



## **Timeline for CPUs**

**80's: CPU and system same speed. Zero wait states.**

**1993: CPUs faster than the rest of the system. Rapid raise of frequency.**

**Late 90's to present: Multi-CPU systems, multi-core CPUs.**

**CPUs are still improving, but going for higher frequency is not as obvious as before.**



Information Coding / Computer Graphics, ISY, LiTH

## **Meanwhile, at the graphics dept**

**80's: Hardware sprites. Push pixels with low-level code.**

**1993: Textured 3D games: Wolf3D, Doom.**

**Early 90's: Professional 3D boards.**

**1996: 3dfx Voodoo1!**

**2001: Programmable shaders.**

**2006: G80, unified architecture. CUDA**

**2009: OpenCL.**

**2010: Fermi architecture**

**2012: Kepler architecture**



## Information Coding / Computer Graphics, ISY, LiTH

	1995	2005	
CPU Frequency (GHz)	.1	3.2	32x
Memory Frequency (GHz)	.03	1.2	40x
Bus Bandwidth (GB/sec)	.1	4	40x
Hard Disk Size (GB)	.5	200	400x



## Information Coding / Computer Graphics, ISY, LiTH

	1995	2005	
CPU Frequency (GHz)	.1	3.2	32x
Memory Frequency (GHz)	.03	1.2	40x
Bus Bandwidth (GB/sec)	.1	4	40x
Hard Disk Size (GB)	.5	200	400x
Pixel Fill Rate (GPixels/sec)	.0004	3.3	8250x
Vertex Rate (GVerts/sec)	.0005	.35	700x
Graphics flops (GFlops/sec)	.001	40	40000x
Graphics Bandwidth (GB/sec)	.3	19	63x
Frame Buffer Size (MB)	2	256	128x



## How about 2005-2016?

	2005	2011		2016	
CPU Frequency (GHz)	3.2	3.8	1.18x x cores?	4.0	1.25x
Memory Frequency (GHz)	1.2	2.0	1.67x	3.2	2.67x
Bus Bandwidth (GB/sec)	4	31	7.75x	?	?
Hard Disk Size (GB)	200	4000	20x	8000	40x



## Information Coding / Computer Graphics, ISY, LiTH

	2005	2011		2016	
CPU Frequency (GHz)	3.2	3.8	1.18x x cores?	4.0	1.25x
Memory Frequency (GHz)	1.2	2.0	1.67x	3.2	2.67x
Bus Bandwidth (GB/sec)	4	31	7.75x	?	?
Hard Disk Size (GB)	200	4000	20x	8000	40x
Pixel Fill Rate (GPixels/sec)	3.3	59	18x	128	38x
Vertex Rate (GVerts/sec)	.35	?	?	?	?
Graphics flops (GFlops/sec)	40	2488	62x	16380	409x
Graphics Bandwidth (GB/sec)	19	327.7	17x	512	27x
Frame Buffer Size (MB)	256	3000	12x	8000	31x



**But is this a fair comparison?  
Let us compare apples with apples:  
GFLOPS for both!**

	<b>GPU</b>	<b>CPU</b>
<b>1995:</b>	<b>0.001</b>	<b>0.09</b>
<b>2005:</b>	<b>40</b>	<b>5.6</b>
<b>2011:</b>	<b>2488</b>	<b>91</b>
<b>2015:</b>	<b>7000</b>	<b>176</b>
<b>2016:</b>	<b>16380</b>	<b>400-700*</b>

\* Theoretical, 16 cores

**(Various sources)**



## How about economy: dollar per GFLOPS?

<b>1961:</b>	<b>8.3 trillion</b>
<b>1984:</b>	<b>42 million</b>
<b>1997:</b>	<b>42000 (CPU cluster)</b>
<b>2000:</b>	<b>836-1300</b>
<b>2007:</b>	<b>52</b>
<b>2012:</b>	<b>0.73 (AMD 7970)</b>
<b>2013:</b>	<b>0.22 (PS4)</b>
<b>2015:</b>	<b>0.08 (Radeon R9 295)</b>

(Wikipedia)

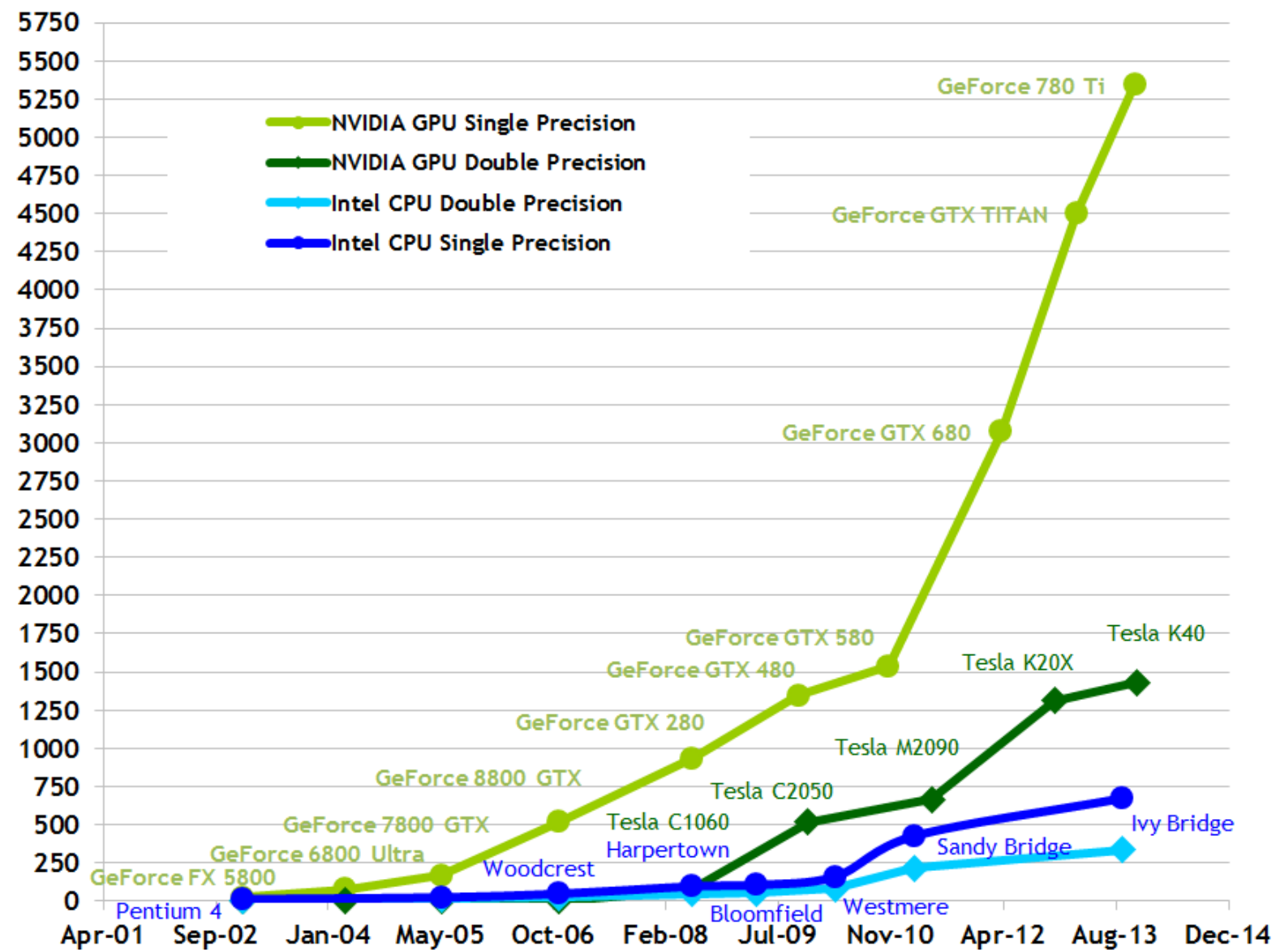




# Information Coding / Computer Graphics, ISY, LiTH

Theoretical GFLOP/s

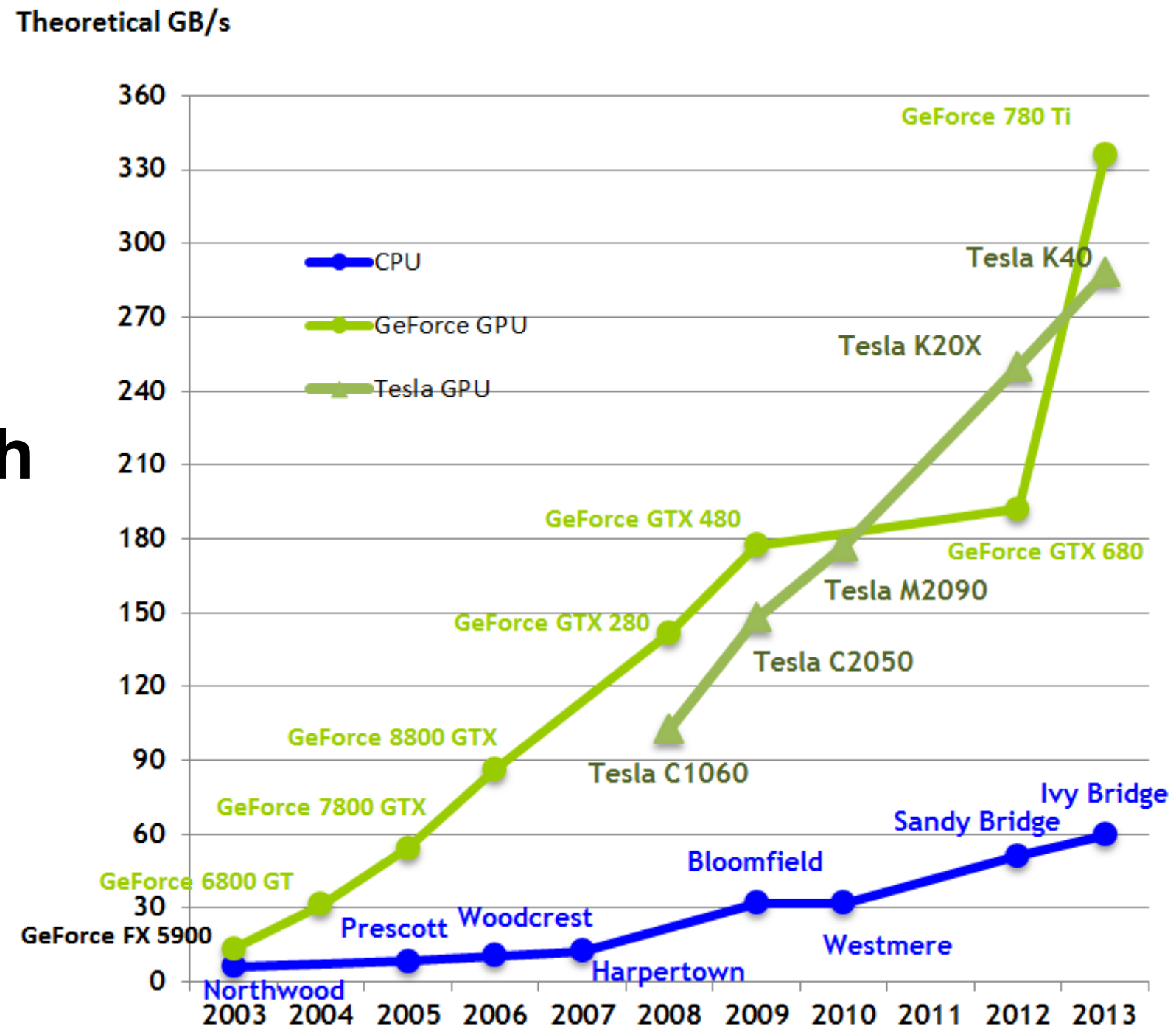
## The GFLOPS race





# Information Coding / Computer Graphics, ISY, LiTH

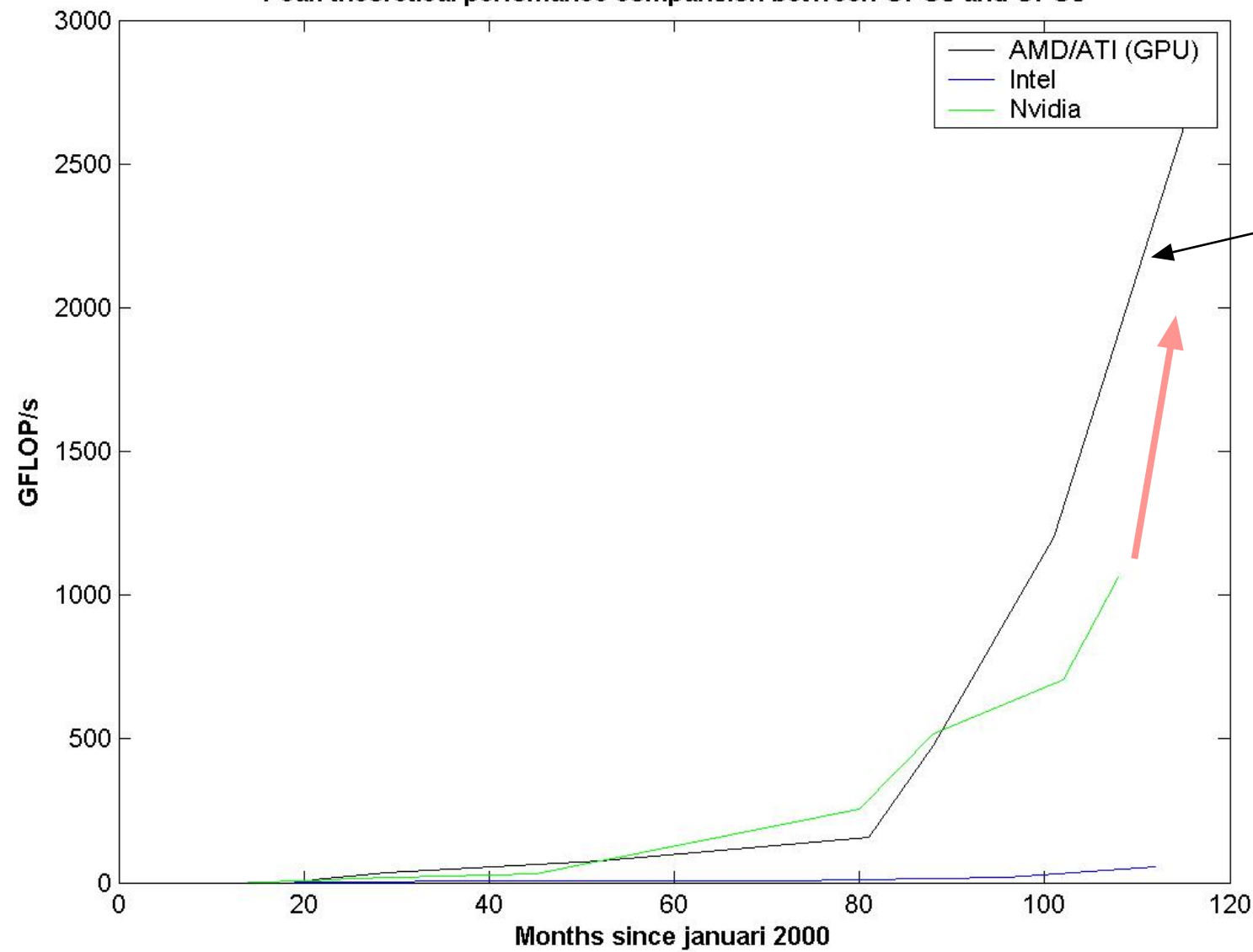
## The bandwidth race





## Another graph, including ATI/AMD

Peak theoretical performance comparison between GPUs and CPUs

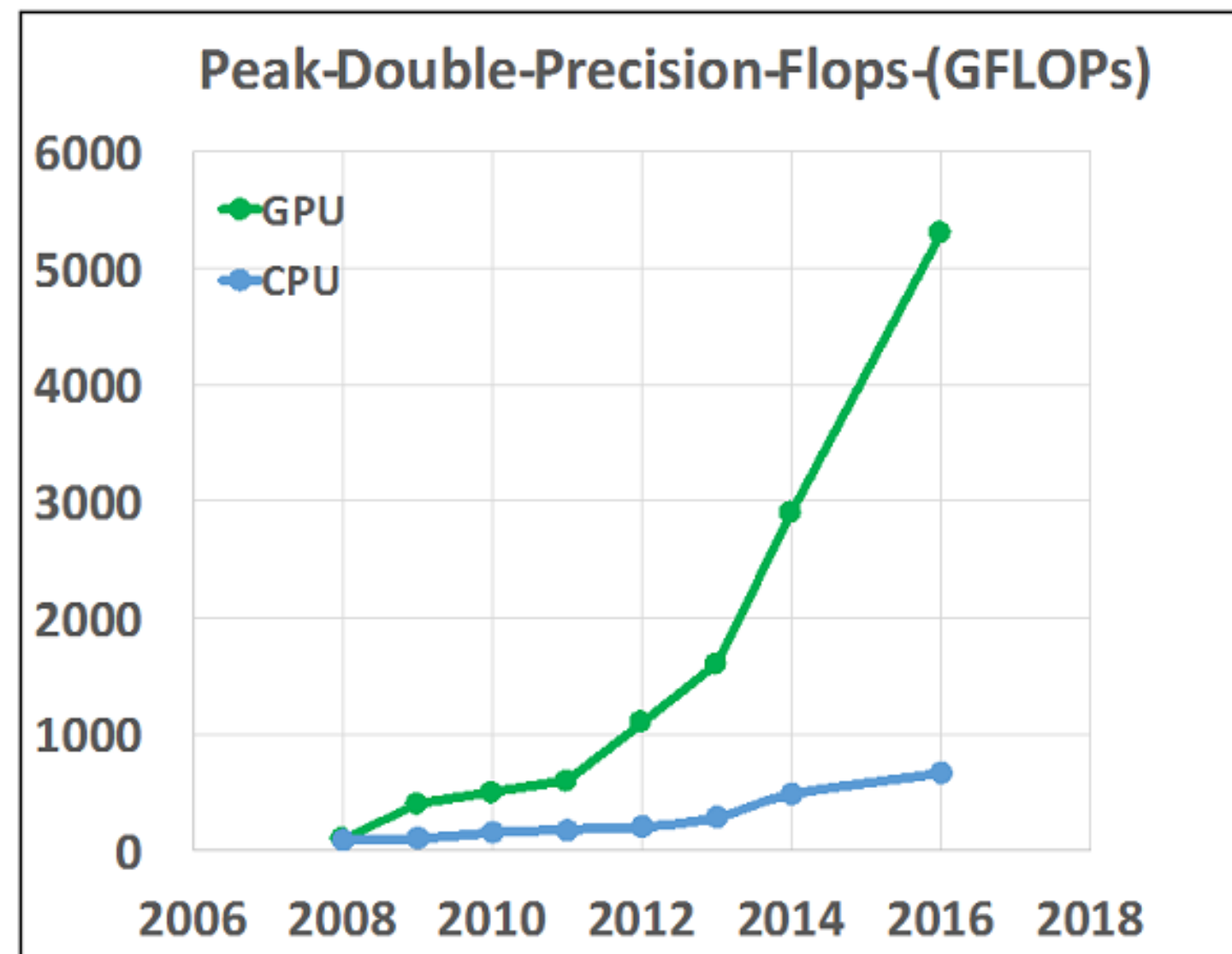


**AMD took the lead in single precision while NVidia was chasing for double with Fermi**



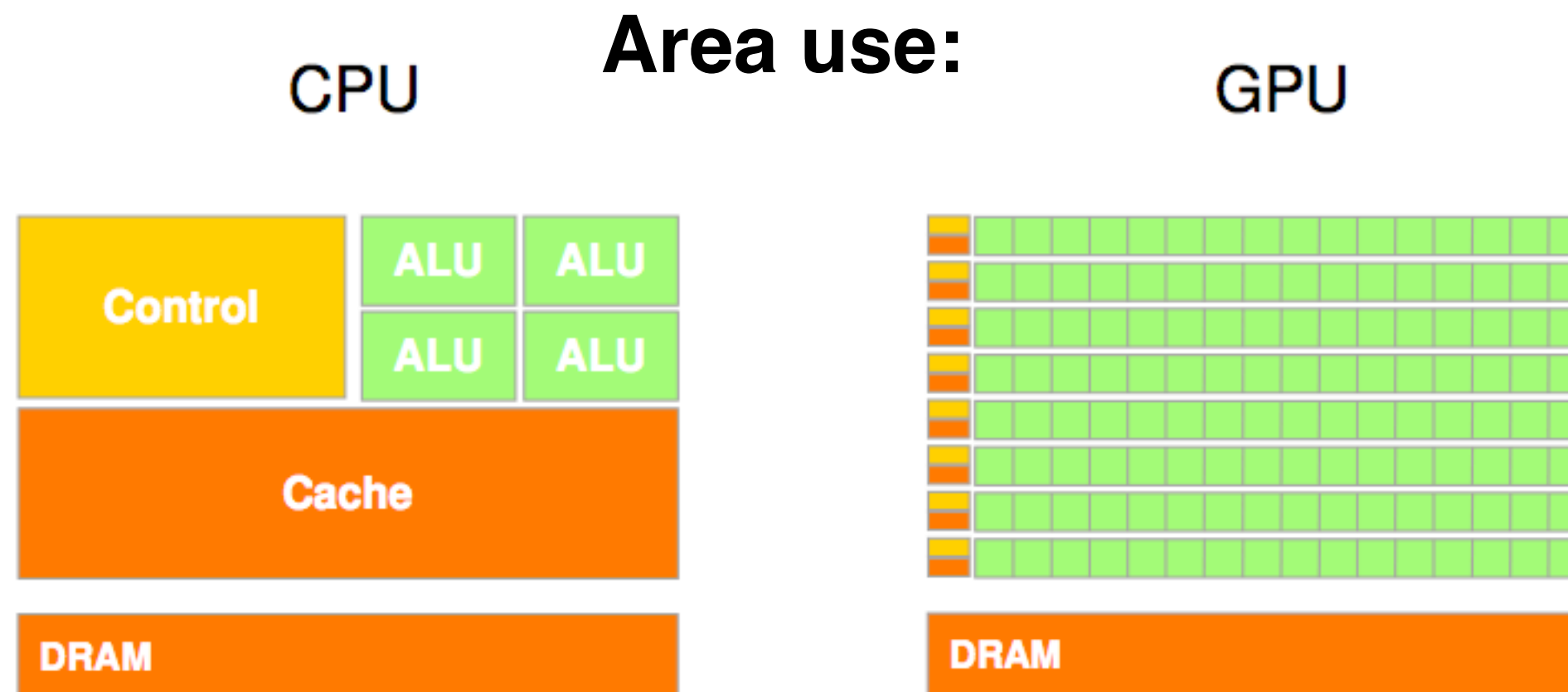
## Information Coding / Computer Graphics, ISY, LiTH

**...up to today.**





## How is this possible?



**But in particular: SIMD architecture**



# Flynn's taxonomy

<b>SISD</b> Single instruction, single data Old single-core systems	<b>MISD</b> Multiple instruction, single data Multiple for redundancy
<b>SIMD</b> Single instruction, multiple data GPUs, vector processors	<b>MIMD</b> Multiple instruction, multiple data Multi-core CPUs

Plus SIMT, single instruction, multiple threads



# SIMD

Single instruction, multiple data

Simplifies instruction handling. All cores get the same instruction.

Excellent for operations where one operation must be made on many data elements.

Is that so common? Yes!

Data best in stored arrays.



# Data Oriented Programming

DOP optimizes for performance.

Data structures selected to fit the computations,  
instead of the programmer!

Optimize for the end user instead for the programmer!

Popular in the game industry - why not elsewhere?





## SIMT - Single Instruction, Multiple Thread

A variant of SIMD.

Parallelism expressed as threads.

A programming model, but also demands that the hardware can handle threads very fast.

Threads dependent - executed in a SIMD processor!

So, why does SIMT fit a graphics processor so well?



## **Why did GPUs get so much performance?**

**Early problem with large amounts of data. (Complex geometry, millions of output pixels.)**

**Graphics pipeline designed for parallelism!**

**Hiding memory latency by parallelism**

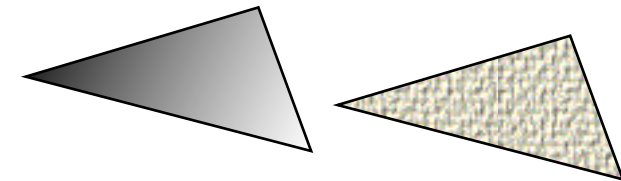
**Volume. 3D graphics boards central component in game industry. Everybody wants one!**

**New games need new impressive features. Many important advancements started as game features.**

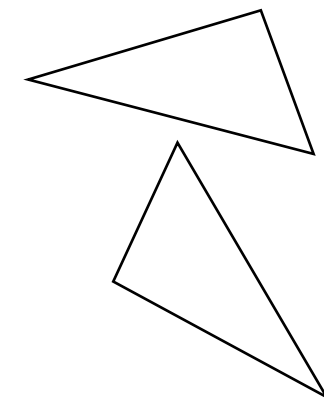


## Information Coding / Computer Graphics, ISY, LiTH

***Must process many pixels fast!***



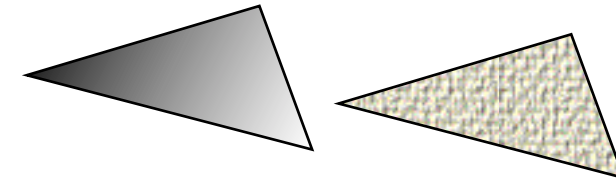
**Early GPUs could draw textured, shaded triangles much faster than the CPU.**





## Information Coding / Computer Graphics, ISY, LiTH

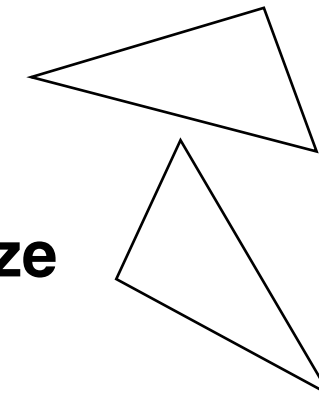
***Must process many pixels fast!***



**Early GPUs could draw textured, shaded triangles much faster than the CPU.**

***Must do matrix multiplication and divisions fast.***

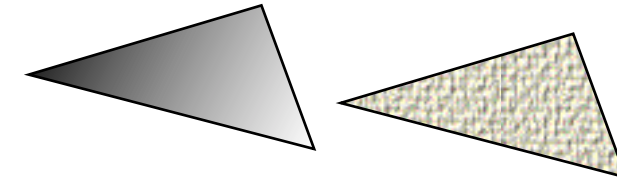
**Next generation could transform vertices and normalize vectors.**





## Information Coding / Computer Graphics, ISY, LiTH

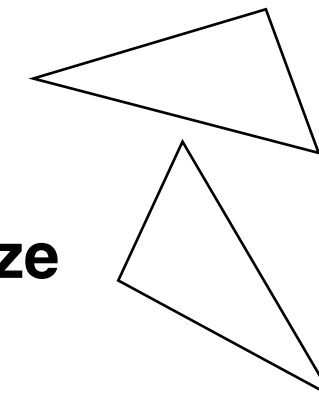
***Must process many pixels fast!***



**Early GPUs could draw textured, shaded triangles much faster than the CPU.**

***Must do matrix multiplication and divisions fast.***

**Next generation could transform vertices and normalize vectors.**



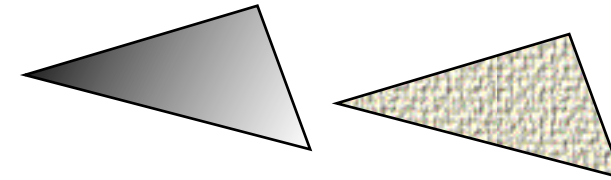
***Must have programmable parts.***

**This was added to make Phong shading and bump mapping.**



## Information Coding / Computer Graphics, ISY, LiTH

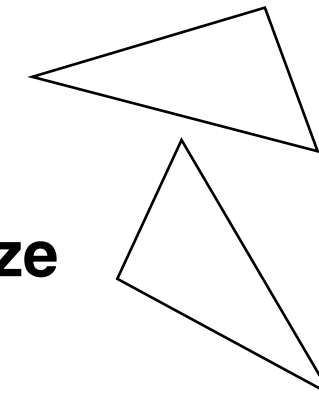
***Must process many pixels fast!***



**Early GPUs could draw textured, shaded triangles much faster than the CPU.**

***Must do matrix multiplication and divisions fast.***

**Next generation could transform vertices and normalize vectors.**



***Must have programmable parts.***

**This was added to make Phong shading and bump mapping.**

***Must work in floating-point!***

**This was for light effects, HDR.**



## **So a GPU should**

- **process vertices, many in parallel, applying the same transformations on each**
- **process pixels (fragments) in parallel, applying the same color/light/texture calculations on each**

**SIMD friendly problem!**

**Less control, control many calculations instead of one**



## **A different kind of threads**

**SIMD threads, all run the same program**

**Group-wise, they execute in parallel, SIMD-style**

**Made for graphics operations: Shader threads calculate  
one pixel or one vertex**

**CUDA/OpenCL threads may calculate anything, but  
typically one part of the output - in order**

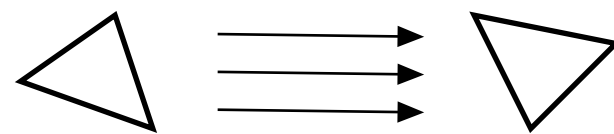




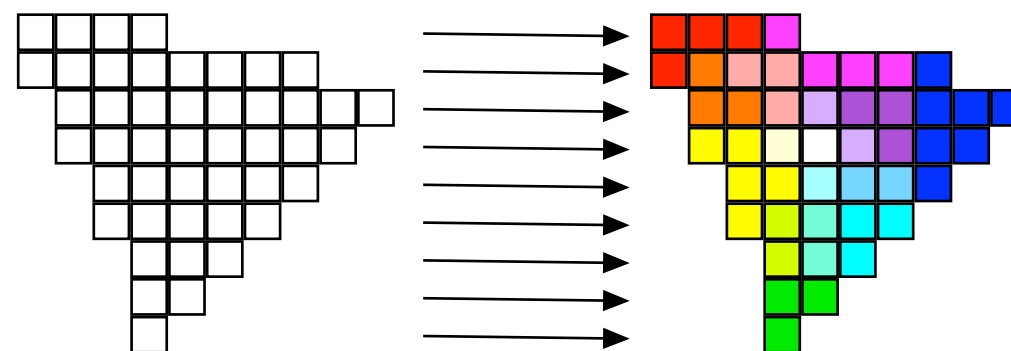
# Information Coding / Computer Graphics, ISY, LiTH

## The main tasks in rendering graphics:

One thread per vertex  
Same operations, same kernel, different data



One thread per pixel (fragment)  
Same operations, same kernel, different data





# The 3D pipeline in the GPU

Low-level operations from vertices to pixel data

