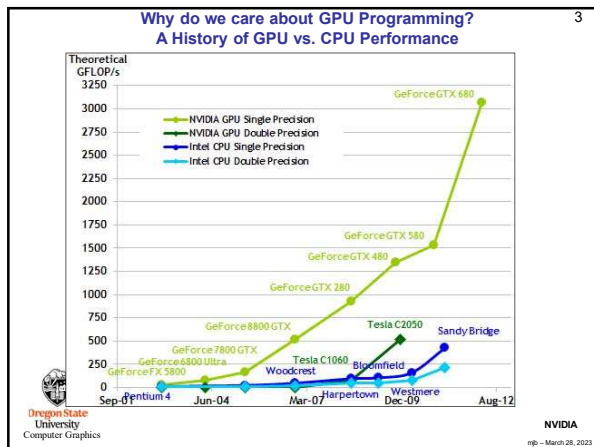
How Have You Been Able to Gain Access to GPU Power?

There have been three ways:

1. Write a graphics display program (≥ 1985)
2. Write an application that looks like a graphics display program, but uses the fragment shader to do some per-node computation (≥ 2002)
3. Write in OpenCL or CUDA, which looks like C++ (≥ 2006)

mjb - March 28, 2023



Why have GPUs Been Outpacing CPUs in Performance?

Another reason is that general CPU chips contain on-chip logic to do **branch prediction** and **out-of-order execution**. This, too, takes up chip die space.

But, CPU chips can handle more general-purpose computing tasks.

So, which is better, a CPU or a GPU?

It depends on what you are trying to do!

Oregon State University
Computer Graphics

mjb - March 28, 2023

Originally, GPU Devices were very task-specific

Oregon State University
Computer Graphics

mjb - March 28, 2023

Today's GPU Devices are not task-specific

Oregon State University
Computer Graphics

mjb - March 28, 2023

Consider the architecture of the NVIDIA Tesla V100's that we have in our DGX System

Oregon State University
Computer Graphics

mjb - March 28, 2023

What is a "Core" in the GPU Sense?

Look closely, and you'll see that NVIDIA really calls these "CUDA Cores"

Look even more closely and you'll see that these CUDA Cores have no control logic – they are **pure compute units**. (The surrounding SM has the control logic.)

Other vendors refer to these as "Lanes". You might also think of them as 64-way SIMD.

Oregon State University
Computer Graphics

mjb - March 28, 2023

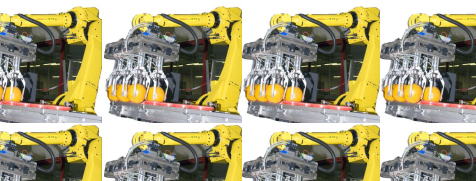
A Mechanical Equivalent...

Oregon State University
Computer Graphics

<http://news.cision.com>

mjb - March 28, 2023

How Many Robots Do You See Here?



12? 72? Depends what you count as a "robot".

Oregon State University
Computer Graphics


© 2020

14

A Spec Sheet Example

NVIDIA Card		Number of CUDA Cores	Size of Power Supply *	Memory Type	Memory Interface	Memory Bandwidth GB/sec	Base Clock Speed	Boost Clock Speed	NOTES
RTX-4090		9728	750 watt	GGDDR6	256 bit	719.5 GB/s	2.21 GHz	2.51 GHz	16 GB of Memory
RTX-4080		16384	850 watt	GGDDR6	256 bit	1008 GB/s	2.23 GHz	2.52 GHz	24 GB of Memory
NVIDIA Card		Number of CUDA Cores	Size of Power Supply *	Memory Type	Memory Interface	Memory Bandwidth GB/sec	Base Clock Speed	Boost Clock Speed	NOTES
3000 Series									
RTX-3050		2560	550 watt	GGDR6	128 bit	224 GB/s	1550 MHz	1670 MHz	Standard with 8 GB of Memory
RTX-3060		3684	650 watt	GGDR6	192 bit	384 GB/s	1500 MHz	1780 MHz	Standard with 12 GB of Memory
RTX-3070		4608	600 watt	GGDR6	256 bit	448 GB/s	1410 MHz	1760 MHz	Standard with 8 GB of Memory
RTX-3070		5688	650 watt	GGDR6	256 bit	448 GB/s	1580 MHz	1770 MHz	Standard with 8 GB of Memory
RTX-3080		6144	750 watt	GGDR6	256 bit	608 GB/s	1500 MHz	1780 MHz	Standard with 8 GB of Memory
RTX-3090		8704	750 watt	GGDR6	320 bit	760 GB/s	1440 MHz	1710 MHz	Standard with 10 GB of Memory
RTX-3090 Ti		10240	750 watt	GGDR6	384 bit	812 GB/s	1370 MHz	1670 MHz	Standard with 12 GB of Memory
RTX-3090		10496	750 watt	GGDR6	384 bit	836 GB/s	1400 MHz	1700 MHz	Standard with 24 GB of Memory
RTX-3090 Ti		10572	850 watt	GGDR6	384 bit	936 GB/s	1670 MHz	1960 MHz	Standard with 24 GB of Memory
NVIDIA Card		Number of CUDA Cores	Size of Power Supply *	Memory Type	Memory Interface	Memory Bandwidth GB/sec	Base Clock Speed	Boost Clock Speed	NOTES
2000 Series									
RTX-2060		1920	500 watt	GGDR6	192 bit	336 GB/s	1365 MHz	1660 MHz	Standard with 6 GB of Memory
RTX-2060 Super		2176	600 watt	GGDR6	192 bit	384 GB/s	1470 MHz	1740 MHz	Standard with 8 GB of Memory
RTX-2070		2304	550 watt	GGDR6	256 bit	448 GB/s	1410 MHz	1620 MHz	Standard with 8 GB of Memory
RTX-2070 Super		2560	600 watt	GGDR6	256 bit	496 GB/s	1500 MHz	1710 MHz	Standard with 8 GB of Memory
RTX-2080		2944	650 watt	GGDR6	256 bit	448 GB/s	1515 MHz	1710 MHz	Standard with 8 GB of Memory
RTX-2080 Super		3072	650 watt	GGDR6	256 bit	496 GB/s	1560 MHz	1915 MHz	Standard with 8 GB of Memory
RTX-2080 Ti		4352	650 watt	GGDR6	352 bit	616 GB/s	1350 MHz	1545 MHz	Standard with 11 GB of Memory
Titan RTX		4608	650 watt	GGDR6	384 bit	672 GB/s	1350 MHz	1770 MHz	Standard with 24 GB of Memory

NVIDIA

 Oregon State University
Computer Graphics

rsb - March 28, 2025

NVIDIA 4090 Spec Sheet		
Graphics Processor	Graphics Card	Relative Performance
GPU Name: ADA4090	Release Date: Sep 2024, 2022	GeForce RTX 4090 100%
GPU Variant: RTX 4090-A1	Availability: Oct 13th, 2022	GeForce RTX 4090 99%
Architecture: Ada Lovelace	Generation: GeForce 40	Radeon RX 8900 67%
Foundry: TSMC	Production: GeForce 30	GeForce RTX 4070 65%
Process Size: 4 nm	Production: Active	GeForce RTX 4060 55%
Transistors: 76,200 mm²	Launch Price: 1,599 USD	GeForce RTX 3090 39%
Memory: 12GB 384 bit	Current Price: \$1,520 / \$1,550	GeForce RTX 4050 38%
Bit Size: 384 bit	PCI-E: 16 Lanes	Radeon RX 7900 37%
	Refused: Multi-use database	GeForce RTX 4040 100%
Clock Speeds	Theoretical Performance	Memory
Base Clock: 2235 MHz	Pixel Rate: 441.2 GPixel/s	Memory Size: 24 GB
Boost Clock: 2520 MHz	Texture Rate: 1,200 GTexture/s	Memory Type: GDDR6X
Memory Clock: 1935 MHz	FP32 Rate: 15.76 TFLOPS (0.1)	Memory Bus: 384 bit
21 Gbps Effective	FP16 Rate: 31.52 TFLOPS	Bandwidth: 1,008 GB/s
	FP16 (double): 1,206 GFLOPS (0.68)	
Board Design	Graphics Features	Render Clocks
Slot Width: Triple slot	Stream: 4.6	Shader Units: 14880
Length: 306 mm	RayTracing: 5.0	TMN: 512
121 mm	Optimus: 1.3	ROPs: 176
Height: 69 mm	VRAM: 1.3	SM Count: 128
3.4 inches	CUDA: 8.9	Tensor Core: 512
Port: 40W	Shader Model: 6.7	RT Count: 128
Suggested Price: \$550		L1 Cache: 128 KB per SM
		L2 Cache: 72 MB
Outputs: 1x HDMI 2.1		
3x DisplayPort 1.4a		
Power Connection: 1x 16pin		
Board Number: RTX 4090 A1 380		

NVIDIA

np - March 28, 2025

The Bottom Line is This	
It is obvious that it is difficult to <i>directly</i> compare a CPU with a GPU. They are optimized to do different things.	
So, let's use the information about the architecture as a way to consider what CPUs should be good at and what GPUs should be good at	
<u>CPU</u>	<u>GPU</u>
General purpose programming	Data parallel programming
Multi-core under user control	Little user control
Irregular data structures	Regular data structures
Irregular flow control	Regular Flow Control

The diagram illustrates the hierarchy of GPU architecture in four stages:


- Platform:** A large yellow rectangle containing two green rectangles labeled "Device #1" and "Device #2".
- Device:** A zoomed-in view of a single green rectangle, showing a 2x3 grid of yellow rectangles labeled "CU" (Compute Unit).
- Compute Unit:** A zoomed-in view of a single yellow rectangle, showing a 3x3 grid of orange rectangles labeled "PE" (Processing Element).
- Processing Element:** A zoomed-in view of a single orange rectangle, showing its internal structure.

Dotted lines indicate the zooming sequence from Platform to Device, then to Compute Unit, and finally to Processing Element.

Thinking ahead to CUDA and OpenCL...

How can GPUs execute General C Code Efficiently?

- Ask them to do what they do best. Unless you have a very intense **Data Parallel** application, don't even think about using GPUs for computing.
- GPU programs expect you to not just have a few threads, but to have **thousands** of them!
- Each thread executes the same program (called the **kernel**), but operates on a different small piece of the overall data
- Thus, you have many, many threads, all waking up at about the same time, all executing the same kernel program, all hoping to work on a small piece of the overall problem.
- CUDA and OpenCL have built-in functions so that each thread can figure out which thread number it is, and thus can figure out what part of the overall job it's supposed to do.
- When a thread gets blocked somehow (a memory access, waiting for information from another thread, etc.), the processor switches to executing another thread to work on.

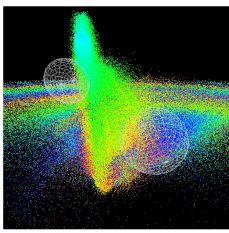
 mp - March 28, 2023

So, the Trick is to Break your Problem into Many, Many Small Pieces


Particle Systems are a great example.

1. Have one thread per *each particle*.
2. Put all of the initial parameters into an array in GPU memory.
3. Tell each thread what the current **Time** is.
4. Each thread then computes its particle's position, color, etc. and writes it into arrays in GPU memory.
5. The CPU program then initiates OpenGL drawing of the information in those arrays.

Note: once setup, the data never leaves GPU memory!



Ben Weiss

 mp - March 28, 2023

Something New – Tensor Cores



NVIDIA


 mp - March 28, 2023

Tensor Cores Accelerate Fused-Multiply-Add Arithmetic

$$D = \begin{pmatrix} A_{11} & A_{12} & A_{13} & A_{14} \\ A_{21} & A_{22} & A_{23} & A_{24} \\ A_{31} & A_{32} & A_{33} & A_{34} \\ A_{41} & A_{42} & A_{43} & A_{44} \end{pmatrix} \begin{pmatrix} B_{11} & B_{12} & B_{13} & B_{14} \\ B_{21} & B_{22} & B_{23} & B_{24} \\ B_{31} & B_{32} & B_{33} & B_{34} \\ B_{41} & B_{42} & B_{43} & B_{44} \end{pmatrix} + \begin{pmatrix} C_{11} & C_{12} & C_{13} & C_{14} \\ C_{21} & C_{22} & C_{23} & C_{24} \\ C_{31} & C_{32} & C_{33} & C_{34} \\ C_{41} & C_{42} & C_{43} & C_{44} \end{pmatrix}$$


FP16 or FP32 FP16 FP16 FP16 or FP32

cuBLAS Mixed-Precision GEMM (FP16 Input, FP32 Compute)



Matrix Size (M=N=K)

■ Tesla P100 ■ Tesla V1 ■ Tensor Cores

 mp - March 28, 2023

What is Fused Multiply-Add?

Many scientific and engineering computations take the form:

$$D = A + (B \cdot C);$$

A "normal" multiply-add would likely handle this as:

```
tmp = B * C;
D = A + tmp;
```


A "fused" multiply-add does it all at once, that is, when the low-order bits of A are ready, they are immediately added into the low-order bits of A at the same time the higher-order bits of B * C are being multiplied.

Consider a Base 10 example: $789 + (123 \cdot 456)$

123
x 456
738
615
492
+ 789
56,877

Can start adding the 9 the moment the 8 is produced!

Note: "Normal" $A + (B \cdot C) \neq \text{"FMA"} A + (B \cdot C)$

 mp - March 28, 2023

There are Two Approaches to Combining CPU and GPU Programs

1. Combine both the CPU and GPU code in the same file. The CPU compiler compiles its part of that file. The GPU compiler compiles just its part of that file.
2. Have two separate programs: a .cpp and a .somethingelse that get compiled separately.

Advantages of Each


1. The CPU and GPU sections of the code know about each others' intents. Also, they can share common structs, #define's, etc.
2. It's potentially cleaner to look at each section by itself. Also, the GPU code can be easily used in combination with other CPU programs.

Who are we Talking About Here?

1 = NVIDIA's CUDA

2 = Khronos's OpenCL

We will talk about each of these separately – stay tuned!


 mp - March 28, 2023


Looking ahead:
If threads all execute the same program,
what happens on flow divergence?

```


if( a > b )
    Do This;
else
    Do That;
  
```

1. The line "if(a > b)" creates a vector of Boolean values giving the results of the if-statement for each thread. This becomes a "mask".
2. Then, the GPU executes all parts of the divergence:
Do This;
Do That;
3. During that execution, anytime a value wants to be stored, the mask is consulted and the storage only happens if that thread's location in the mask is the right value.

 mp - March 28, 2023





- GPUs were originally designed for the streaming-ness of computer graphics
- Now, GPUs are also used for the streaming-ness of data-parallel computing
- GPUs are better for some things. CPUs are better for others.

 mp - March 28, 2023

Dismantling a Graphics Card


This is an Nvidia 1080 ti card – one that died on us. It willed its body to education.




 mp - March 28, 2023

Dismantling a Graphics Card

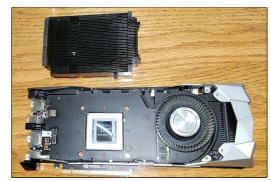

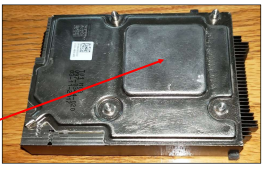
Removing the covers:




 mp - March 28, 2023

Dismantling a Graphics Card

Removing the heat sink:

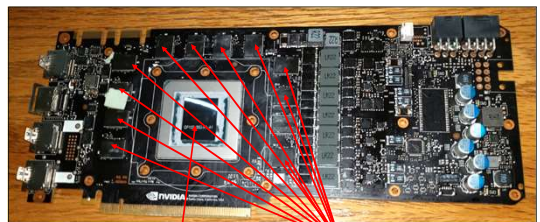




This transfers heat from the GPU Chip to the cooling fins


 mp - March 28, 2023

Dismantling a Graphics Card

Removing the fan assembly reveals the board:



GPU Chip Memory

 mp - March 28, 2023

Dismantling a Graphics Card 31

Power half of the board:

Power distribution **Power input**

mp - March 28, 2023

Dismantling a Graphics Card 32

Graphics half of the board:

Video out **GPU Chip**

This one contains 7.2 billion transistors!
The newer cards contain 70+ billion transistors.
(Thank you, Moore's Law)

mp - March 28, 2023

Dismantling a Graphics Card 33

Underside of the board:

mp - March 28, 2023

Dismantling a Graphics Card 34

Underside of where the GPU chip attaches:

Here is a fun video of someone explaining the different parts of this same card:
<https://www.youtube.com/watch?v=dSCNf9DIBGE>

mp - March 28, 2023