# GPU 101

**Mike Bailey**

mjb@cs.oregonstate.edu

**Oregon State University**
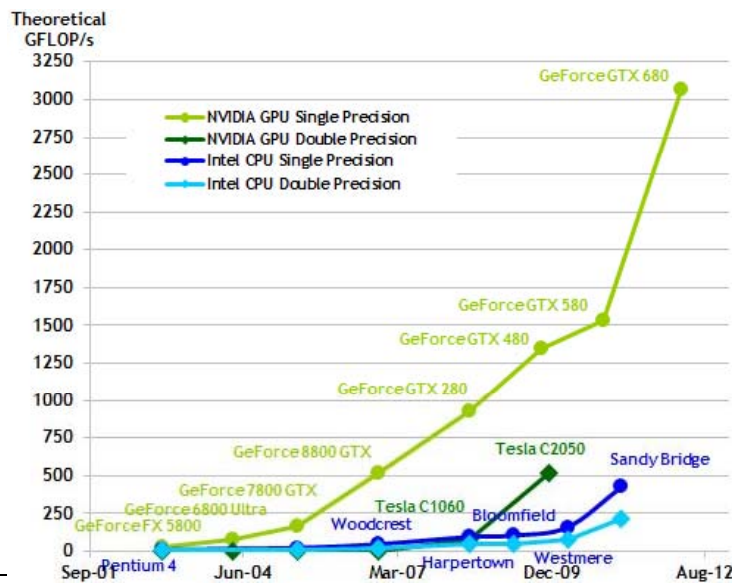
---

# Why do we care about GPU Programming?
# A History of GPU Performance vs. CPU Performance



Theoretical GFLOP/s

- NVIDIA GPU Single Precision
- NVIDIA GPU Double Precision
- Intel CPU Single Precision
- Intel CPU Double Precision

GeForce GTX 680
GeForce GTX 580
GeForce GTX 480
GeForce GTX 280
GeForce 8800 GTX
GeForce 7800 GTX
GeForce 6800 Ultra
GeForce FX 5800
Tesla C2050
Tesla C1060
Sandy Bridge
Woodcrest
Bloomfield
Harpertown
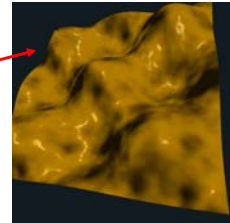Westmere
Pentium 4

Sep-01    Jun-04    Mar-07    Dec-09    Aug-12
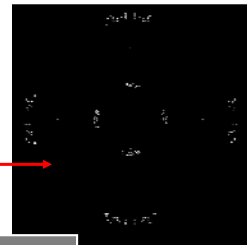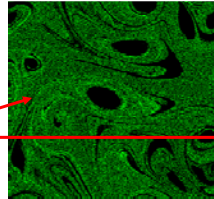
Source: NVIDIA

1

## How Can You Gain Access to GPU Power?

**There are 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 computation (≥ 2002)



3. Write in OpenCL (or CUDA), which looks like C++ (≥ 2006)

---

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

## Why have GPUs Been Outpacing CPUs in Performance?

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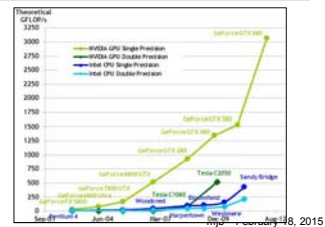
Another reason is that GPU chips do not need the significant amount of **cache** space that occupies much of the real estate on general-purpose CPU chips. The GPU die real estate can then be re-targeted to hold more cores and thus to produce more processing power.

Another reason is that general CPU chips contain on-chip logic to do **branch prediction.** This, too, takes up chip die space.
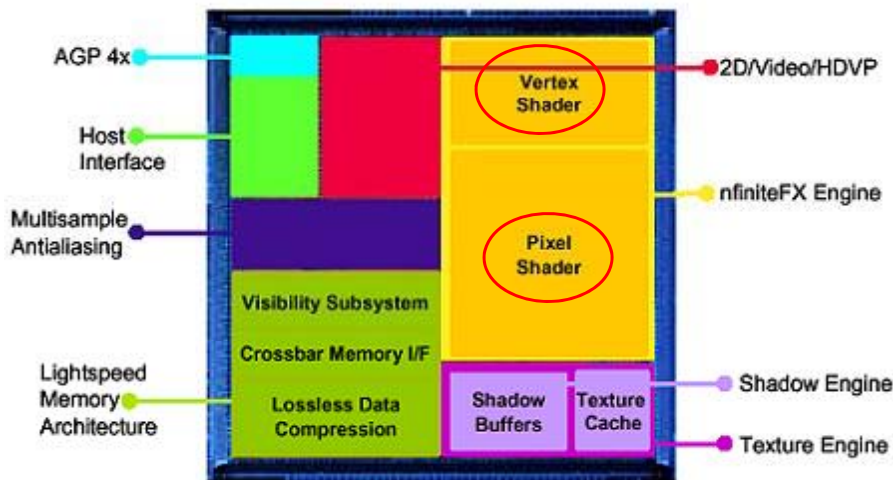
Another reason is that general CPU chips contain on-chip logic to process instructions **out-of-order** if the CPU is blocked and is waiting on something (e.g., a memory fetch). This, too, takes up chip die space.

So, which is bette, CPU or GPUr?
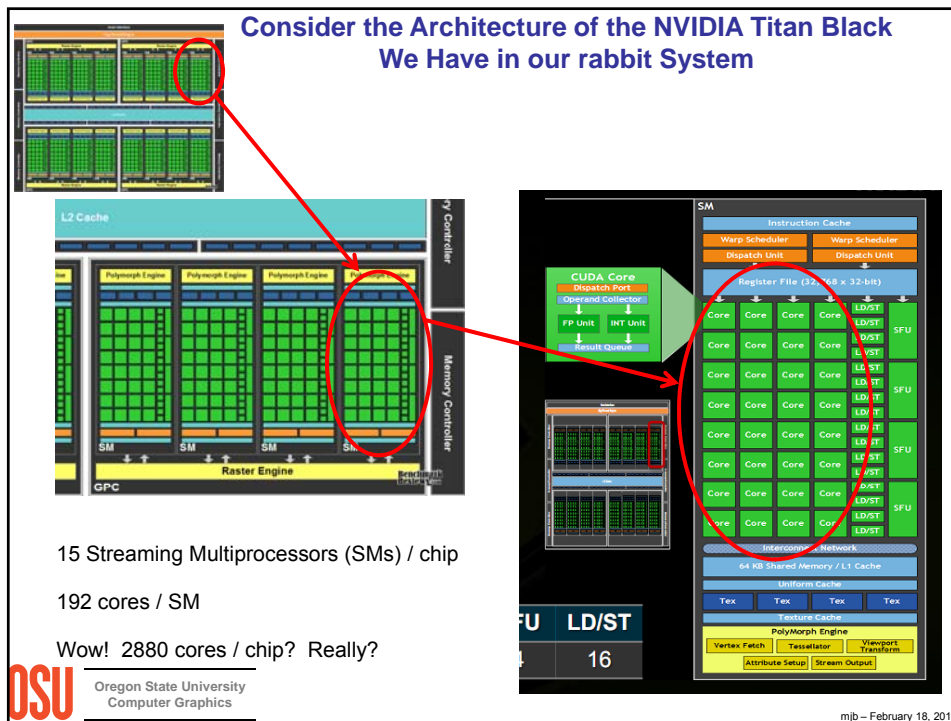
*It depends on what you are trying to do!*

mjb – February 18, 2015

---

## Originally, GPU Devices were very task-specific

mjb – February 18, 2015

## Today's GPU Devices are less task-specific

## Consider the Architecture of the NVIDIA Titan Black We Have in our rabbit System



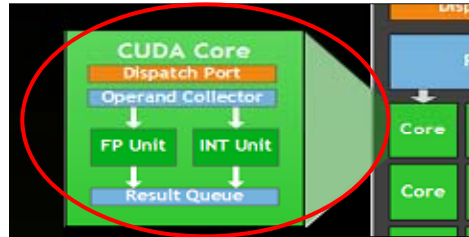15 Streaming Multiprocessors (SMs) / chip

192 cores / SM

Wow!  2880 cores / chip?  Really?

Oregon State University
Computer Graphics

## 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.)
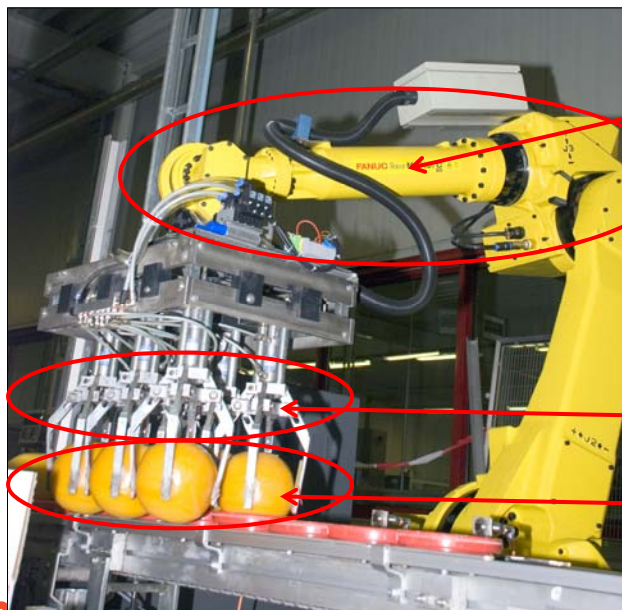
Other vendors refer to these as "Lanes". You might also think of them as 32-way SIMD.

## A Mechanical Equivalent…



"Streaming Multiprocessor"

"CUDA Cores"

"Data"

http://news.cision.com

## The Bottom Line is This

So, the Titan Black has 15 processors per chip, each of which is optimized to do 192-way SIMD. This is an amazing achievement in computing power. But, it is obvious that 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

| **CPU** | **GPU** |
| --- | --- |
| General purpose programming | Data parallel programming |
| Multi-core under user control | Little user control |
| Irregular data structures | Regular data structures |

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**.

---

## How Many Robots Do You See Here?

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

## Compute Units and Processing Elements are Arranged in Grids

Platform

Device #0

Device #1

**Device**

| CU | CU | CU |
| CU | CU | CU |

● ● ●

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, a Fermi Device has 16 Compute Units and each Compute Unit has 32 Processing Elements.

**Compute Unit**

| PE | PE | PE | PE | PE |
| PE | PE | PE | PE | PE |
| PE | PE | PE | PE | PE |

● ● ●

Computer Graphics

---

## Thinking ahead to 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 *hundreds* or *thousands* of them!

• Each thread executes the same program (called the *kernel*), but operates on different small pieces 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.

• OpenCL has built-in functions so that each thread can figure out which number it is, and thus figure out what its job is (i.e., SPMD)

• 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.

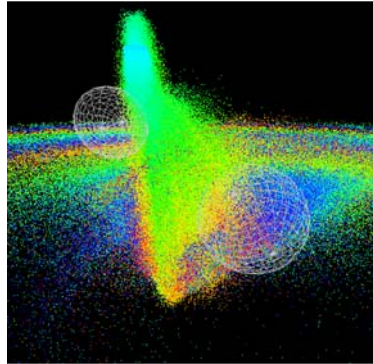**Oregon State University**
**Computer Graphics**

### 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