




# 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





Computer Graphics

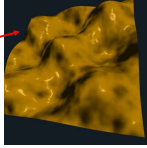
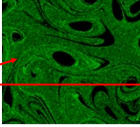
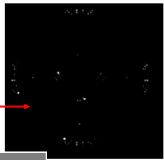

gpu101.pdf

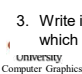
mb - March 15, 2024

1

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

There have been three ways:

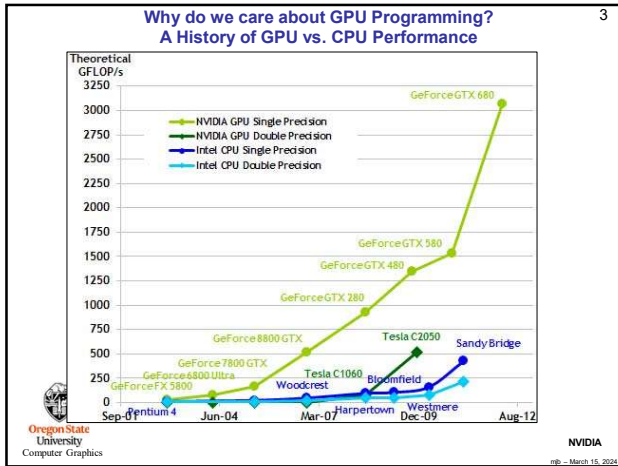
1. Write a graphics display program ( $\geq 1985$ ) 
2. Write an application that looks like a graphics display program, but uses the fragment shader to do some per-node computation ( $\geq 2002$ )  
3. Write in OpenCL or CUDA, which looks like C++ ( $\geq 2006$ ) 



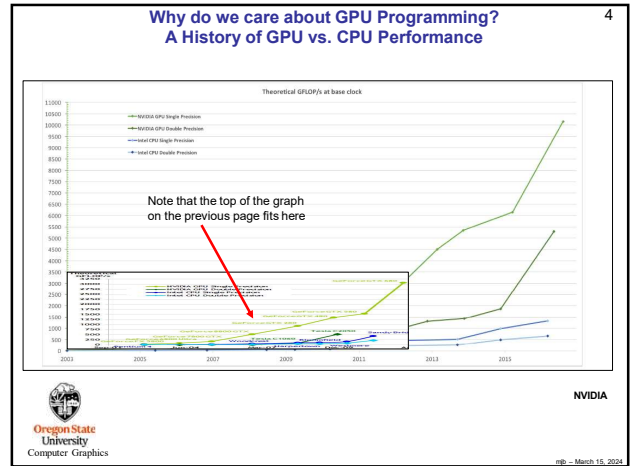
Computer Graphics

mb - March 15, 2024

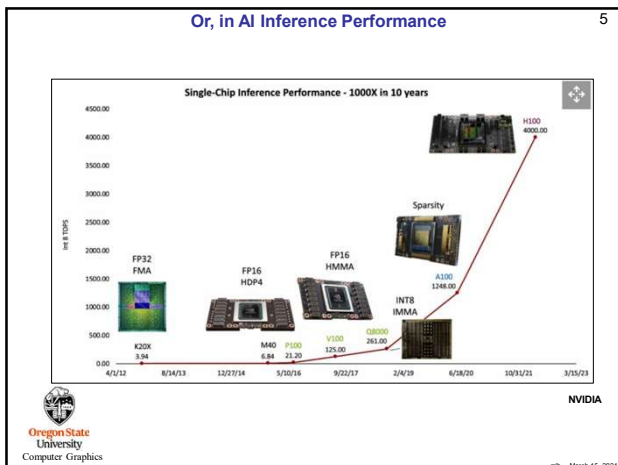
2



3



4



5

### The "Core-Score". How can this be?



Computer Graphics

NVIDIA

mb - March 15, 2024

6

### Why have GPUs Been Outpacing CPUs in Performance?

Due to the nature of graphics computations, GPU chips are customized to stream **regular data**. General CPU chips must be able to handle **irregular data**.

Another reason is that GPU chips do not need the significant amount of **cache** space that occupies much of the real estate on general-purpose CPU chips. The GPU die real estate can then be re-targeted to hold more cores and thus to produce more processing power.

**CPU** vs **GPU** (NVIDIA)

Oregon State University Computer Graphics  
mb - March 15, 2024

7

### 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  
mb - March 15, 2024

8

### Originally, Parts of GPU Chips were very Task-specific

Oregon State University Computer Graphics  
mb - March 15, 2024

9

### Today's GPU Devices are not Task-specific – They Can Be Dynamically Re-purposed for any GPU Function

Oregon State University Computer Graphics  
mb - March 15, 2024

10

### Consider the architecture of the NVIDIA 4090

**128 Streaming Multiprocessors (SMs) / chip**  
**128 cores / SM**  
**Wow! 16,384 cores / chip? Really?**

Oregon State University Computer Graphics  
mb - March 15, 2024

11

### 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  
mb - March 15, 2024

12

### A Mechanical Equivalent...

“Streaming Multiprocessor”

“CUDA Cores”

“Data”

University  
Computer Graphics

<http://news.cision.com>

mp - March 15, 2024

13

### A Spec Sheet Example

| NVIDIA Card 4000 Series | Number of CUDA Cores | Size of Power Supply** | Memory Type | Memory Interface Width | Memory Bandwidth GB/sec | Base Clock Speed | Boost Clock Speed | NOTES           |
|-------------------------|----------------------|------------------------|-------------|------------------------|-------------------------|------------------|-------------------|-----------------|
| RTX 4080                | 9728                 | 750 watt               | GDDR6X      | 256 bit                | 736.8 GB/s              | 2.34 GHz         | 2.51 GHz          | 16 GB of Memory |
| RTX 4090                | 16384                | 850 watt               | GDDR6X      | 384 bit                | 1008 GB/s               | 2.23 GHz         | 2.52 GHz          | 24 GB of Memory |

| NVIDIA Card 3000 Series | Number of CUDA Cores | Size of Power Supply** | Memory Type | Memory Interface Width | Memory Bandwidth GB/sec | Base Clock Speed | Boost Clock Speed | NOTES                         |
|-------------------------|----------------------|------------------------|-------------|------------------------|-------------------------|------------------|-------------------|-------------------------------|
| RTX-3050                | 2560                 | 550 watt               | GDDR6       | 128 bit                | 224 GB/s                | 1050 MHz         | 1780 MHz          | Standard with 8 GB of Memory  |
| RTX-3060                | 3584                 | 500 watt               | GDDR6       | 192 bit                | 384 GB/s                | 1320 MHz         | 1780 MHz          | Standard with 12 GB of Memory |
| RTX-3060 Ti             | 4864                 | 600 watt               | GDDR6       | 256 bit                | 448 GB/s                | 1410 MHz         | 1670 MHz          | Standard with 8 GB of Memory  |
| RTX-3070                | 5888                 | 650 watt               | GDDR6       | 256 bit                | 448 GB/s                | 1680 MHz         | 1770 MHz          | Standard with 8 GB of Memory  |
| RTX-3070 Ti             | 6144                 | 750 watt               | GDDR6X      | 256 bit                | 608 GB/s                | 1500 MHz         | 1730 MHz          | Standard with 8 GB of Memory  |
| RTX-3080                | 8704                 | 750 watt               | GDDR6X      | 320 bit                | 760 GB/s                | 1440 MHz         | 1710 MHz          | Standard with 10 GB of Memory |
| RTX-3080 Ti             | 10240                | 750 watt               | GDDR6X      | 384 bit                | 912 GB/s                | 1370 MHz         | 1670 MHz          | Standard with 12 GB of Memory |
| RTX-3090                | 10496                | 750 watt               | GDDR6X      | 384 bit                | 936 GB/s                | 1400 MHz         | 1700 MHz          | Standard with 24 GB of Memory |
| RTX-3090 Ti             | 10572                | 850 watt               | GDDR6X      | 384 bit                | 936 GB/s                | 1670 MHz         | 1960 MHz          | Standard with 24 GB of Memory |

| NVIDIA Card 2000 Series | Number of CUDA Cores | Size of Power Supply** | Memory Type | Memory Interface Width | Memory Bandwidth GB/sec | Base Clock Speed | Boost Clock Speed | NOTES                         |
|-------------------------|----------------------|------------------------|-------------|------------------------|-------------------------|------------------|-------------------|-------------------------------|
| RTX-2050                | 1920                 | 300 watt               | GDDR6       | 192 bit                | 336 GB/s                | 1350 MHz         | 1680 MHz          | Standard with 6 GB of Memory  |
| RTX-2060 Super          | 2176                 | 550 watt               | GDDR6       | 256 bit                | 448 GB/s                | 1470 MHz         | 1650 MHz          | Standard with 8 GB of Memory  |
| RTX-2070                | 2304                 | 550 watt               | GDDR6       | 256 bit                | 448 GB/s                | 1410 MHz         | 1620 MHz          | Standard with 8 GB of Memory  |
| RTX-2070 Super          | 2560                 | 650 watt               | GDDR6       | 256 bit                | 448 GB/s                | 1605 MHz         | 1770 MHz          | Standard with 8 GB of Memory  |
| RTX-2080                | 2944                 | 650 watt               | GDDR6       | 256 bit                | 448 GB/s                | 1615 MHz         | 1710 MHz          | Standard with 8 GB of Memory  |
| RTX-2080 Super          | 3072                 | 650 watt               | GDDR6       | 256 bit                | 496 GB/s                | 1650 MHz         | 1815 MHz          | Standard with 8 GB of Memory  |
| RTX-2080 Ti             | 4352                 | 650 watt               | GDDR6       | 352 bit                | 616 GB/s                | 1300 MHz         | 1545 MHz          | Standard with 11 GB of Memory |
| Titan RTX               | 4608                 | 650 watt               | GDDR6       | 384 bit                | 672 GB/s                | 1350 MHz         | 1770 MHz          | Standard with 24 GB of Memory |

University  
Computer Graphics

NVIDIA

mp - March 15, 2024

14

### Another Spec Sheet Example

| Graphics Card               | RTX-4070  | RTX-4080  | RTX-4070 Ti |
|-----------------------------|-----------|-----------|-------------|
| Architecture                | AD104     | AD103     | AD104       |
| TSMC 4N                     | TSMC 4N   | TSMC 4N   | TSMC 4N     |
| Transistors (Billion)       | 32        | 45.9      | 35.8        |
| Die size (mm <sup>2</sup> ) | 294.5     | 378.6     | 294.5       |
| SMs                         | 46        | 76        | 60          |
| GPU Cores (Shaders)         | 5888      | 9728      | 7680        |
| Tensor Cores                | 184       | 304       | 240         |
| Ray Tracing "Cores"         | 46        | 76        | 60          |
| Boost Clock (MHz)           | 2475      | 2505      | 2610        |
| VRAM Speed (Gbps)           | 21        | 22.4      | 21          |
| VRAM (GB)                   | 12        | 16        | 12          |
| VRAM Bus Width              | 192       | 256       | 192         |
| L2 Cache (MB)               | 36        | 64        | 48          |
| SDPs                        | 64        | 112       | 80          |
| TMUs                        | 184       | 304       | 240         |
| TFLOPS FP32 (Boost)         | 29.1      | 48.7      | 40.1        |
| TFLOPS FP16 (FP8)           | 233 (466) | 396 (790) | 321 (641)   |
| Bandwidth (GB/s)            | 504       | 717       | 504         |
| TGP (watts)                 | 200       | 320       | 285         |
| Launch Date                 | Nov 2022  | Nov 2022  | Jan 2023    |
| Launch Price                | \$599     | \$1,199   | \$799       |

Tom's Hardware

University  
Computer Graphics

mp - March 15, 2024

15

### NVIDIA's Ampere Chip

University  
Computer Graphics

mp - March 15, 2024

16

### The Bottom Line is This

It is difficult to *directly* 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

|   |   |
|---|---|
| <b>CPU</b>  | <b>GPU</b>  |
| <ul style="list-style-type: none"> <li>General purpose programming</li> <li>Multi-core under user control</li> <li>Irregular data structures</li> <li>Irregular flow control</li> </ul> | <ul style="list-style-type: none"> <li>Data parallel programming</li> <li>Little user control</li> <li>Regular data structures</li> <li>Regular Flow Control</li> </ul> |

BTW,  
The general term in the OpenCL world for an SM is a **Compute Unit**.  
The general term in the OpenCL world for a CUDA Core is a **Processing Element**.

University  
Computer Graphics

mp - March 15, 2024

17

### Compute Units and Processing Elements are Arranged in Grids

A GPU Platform can have one or more Devices.

A GPU Device is organized as a grid of Compute Units.

Each Compute Unit is organized as a grid of Processing Elements.

So, in NVIDIA terms, their 4090 GPU has 128 Compute Units, each of which has 128 Processing Elements, for a grand total of 16,384 Processing Elements.

University  
Computer Graphics

mp - March 15, 2024

18

### Thinking ahead to CUDA and OpenCL...

19

#### 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.

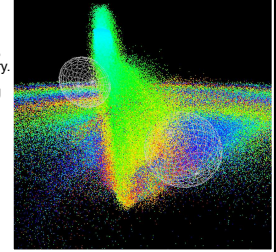
19

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

20

**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.



Ben Weiss

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

20

### Something New – Tensor Cores

21



21

### Tensor Cores Accelerate Fused-Multiply-Add Arithmetic

22

$$D = \begin{pmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{pmatrix} \begin{pmatrix} B_{11} & B_{12} & B_{13} \\ B_{21} & B_{22} & B_{23} \\ B_{31} & B_{32} & B_{33} \end{pmatrix} + \begin{pmatrix} C_{11} & C_{12} & C_{13} \\ C_{21} & C_{22} & C_{23} \\ C_{31} & C_{32} & C_{33} \end{pmatrix}$$

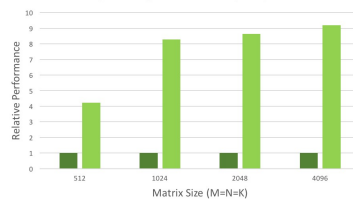
FP16 or FP32

FP16

FP16

FP16 or FP32

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



22

### What is Fused Multiply-Add?

23

Many scientific and engineering computations take the form:  
 $D = A + (B \cdot C);$

A "normal" multiply-add would likely handle this as:  
 $tmp = B \cdot C;$   
 $D = A + tmp;$

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

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

```

123
x 456
----
 738
 615
 492
+ 789
----
56,877
    
```

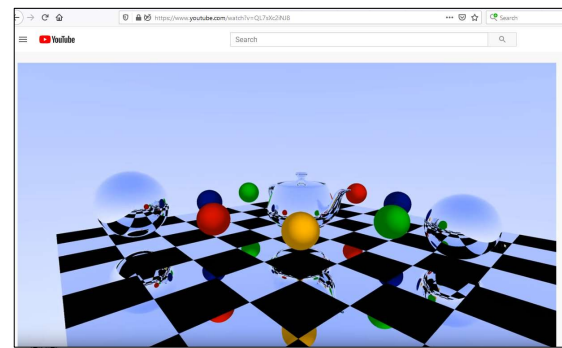
Can start adding the 9 the moment the 8 is produced!

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

23

### Something Even Newer – Ray-Trace Cores

24



<https://www.youtube.com/watch?v=QL7sXc2iNJ8>

24

**There are Two Approaches to Combining CPU and GPU Programs** 25

1. Combine both the CPU and GPU code in the same code file. Somehow mark what part is CPU code and what part is GPU code. The CPU compiler compiles just 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 by a CPU compiler and a GPU compiler.

**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!**

mb - March 15, 2024

25

**Looking ahead:  
If threads all execute the same program,  
what happens on flow divergence?** 26

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

1. On a GPU, the line "if (a > b )" creates a vector of 0/1 Boolean values giving the results of the if-statement for each thread. This becomes a "bitmask".
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.



mb - March 15, 2024

26



- GPUs were originally designed for the streaming-ness of computer graphics
- That same streaming-ness can now also be applied to data-parallel computing
- GPUs are better for some things. CPUs are better for others.



mb - March 15, 2024

27

**Dismantling a Graphics Card** 28

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



mb - March 15, 2024

28

**Dismantling a Graphics Card** 29

Removing the covers:

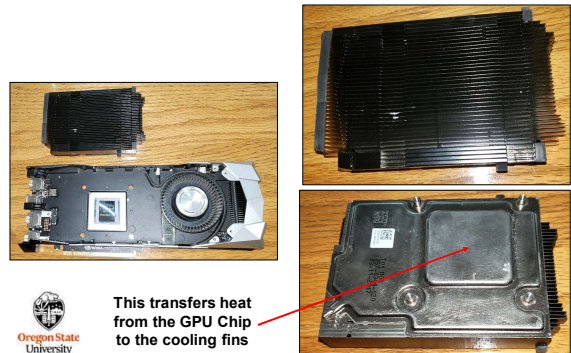


mb - March 15, 2024

29

**Dismantling a Graphics Card** 30

Removing the heat sink:



**This transfers heat from the GPU Chip to the cooling fins**

mb - March 15, 2024

30

**Dismantling a Graphics Card** 31

Removing the fan assembly reveals the board:

GPU Chip      Memory

Oregon State University Computer Graphics

mb - March 15, 2024

31

**Dismantling a Graphics Card** 32

Power half of the board:

Power distribution      Power input

Oregon State University Computer Graphics

mb - March 15, 2024

32

**Dismantling a Graphics Card** 33

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)

Oregon State University Computer Graphics

mb - March 15, 2024

33

**Dismantling a Graphics Card** 34

Underside of the board:

Oregon State University Computer Graphics

mb - March 15, 2024

34

**Dismantling a Graphics Card** 35

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>

Oregon State University Computer Graphics

mb - March 15, 2024

35