# HIGH-PERFORMANCE COMPUTING WITH NVIDIA TESLA GPUS

## Timothy Lanfear, NVIDIA

# WHY GPU COMPUTING?

# Science is Desperate for Throughput



Gigaflops

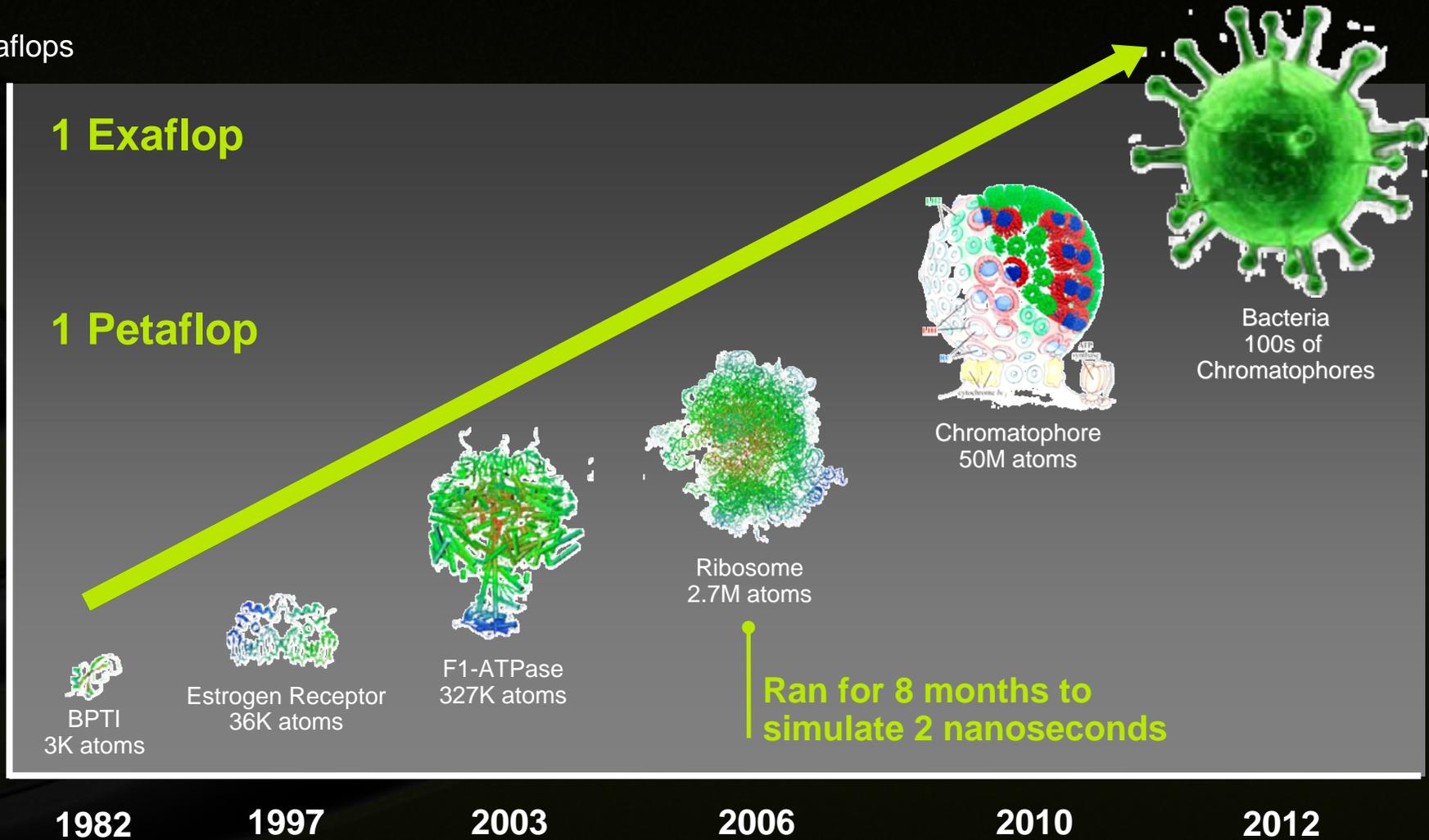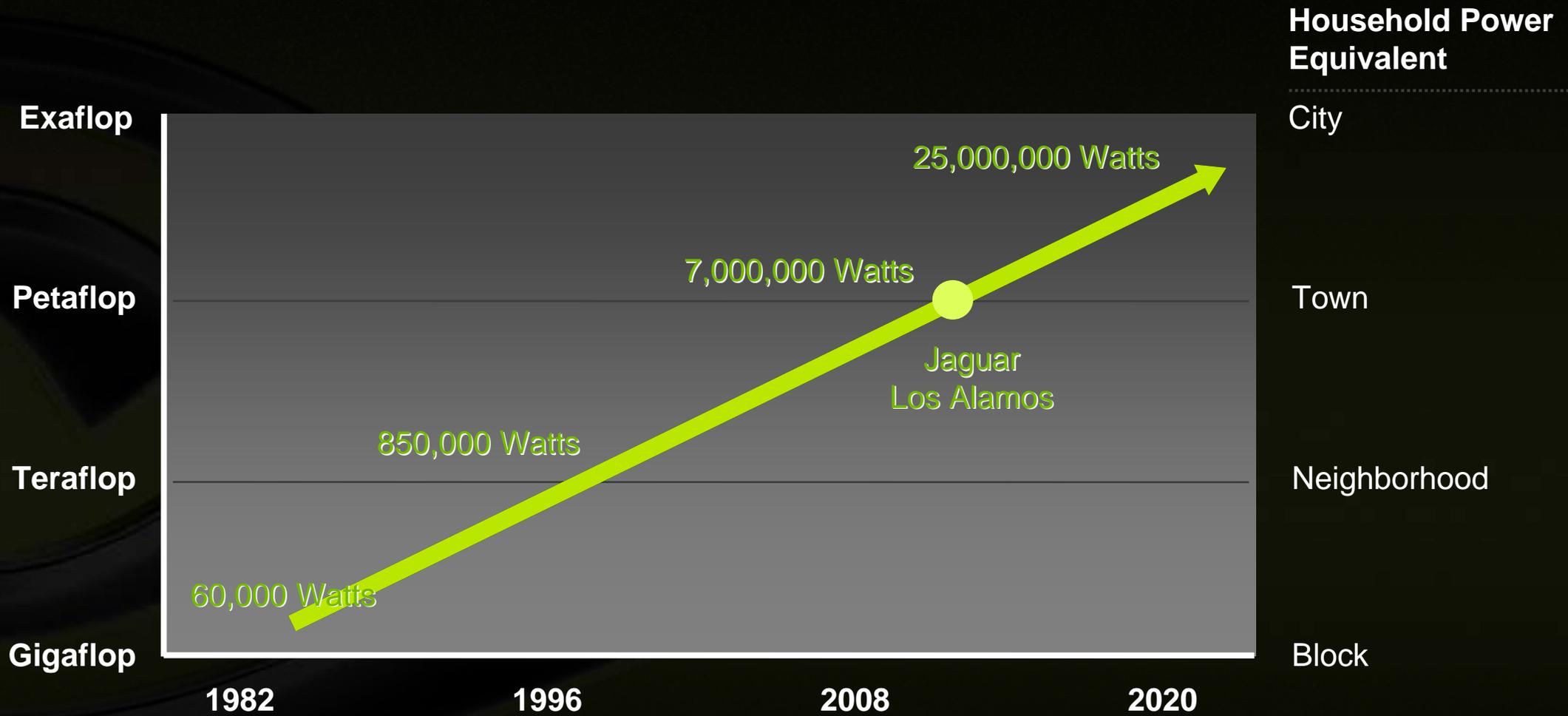| | | | | | | |
|---|---|---|---|---|---|---|
| 1,000,000,000 | 1 Exaflop | | | | | |
| 1,000,000 | 1 Petaflop | | | | | |
| 1,000 | | | | | | |
| 1 | | | | | | |

BPTI
3K atoms

Estrogen Receptor
36K atoms

F1-ATPase
327K atoms

Ribosome
2.7M atoms

Chromatophore
50M atoms

Bacteria
100s of
Chromatophores

**Ran for 8 months to simulate 2 nanoseconds**

| 1982 | 1997 | 2003 | 2006 | 2010 | 2012 |
|---|---|---|---|---|---|

# Power Crisis in Supercomputing

**Household Power Equivalent**

Exaflop — City

25,000,000 Watts

Petaflop — Town

7,000,000 Watts

Jaguar
Los Alamos

850,000 Watts

Teraflop — Neighborhood

60,000 Watts

Gigaflop — Block

1982    1996    2008    2020

# Double Performance per Watt



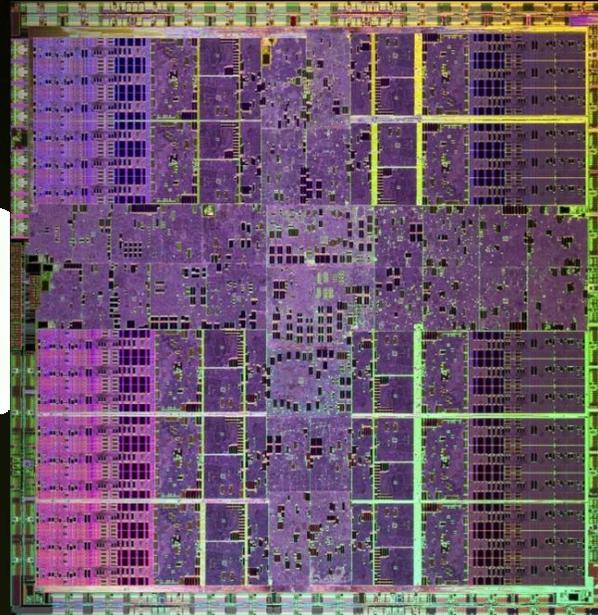© NVIDIA Corporation 2009

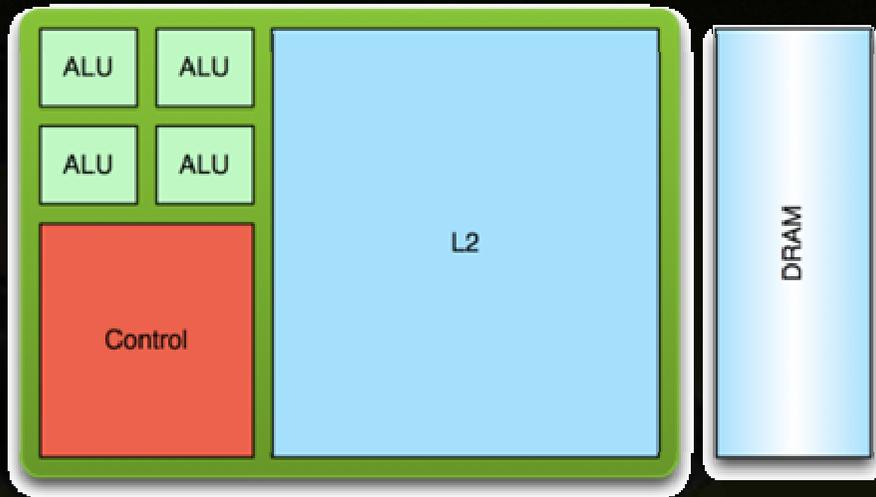# What is GPU Computing?
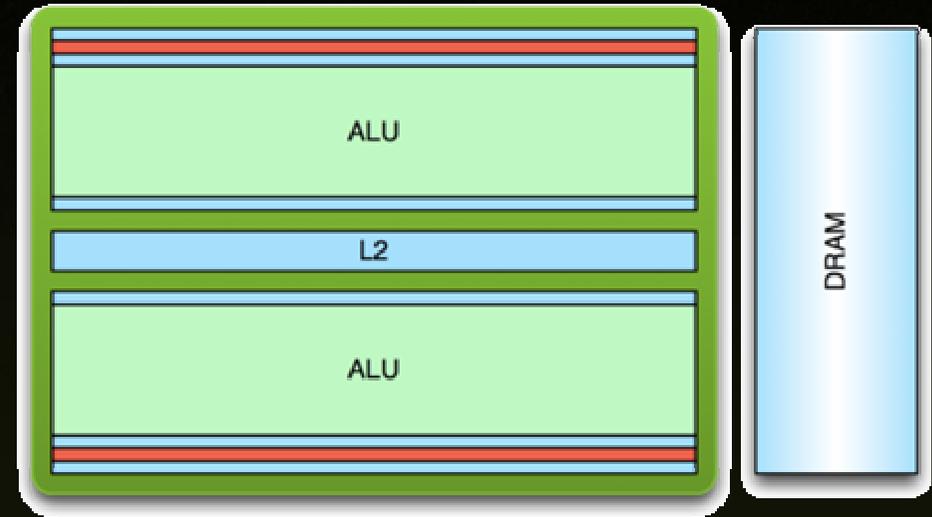
x86          PCIe bus          GPU

**Computing with CPU + GPU**
*Heterogeneous Computing*

# Low Latency or High Throughput?



**CPU**

- Optimised for low-latency access to cached data sets
- Control logic for out-of-order and speculative execution

**GPU**

- Optimised for data-parallel, throughput computation
- Architecture tolerant of memory latency
- More transistors dedicated to computation

# Why Didn't GPU Computing Take Off Sooner?

- **GPU Architecture**
  - **Gaming oriented, process pixel for display**
  - **Single threaded operations**
  - **No shared memory**

- **Development Tools**
  - **Graphics oriented (OpenGL, GLSL)**
  - **University research (Brook)**
  - **Assembly language**

- **Deployment**
  - **Gaming solutions with limited lifetime**
  - **Expensive OpenGL professional graphics boards**
  - **No HPC compatible products**
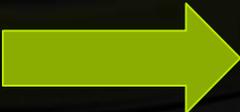
# NVIDIA Invested in GPU Computing in 2004

- **Strategic move for the company**
  - **Expand GPU architecture beyond pixel processing**
  - **Future platforms will be hybrid, multi/many cores based**

- **Hired key industry experts**
  - **x86 architecture**
  - **x86 compiler**
  - **HPC hardware specialist**

**Create a GPU based Compute Ecosystem by 2008**

# NVIDIA : Leadership in GPU Computing

## Over 240 Universities Teaching CUDA

UIUC
MIT
Harvard
Berkeley
Cambridge
Oxford
. . .

IIT Delhi
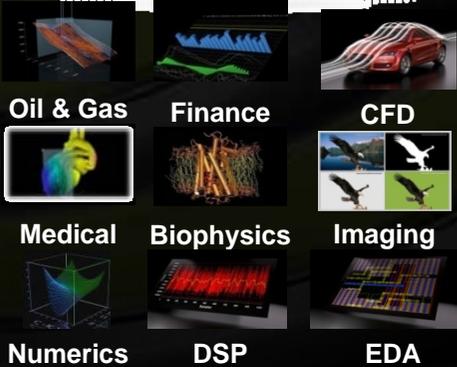Tsinghua
Dortmundt
ETH Zurich
Moscow
NTNU
. . .

## Languages

C, C++
DirectX
Fortran
Java
OpenCL
Python

## Tools

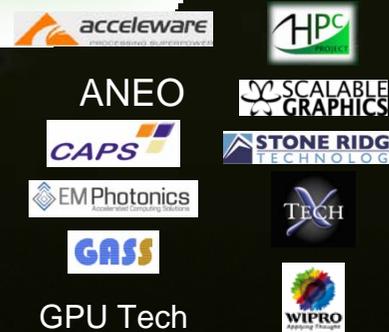PGI Fortran
CAPs HMPP
Nexus
MCUDA
MPI
NOAA Fortran2C
OpenMP

## Applications

Oil & Gas    Finance    CFD

Medical    Biophysics    Imaging

Numerics    DSP    EDA

## Libraries

FFT
BLAS
LAPACK
Image processing
Video processing
Signal processing
Vision

## Consultants

acceleware

ANEO

CAPS

EM Photonics

GASS

GPU Tech

HPC PROJECT

SCALABLE GRAPHICS

STONE RIDGE TECHNOLOGY

TECH

WIPRO
Applying Thought

## OEMs

DELL

CRAY

sgi

APPRO

hp

IBM

FUJITSU

lenovo 联想

PENGUIN COMPUTING

SUPERMICRO

BULL

NEC

# NVIDIA GPU Product Families

**GeForce®**
Entertainment

**Tesla™**
High-Performance Computing

**Quadro®**
Design & Creation
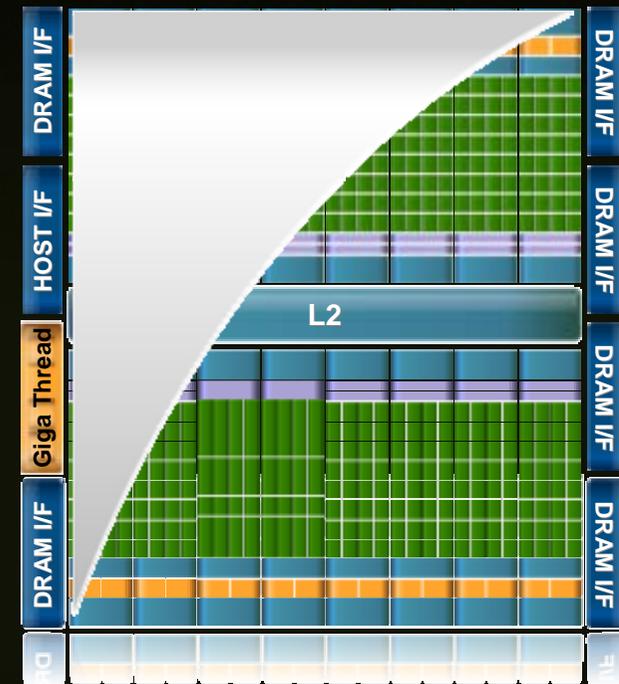
# Fermi: The Computational GPU

**Performance**
- More than ½ Teraflop 64-bit performance
- IEEE 754-2008 SP and DP Floating Point

**Flexibility**
- Increased Shared Memory from 16 KB to 64 KB
- Added L1 and L2 Caches
- ECC on all Internal and External Memories
- Enable up to 1 TeraByte of GPU Memories
- High Speed GDDR5 Memory Interface

**Usability**
- Multiple Simultaneous Tasks on GPU
- 10× Faster Atomic Operations
- C++ Support
- System Calls, printf support

# NVIDIA Tesla GPU Computing Products



## Data Center Products

## Workstation

| | Tesla M1060 | Tesla M2050 | Tesla S2070 | Tesla 2050 | Tesla S1070 | Tesla C2070 | Tesla C2050 | Tesla C1060 |
|---|---|---|---|---|---|---|---|---|
| **GPUs** | 1 T10 GPU | 1 T20 GPU | 4 T20 GPUs | | 4 T10 GPUs | 1 T20 GPU | | 1 T10 GPU |
| Single Precision | 933 GFlops | 1030 GFlops | 4120 GFlops | | 4140 GFlops | 1030 Gflops | | 933 GFlops |
| Double Precision | 78 GFlops | 515 GFlops | 2060 GFlops | | 346 GFlops | 515 Gflops | | 78 GFlops |
| Memory | 4 GB | 3 GB | 12 GB 3 GB / GPU | 24 GB 6 GB / GPU | 16 GB 4 GB / GPU | 6 GB | 3 GB | 4 GB |
| Mem BW | 102 GB/s | 148.4 GB/s | 148.4 GB/s | | 102 GB/s | 144 GB/s | | 102 GB/s |
| Display | No display IO | | No display IO | | | Single dual-link DVI | | No display IO |

# OEM Servers with Tesla M2050 GPUs
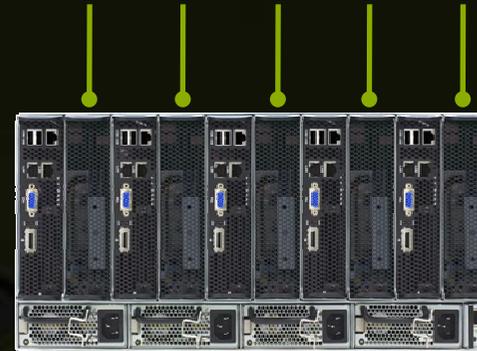
Announced on May 4th, 2010

SUPERMICRO

APPRO

TYAN®

**2 Tesla M2050 GPUs**

**4 Tesla M2050 GPUs**

**10 Tesla M2050 GPUs**

**8 Tesla M2050 GPUs**

SuperServer 6016GT-TF
2 CPUs + 2 GPUs in 1U

Appro Tetra
2 CPUs + 4 GPUs in 1U

Appro GreenBlade
10 CPUs + 10 GPUs in 5U

Tyan B7015
2 CPUs + 8 GPUs in 4U

….. many more coming soon ……

# OEM Servers with Tesla M1060 GPUs

**SUPERMICRO**

**CRAY** THE SUPERCOMPUTER COMPANY

**BULL**

**TYAN®**

**2 Tesla M1060 GPUs**

**Upto 18 Tesla M1060 GPUs**

**8 Tesla M1060 GPUs**

## SuperServer 6016GT-TF
2 CPUs + 2 GPUs in 1U

## Cray CX1000 and Bull Bullx
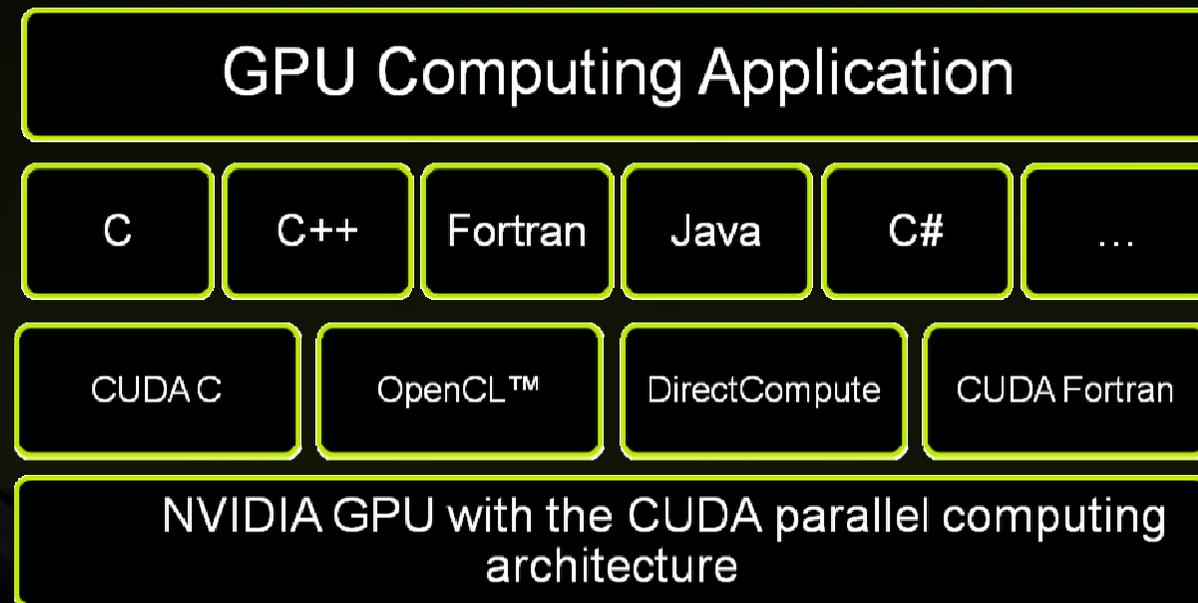36 CPUs + 18 GPUs in 7U

## Tyan B7015
2 CPUs + 8 GPUs in 4U

# CUDA Parallel Computing Architecture

- **Parallel computing architecture and programming model**

- **Includes a CUDA C compiler, support for OpenCL and DirectCompute**

- **Architected to natively support multiple computational interfaces (standard languages and APIs)**



GPU Computing Application

| C | C++ | Fortran | Java | C# | ... |

| CUDA C | OpenCL™ | DirectCompute | CUDA Fortran |

NVIDIA GPU with the CUDA parallel computing architecture

# NVIDIA Parallel Nsight™

**The first development environment for massively parallel applications.**

**Hardware** GPU Source Debugging

**Platform-wide** Analysis

Complete **Visual Studio integration**

**Register for the Beta**
http://developer.nvidia.com/nsight

**Parallel Source Debugging**

**Platform Trace**

**Graphics Inspector**

**GPU Debugging**

**Making it easy**

**Allinea DDT — CUDA Enabled**

# CUDA Zone: www.nvidia.com/CUDA

- **CUDA Toolkit**
  - **Compiler**
  - **Libraries**

- **CUDA SDK**
  - **Code samples**

- **CUDA Profiler**

- **Forums**

- **Resources for CUDA developers**

# Wide Developer Acceptance and Success



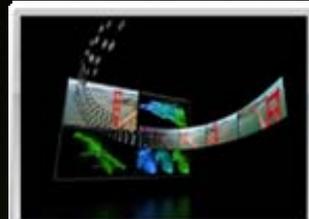| 146X | 36X | 19X | 17X | 100X |
|------|-----|-----|-----|------|
| Interactive visualization of volumetric white matter connectivity | Ion placement for molecular dynamics simulation | Transcoding HD video stream to H.264 | Simulation in Matlab using .mex file CUDA function | Astrophysics N-body simulation |

| 149X | 47X | 20X | 24X | 30X |
|------|-----|-----|-----|------|
| Financial simulation of LIBOR model with swaptions | GLAME@lab: An M-script API for linear Algebra operations on GPU | Ultrasound medical imaging for cancer diagnostics | Highly optimized object oriented molecular dynamics | Cmatch exact string matching to find similar proteins and gene sequences |

# What We Did in the Past Three Years

- **2006**
  - **G80, first GPU with built-in compute features, 128 core multi-threaded, scalable architecture**
  - **CUDA SDK Beta**
- **2007**
  - **Tesla HPC product line**
  - **CUDA SDK 1.0, 1.1**
- **2008**
  - **GT200, second GPU generation, 240 core, 64-bit**
  - **Tesla HPC second generation**
  - **CUDA SDK 2.0**
- **2009 …**

# NEXT-GENERATION GPU ARCHITECTURE — 'FERMI'

# Introducing the 'Fermi' Architecture
## *The Soul of a Supercomputer in the body of a GPU*

- 3 billion transistors
- Over 2× the cores (512 total)
- 8× the peak DP performance
- ECC
- L1 and L2 caches
- ~2× memory bandwidth (GDDR5)
- Up to 1 Terabyte of GPU memory
- Concurrent kernels
- Hardware support for C++

# Design Goal of Fermi



**Data Parallel**

**Instruction Parallel**

GPU

CPU

**Many Decisions**

**Large Data Sets**

- Expand performance sweet spot of the GPU

- Bring more users, more applications to the GPU

# Streaming Multiprocessor Architecture

- **32 CUDA cores per SM (512 total)**

- **2:1 ratio SP:DP floating-point performance**

- **Dual Thread Scheduler**

- **64 KB of RAM for shared memory and L1 cache (configurable)**



Instruction Cache

| Scheduler | Scheduler |
| --- | --- |
| Dispatch | Dispatch |

Register File

| Core | Core | Core | Core |
| --- | --- | --- | --- |
| Core | Core | Core | Core |
| Core | Core | Core | Core |
| Core | Core | Core | Core |
| Core | Core | Core | Core |
| Core | Core | Core | Core |
| Core | Core | Core | Core |
| Core | Core | Core | Core |

Load/Store Units × 16

Special Func Units × 4

Interconnect Network

64K Configurable Cache/Shared Mem

Uniform Cache

# CUDA Core Architecture

- New IEEE 754-2008 floating-point standard, surpassing even the most advanced CPUs

- Fused multiply-add (FMA) instruction for both single and double precision

- Newly designed integer ALU optimized for 64-bit and extended precision operations

**CUDA Core**

Dispatch Port

Operand Collector

FP Unit | INT Unit

Result Queue

Instruction Cache

Scheduler | Scheduler

Dispatch | Dispatch

Register File

Core | Core | Core | Core

Core | Core | Core | Core

Core | Core | Core | Core

Core | Core | Core | Core

Core | Core | Core | Core

Core | Core | Core | Core

Core | Core | Core | Core

Core | Core | Core | Core

Load/Store Units x 16

Special Func Units x 4

Interconnect Network

64K Configurable Cache/Shared Mem

Uniform Cache

# Cached Memory Hierarchy

- First GPU architecture to support a true cache hierarchy in combination with on-chip shared memory

- L1 Cache per SM (32 cores)
  - Improves bandwidth and reduces latency

- Unified L2 Cache (768 KB)
  - Fast, coherent data sharing across all cores in the GPU

## Parallel DataCache™ Memory Hierarchy

# Larger, Faster Memory Interface

- **GDDR5 memory interface**
  - **2× speed of GDDR3**

- **Up to 1 Terabyte of memory attached to GPU**
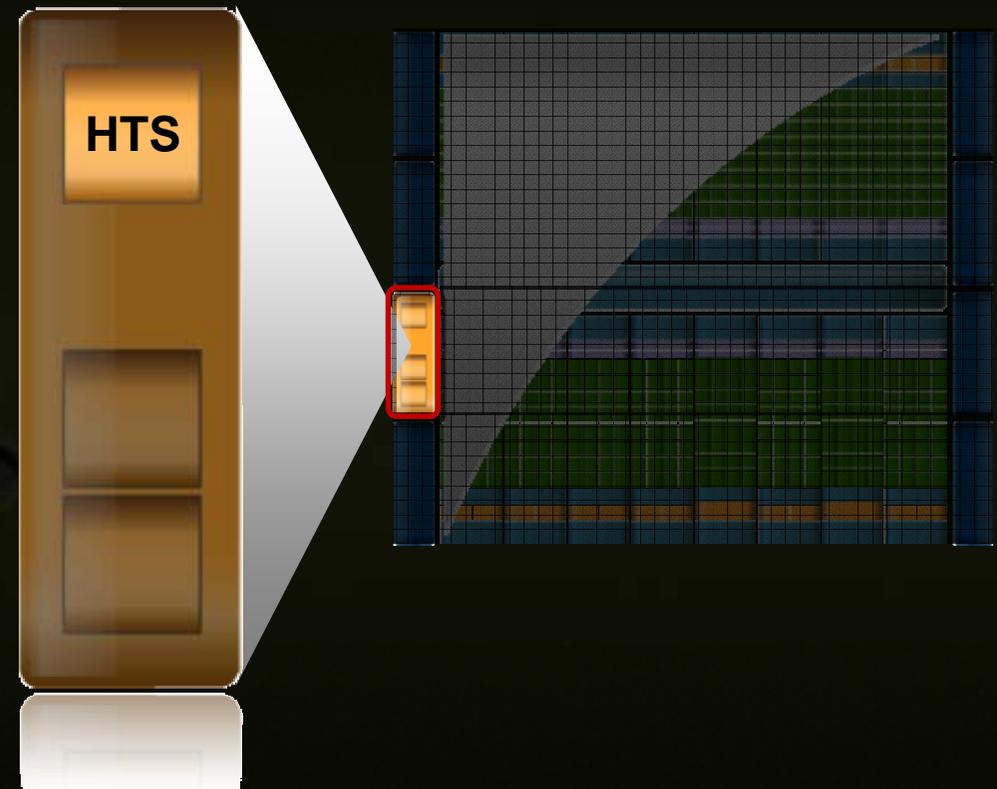  - **Operate on large data sets**

# Error Correcting Code

- **ECC protection for**
  - **DRAM**
    - ECC supported for GDDR5 memory

  - **All major internal memories are ECC protected**
    - Register file, L1 cache, L2 cache
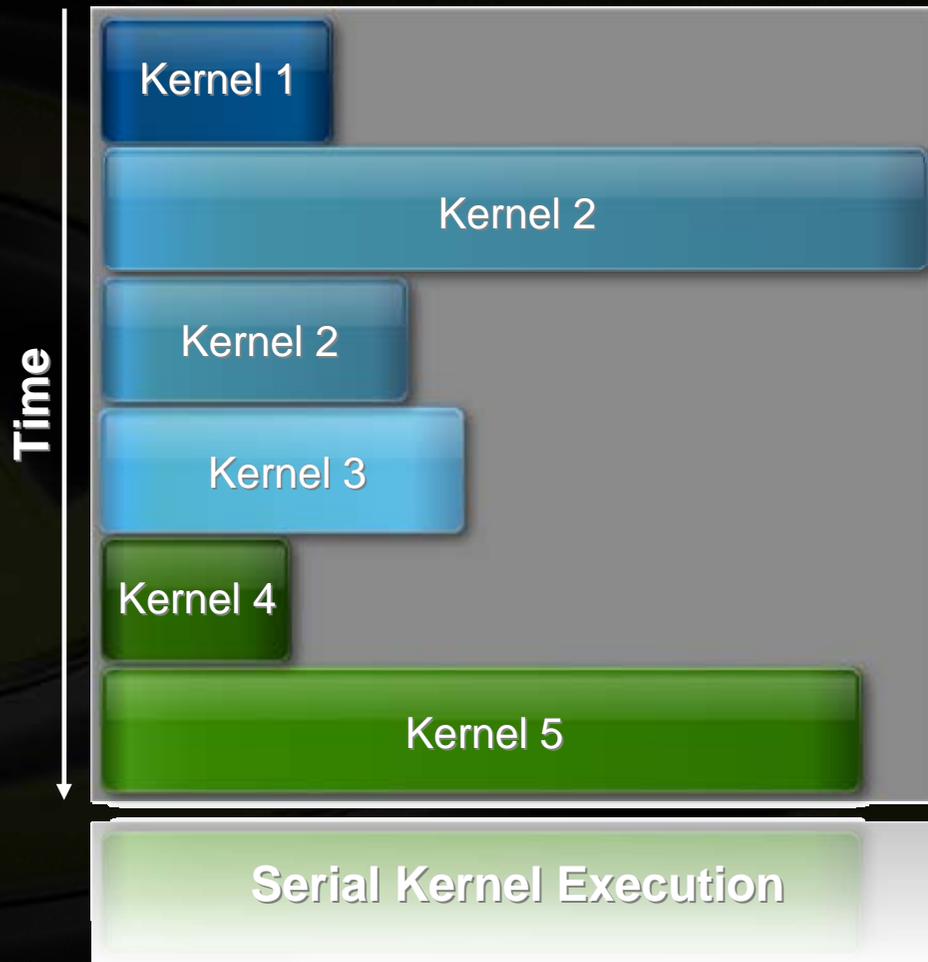
# GigaThread<sup>TM</sup> Hardware Thread Scheduler

- Hierarchically manages thousands of simultaneously active threads

- 10× faster application context switching

- Concurrent kernel execution

HTS

# GigaThread Hardware Thread Scheduler

## Concurrent Kernel Execution + Faster Context Switch



**Serial Kernel Execution**
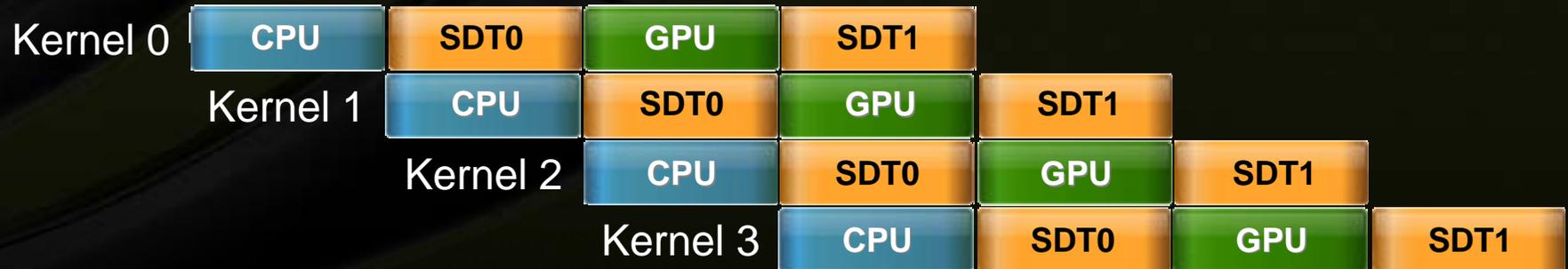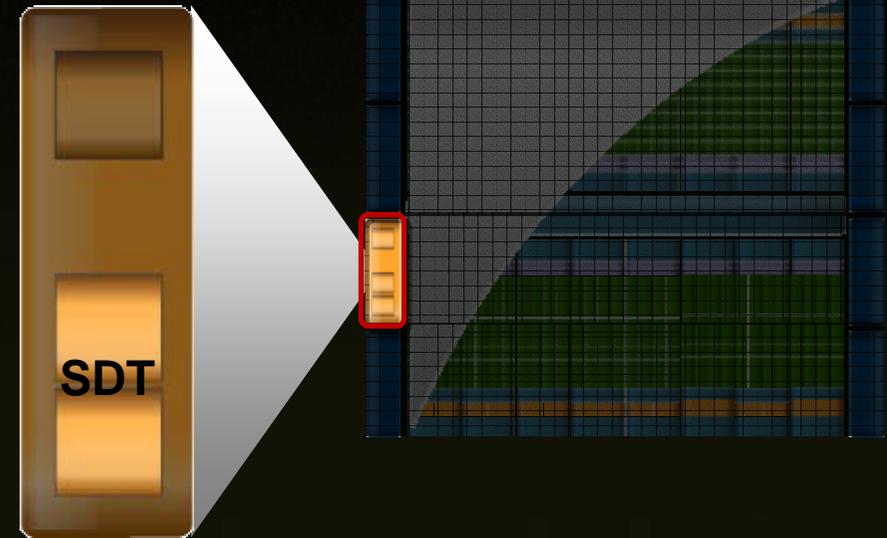
**Parallel Kernel Execution**

# GigaThread Streaming Data Transfer Engine

- **Dual DMA engines**
  - Simultaneous CPU→GPU and GPU→CPU data transfer
  - Fully overlapped with CPU and GPU processing time

- **Activity Snapshot:**

| | CPU | SDT0 | GPU | SDT1 |
|---|---|---|---|---|
| Kernel 0 | CPU | SDT0 | GPU | SDT1 |
| Kernel 1 | CPU | SDT0 | GPU | SDT1 |
| Kernel 2 | CPU | SDT0 | GPU | SDT1 |
| Kernel 3 | CPU | SDT0 | GPU | SDT1 |

SDT

# Enhanced Software Support

○ **Full C++ Support**

   ○ **Virtual functions**

   ○ **Try/Catch hardware support**

○ **System call support**

   ○ **Support for pipes, semaphores, printf, etc**

○ **Unified 64-bit memory addressing**

"Fermi is the world's first complete GPU computing architecture."

**Peter Glaskowsky**
Technology Analyst
The Envisioneering Group

"The convergence of new, fast GPUs optimized for computation as well as 3-D graphics acceleration and industry-standard software development tools marks the real beginning of the GPU computing era. Gentlemen, start your GPU computing engines."

**Nathan Brookwood**
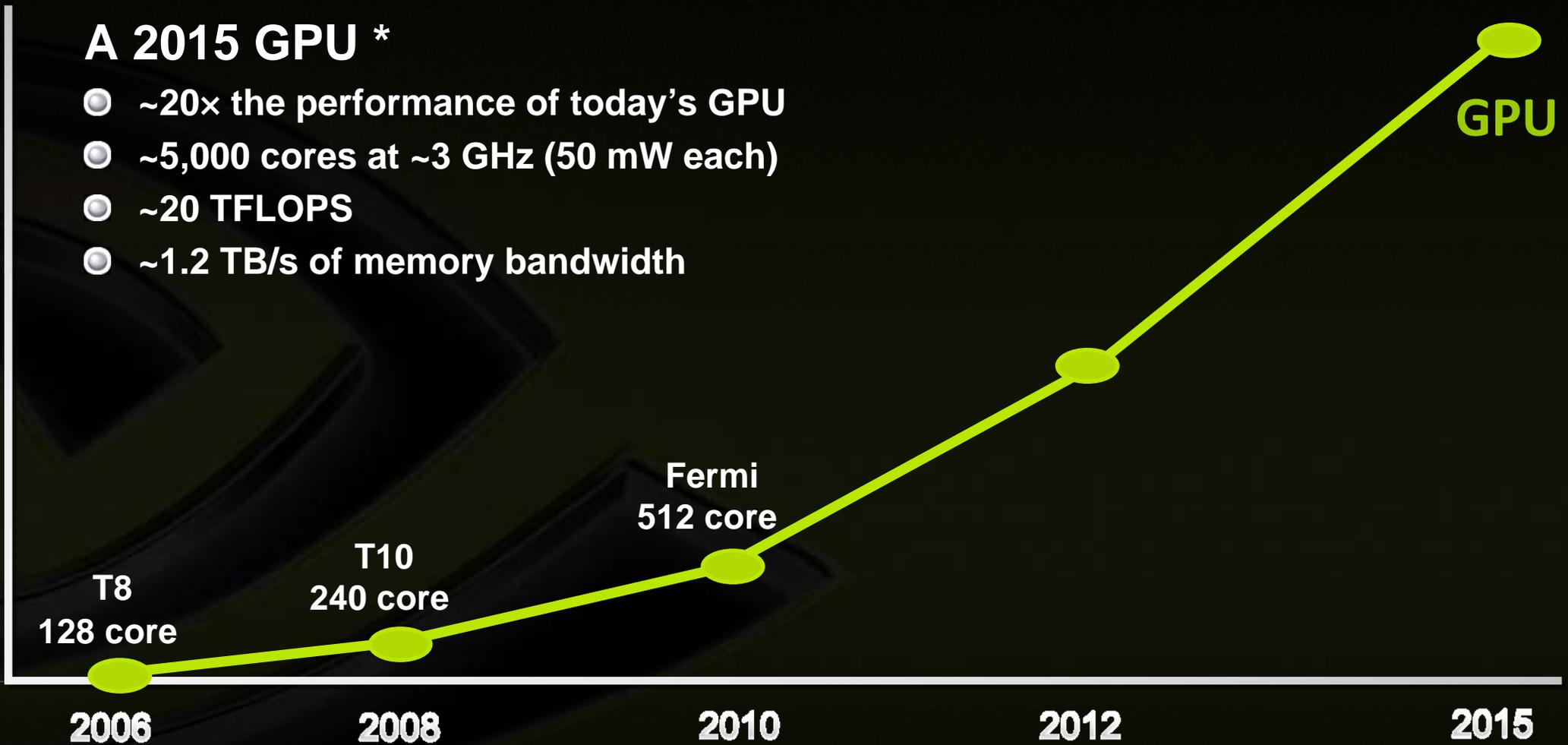Principle Analyst & Founder
Insight 64

# GPU Revolutionizing Computing

**GFlops**

## A 2015 GPU *

- ~20× the performance of today's GPU
- ~5,000 cores at ~3 GHz (50 mW each)
- ~20 TFLOPS
- ~1.2 TB/s of memory bandwidth

**GPU**

**Fermi
512 core**

**T10
240 core**

**T8
128 core**

2006     2008     2010     2012     2015

* This is a sketch of a what a GPU in 2015 might look like; it does not reflect any actual product plans