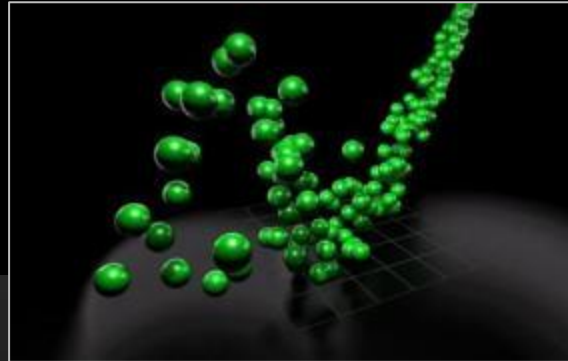
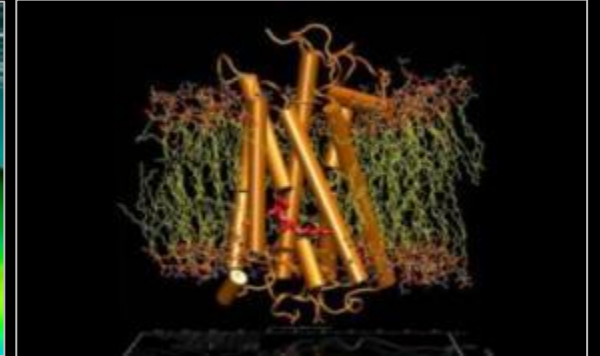
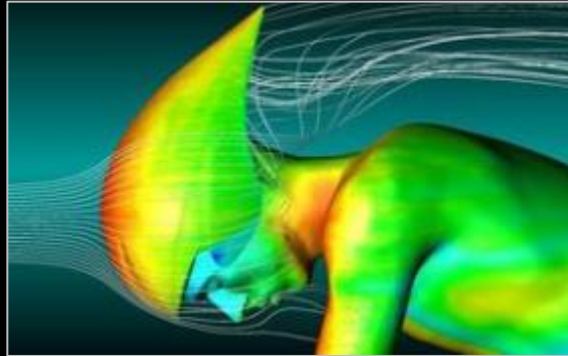


# TESLA

## GPU Computing



Supercomputing at 1/10<sup>th</sup> the Cost

<http://www.nvidia.com/tesla>

# Mainstream Applications Going Parallel

## *CUDA Accelerates Adobe Mercury Playback Engine*



Amazingly fluid,  
real-time video editing

---

Quick preview of real time  
edits and effects

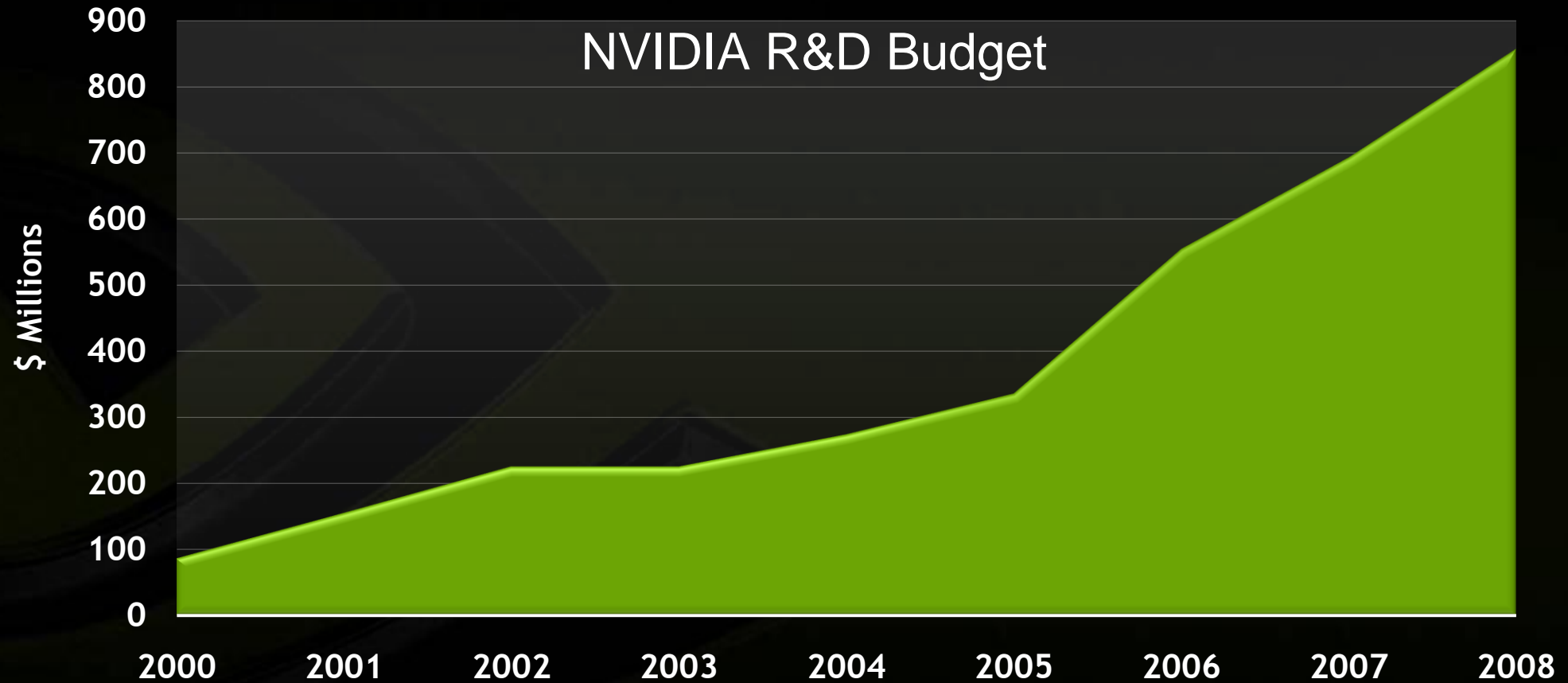
---

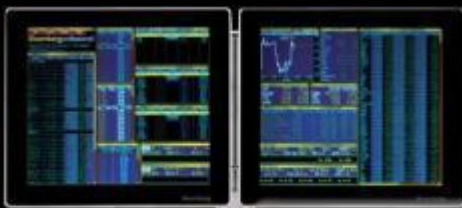
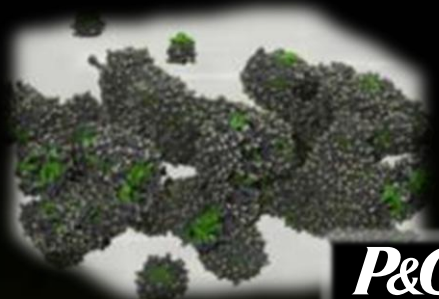
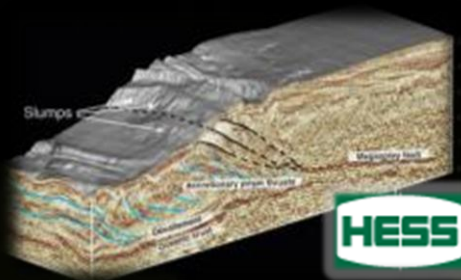
Realistic preview of final  
content

---

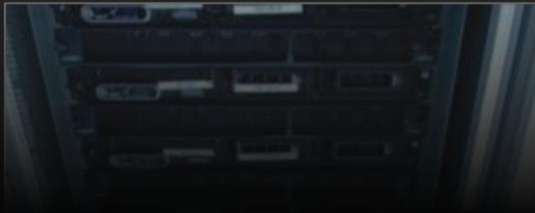
Faster encoding

# GPU Innovation Accelerating

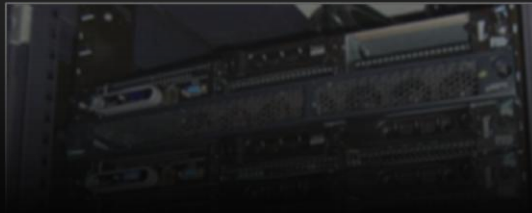




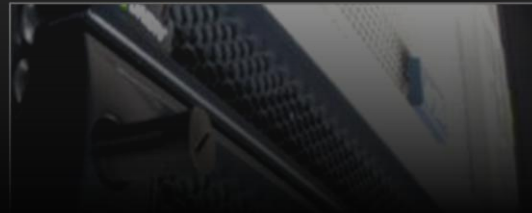
**Max Planck  
Institute**



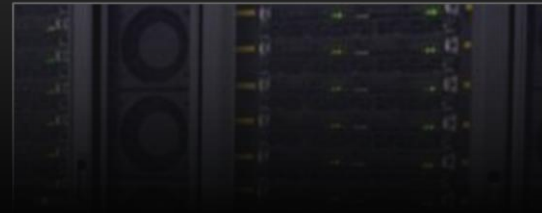
**National Center for  
Supercomputing Applications**



**Tokyo Institute  
of Technology**



**CSIRO**

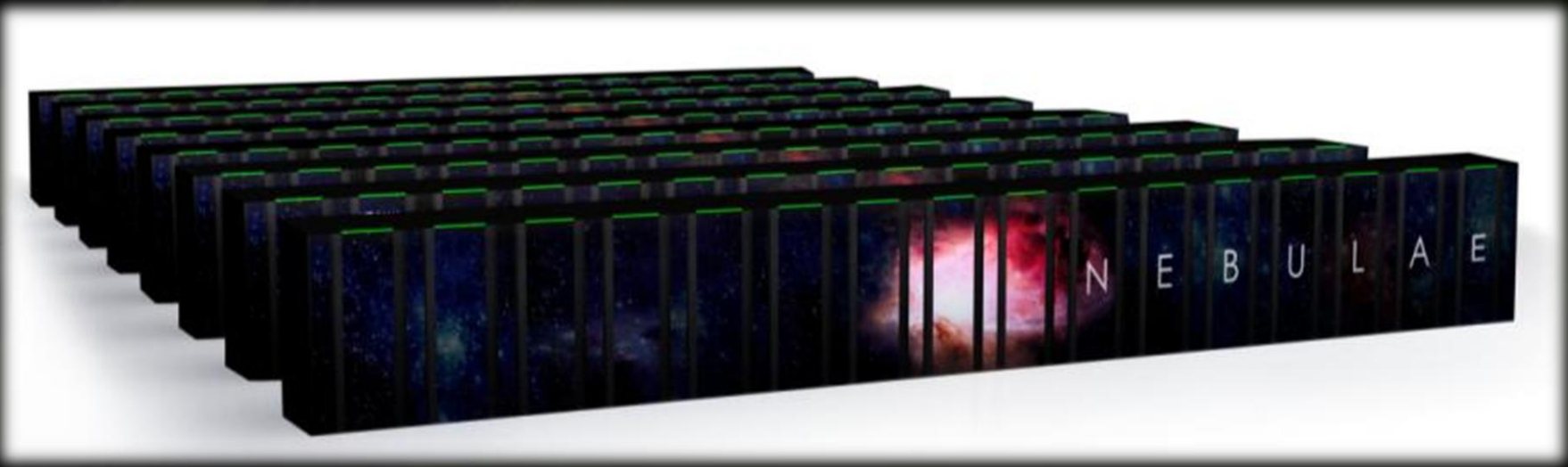


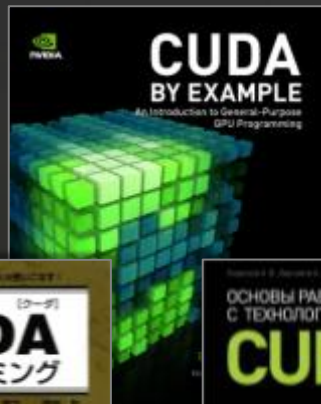
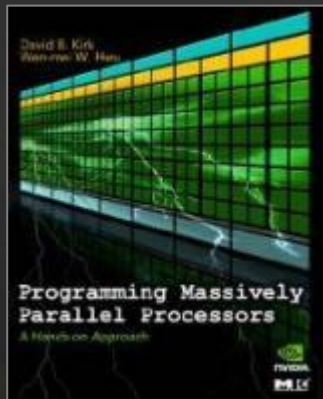
# Dawning Nebulae

Second Fastest Supercomputer in the World

1.27 Petaflop

4640 Tesla GPUs





# 1000+ GPU Clusters Around the World





# NVIDIA Tesla 20-Series Products

Data Center

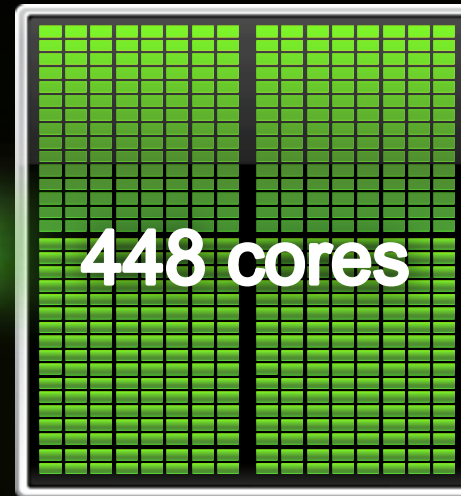
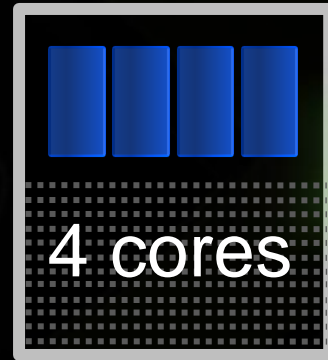


Workstation



# GPU Computing

CPU + GPU Co-Processing



**CPU**

48 GigaFlops (DP)

**GPU**

515 GigaFlops (DP)

# NVIDIA Developer Eco-System

## Numerical Packages

MATLAB  
Mathematica  
NI LabView  
pyCUDA

## Debuggers & Profilers

cuda-gdb  
NV Visual Profiler  
Parallel Nsight  
Visual Studio  
Allinea  
TotalView

## GPU Compilers

C  
C++  
Fortran  
OpenCL  
DirectCompute  
Java  
Python

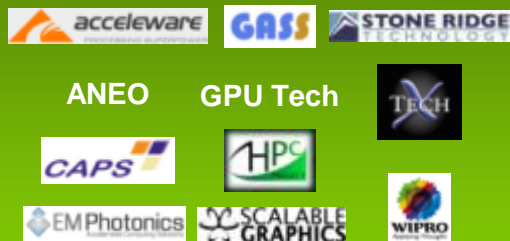
## Parallelizing Compilers

PGI Accelerator  
CAPS HMPP  
mCUDA  
OpenMP

## Libraries

BLAS  
FFT  
LAPACK  
NPP  
Video  
Imaging  
GPULib

## GPGPU Consultants & Training



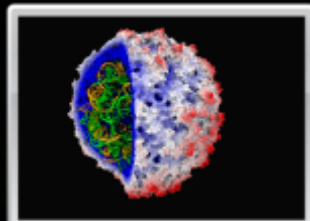
## OEM Solution Providers





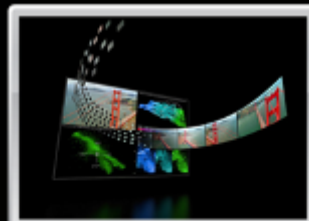
**146X**

Medical Imaging  
U of Utah



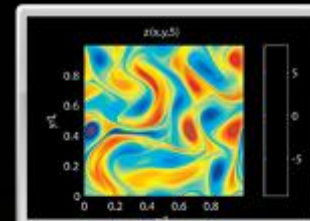
**36X**

Molecular Dynamics  
U of Illinois, Urbana



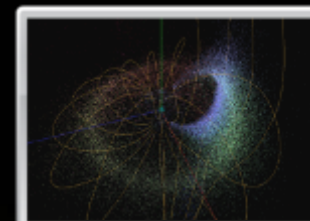
**18X**

Video Transcoding  
Elemental Tech



**50X**

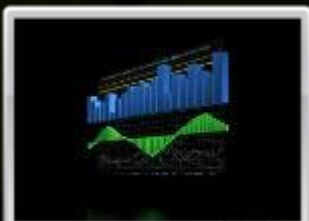
Matlab Computing  
AccelerEyes



**100X**

Astrophysics  
RIKEN

**50x – 150x**



**149X**

Financial simulation  
Oxford



**47X**

Linear Algebra  
Universidad Jaime



**20X**

3D Ultrasound  
Techniscan



**130X**

Quantum Chemistry  
U of Illinois, Urbana



**30X**

Gene Sequencing  
U of Maryland

# Tesla Data Center & Workstation GPU Solutions



**Tesla M-series GPUs**  
M2070 M2050 M1060



**Tesla S-series 1U Systems**  
S2050 S1070



**Tesla C-series GPUs**  
C2070 C2050 C1060



**Integrated CPU-GPU  
Servers & Blades**



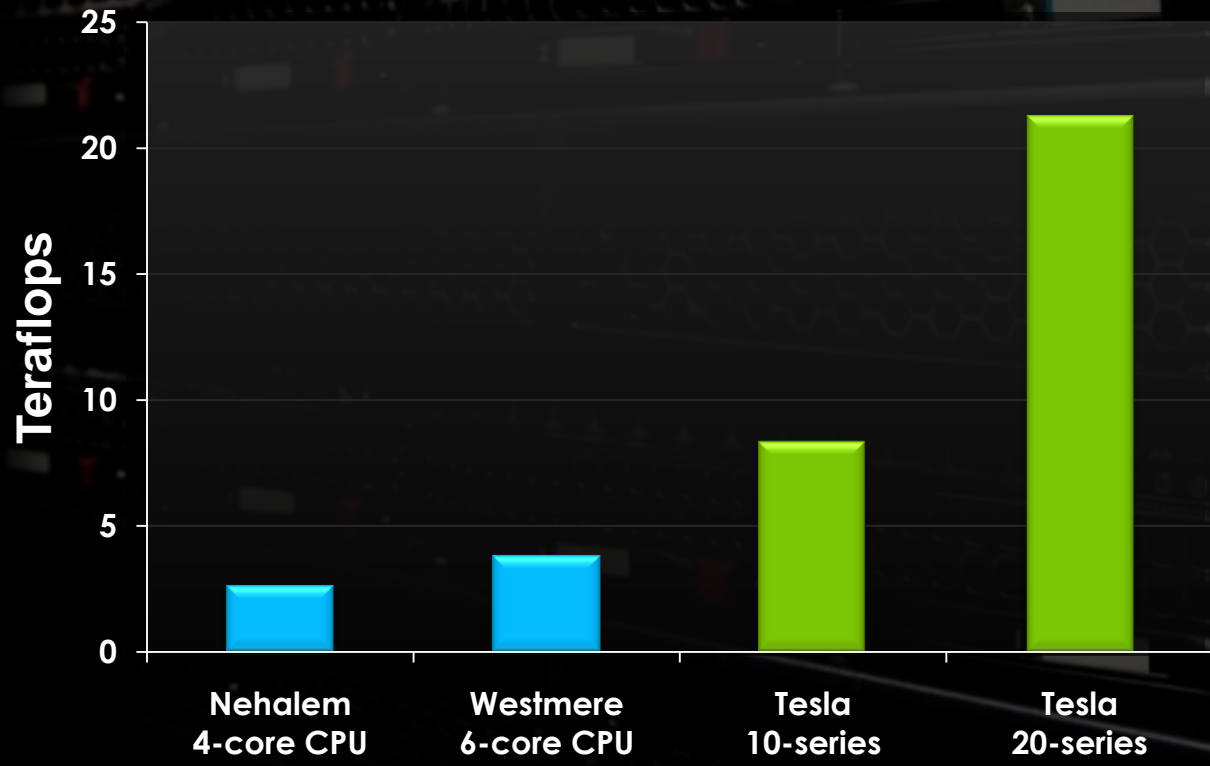
**OEM CPU Server +  
Tesla S-series 1U**



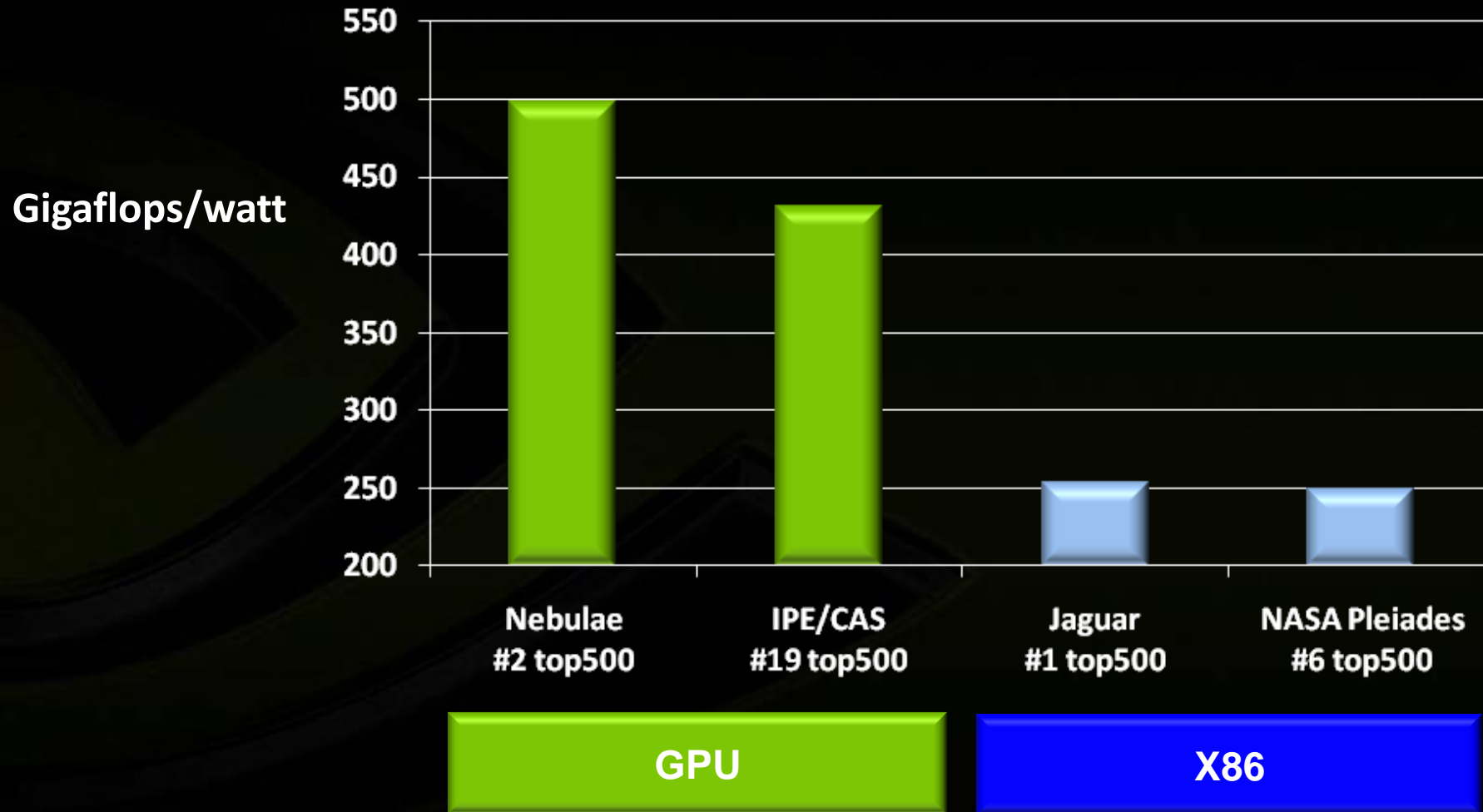
**Workstations  
2 to 4 Tesla GPUs**

# Soul of a Supercomputer

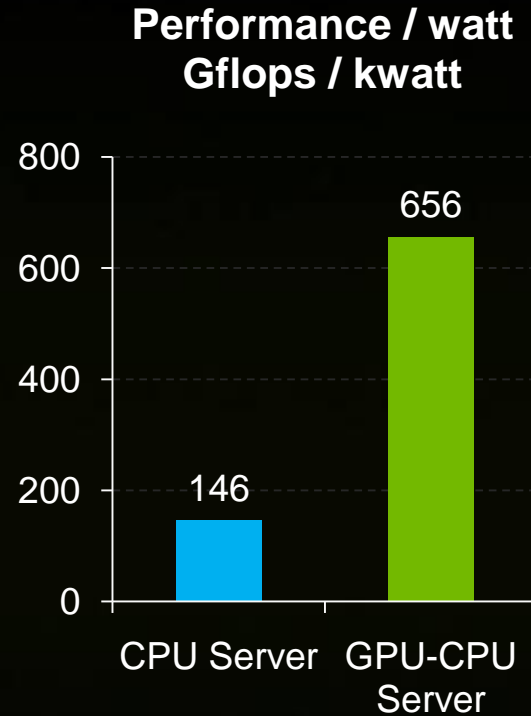
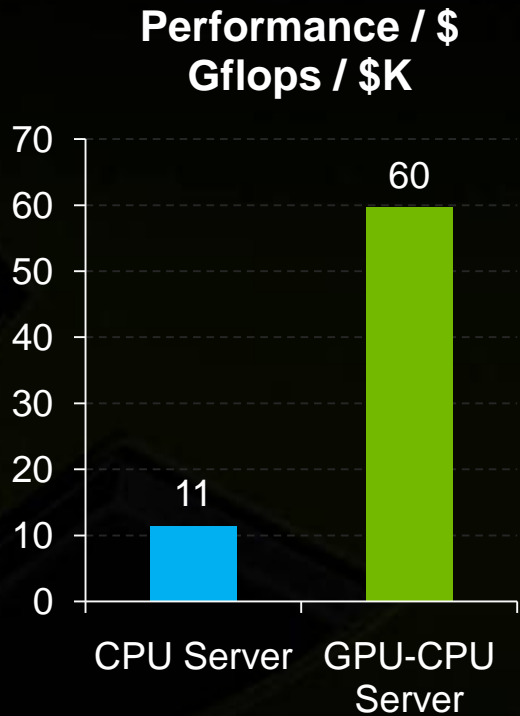
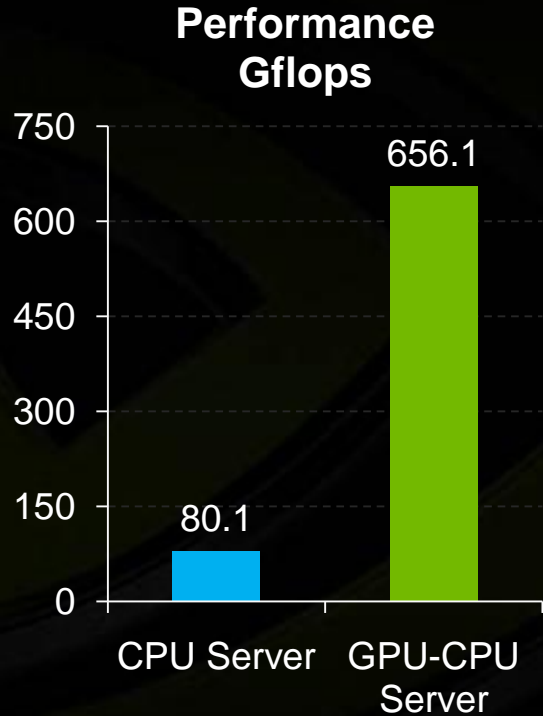
## Linpack Teraflops per Rack



# Heterogeneous = 2x Performance / Watt



# 8x Higher Linpack



**CPU 1U Server: 2x Intel Xeon X5550 (Nehalem) 2.66 GHz, 48 GB memory, \$7K, 0.55 kw**  
**GPU-CPU 1U Server: 2x Tesla C2050 + 2x Intel Xeon X5550, 48 GB memory, \$11K, 1.0 kw**



# 4.5x Lower Power & Cooling Costs

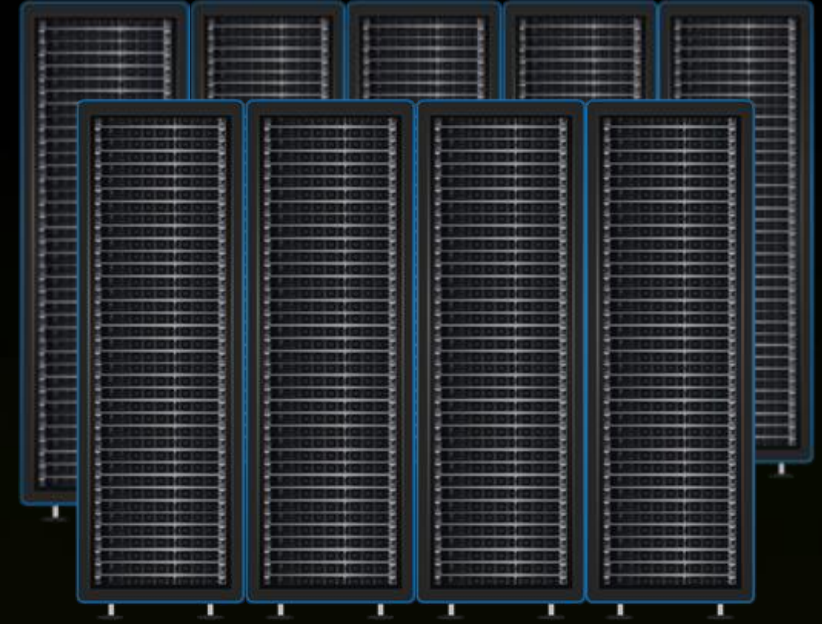
## 37 TeraFlop System : Top 150 System



2 Racks of GPU+CPUs

\$740 K

\$117 K



15 Racks of CPUs

\$3.8 M

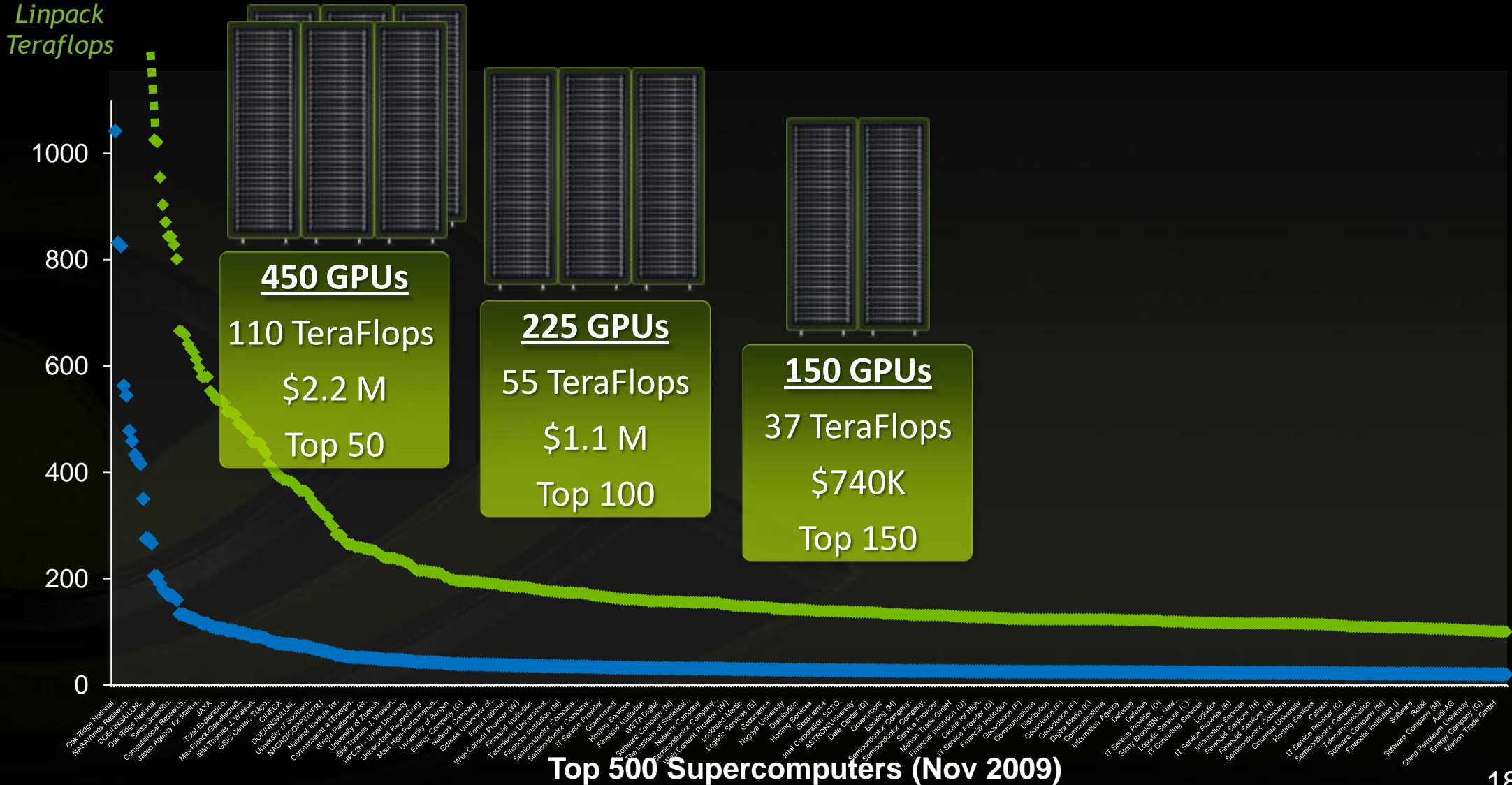
\$524 K

**7x Less Space Required**

**5x Lower Cost**

**4.5x Power Savings every Year**

# What if Every Supercomputer Had Fermi?



# 20+ Oil & Gas Companies Porting to CUDA

## Successful Customers



## Oil & Gas ISVs



## GPU vs CPU Improvements



Performance / Watt	18x - 27x	12x - 17x
Performance / Space	20x - 31x	15x - 20x
Performance / Cost	15x - 20x	10x - 12x

# Finance: 10+ Banks Porting to CUDA

## Successful Customers



**Bloomberg**

Several unannounced

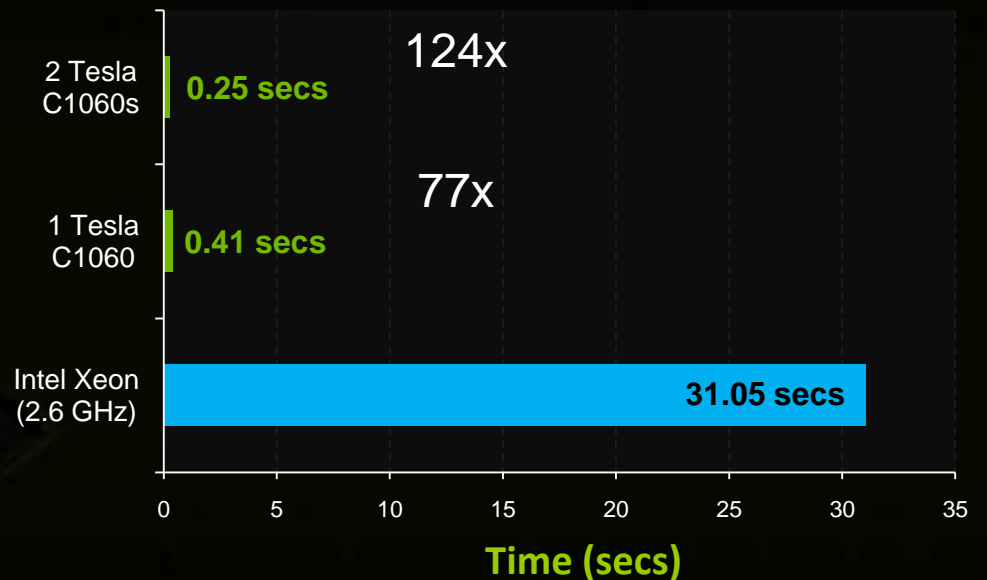
## Finance ISVs



**UnRisk**



## Derivative Pricing using SciFinance Basket Equity-Linked Structured Note



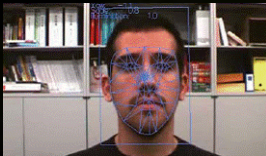
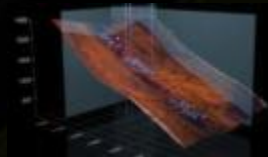
# Defense / Federal Agencies

## Software Available

- **Opportunities**

- Defense Contractors
- Federal Agencies
- Defense services

- **Speedups 10x-50x**



- **GIS**

- Manifold, PCI Geomatics, DigitalGlobe

- **Signal Processing**

- GPU VSIPL

- **MATLAB**

- GPU Plugin available

- **UAV video analysis**

- MotionDSP Ikena

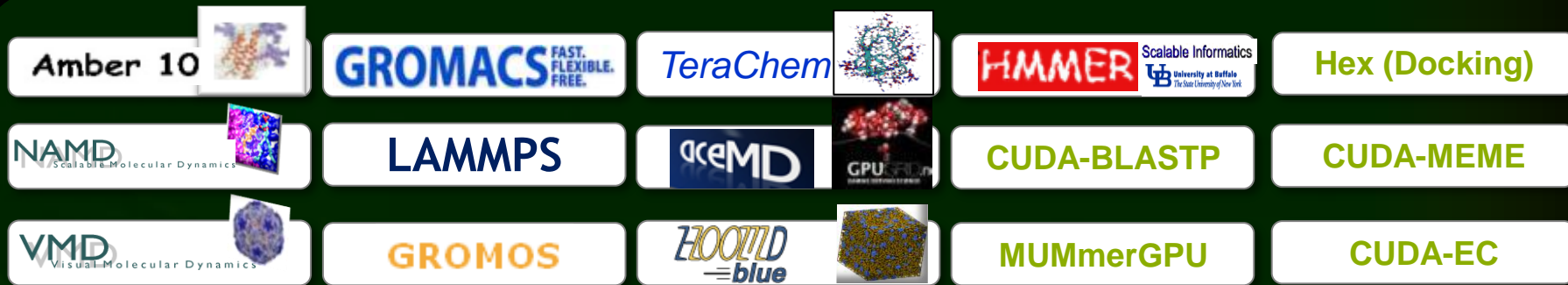
- **Virtual Prototyping**

- RealityServer

- **Surveillance, Cryptography**

# Tesla Bio WorkBench : Bio-Chemistry & Bio-Informatics

## Applications



## Community

Download,  
Documentation

Technical  
papers

Discussion  
Forums

Benchmarks  
& Configurations

## Platforms

Tesla Personal Supercomputer



Tesla GPU Clusters



# Increasing Number of CUDA Applications

	Already Available			2010					
				Q2			Q3	Q4	
Tools & Libraries	CUDA LAPACK Library	CUDA C/C++, PGI Fortran	Nsight Visual Studio IDE	Allinea Debugger	TotalView Debugger		PGI Accelerator Enhancements		
	Thrust: C++ Template Lib	Jacket: MATLAB Plugin	NPP Performance Primitives (NPP)	Platform Cluster Management	Bright Cluster Management	Mathematica	CAPS HMPP Enhancements	Mathworks MATLAB	
Oil & Gas	Seismic Analysis: ffa, HeadWave	Seismic Analysis: Geostar	Seismic City, Acceleware			Seismic Interpretation	Reservoir Simulation 1	Reservoir Simulation 2	
Bio- Chemistry	AMBER, GROMACS, GROMOS, HOOMD, LAMMPS, NAMD, VMD		BigDFT, ABINIT, TeraChem		Quantum Chem Code 1	Other Popular MD code	Quantum Chem Code 2		
Bio- Informatics	Hex Protein Docking	CUDA-BLASTP, GPU-HMMER, MUMmerGPU, MEME, CUDA-EC		Protein Docking		Short-read seq analysis			
Video & Rendering	Fraunhofer JPEG2K	OptiX Ray Tracing Engine	mental ray with iray	Main Concept Video Encoder	Elemental Video Live	3D CAD SW with iray			
Finance	NAG: RNGs	NumeriX: CounterParty Risk	Scicomp SciFinance	Risk Analysis 1	Risk Analysis 2	Credit Risk Analysis ISV	Trading Platform ISV		
CAE	AutoDesk Moldflow	OpenCurrent: CFD/PDE Library	Moldex3D	Acusim AcuSolve CFD	Structural Mechanics ISV	MSC MARC	MSC Nastran	Several CAE ISVs	
EDA	Electro-magnetics: Agilent, CST, Remcom, SPEAG		Agilent ADS Spice Simulator			Verilog Simulator	Lithography Products	SPICE Simulator	

Released  
Product

Announced  
Product

Unannounced  
Product

# Bloomberg: Bond Pricing



48 GPUs

\$144K

\$31K / year



**42x Lower Space**

**28x Lower Cost**

**38x Lower Power Cost**



2000 CPUs

\$4 Million

\$1.2 Million / year



# Oil & Gas: Seismic Processing



1

32 Tesla S1070s

~\$400 K

45 kWatts



**Equal Performance**

**31x Less Space**

**20x Lower Cost**

**27x Lower Power**



1

2000 CPU Servers

~\$8 M

1200 kWatts

# Finding a Better Shampoo

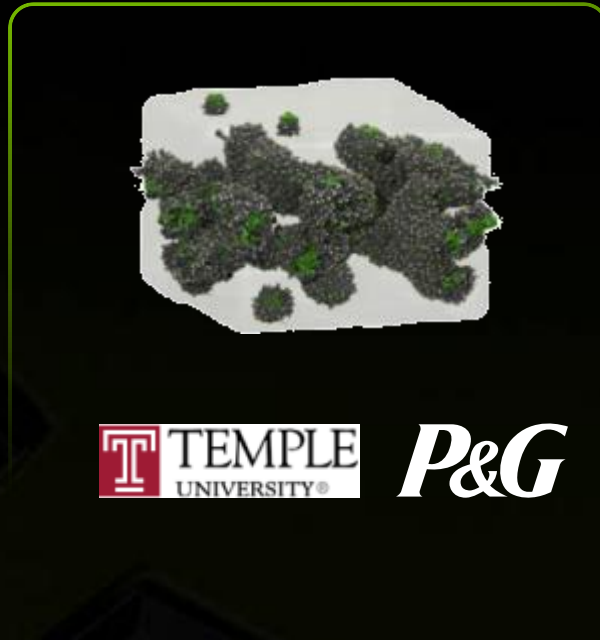


Tesla PSC

1 kWatt

\$7 K

\$2 K



Equal Performance

No Data Center Required

13x Lower System Cost

19x Lower Power & Cooling Cost



32 CPU Servers

21 kWatts

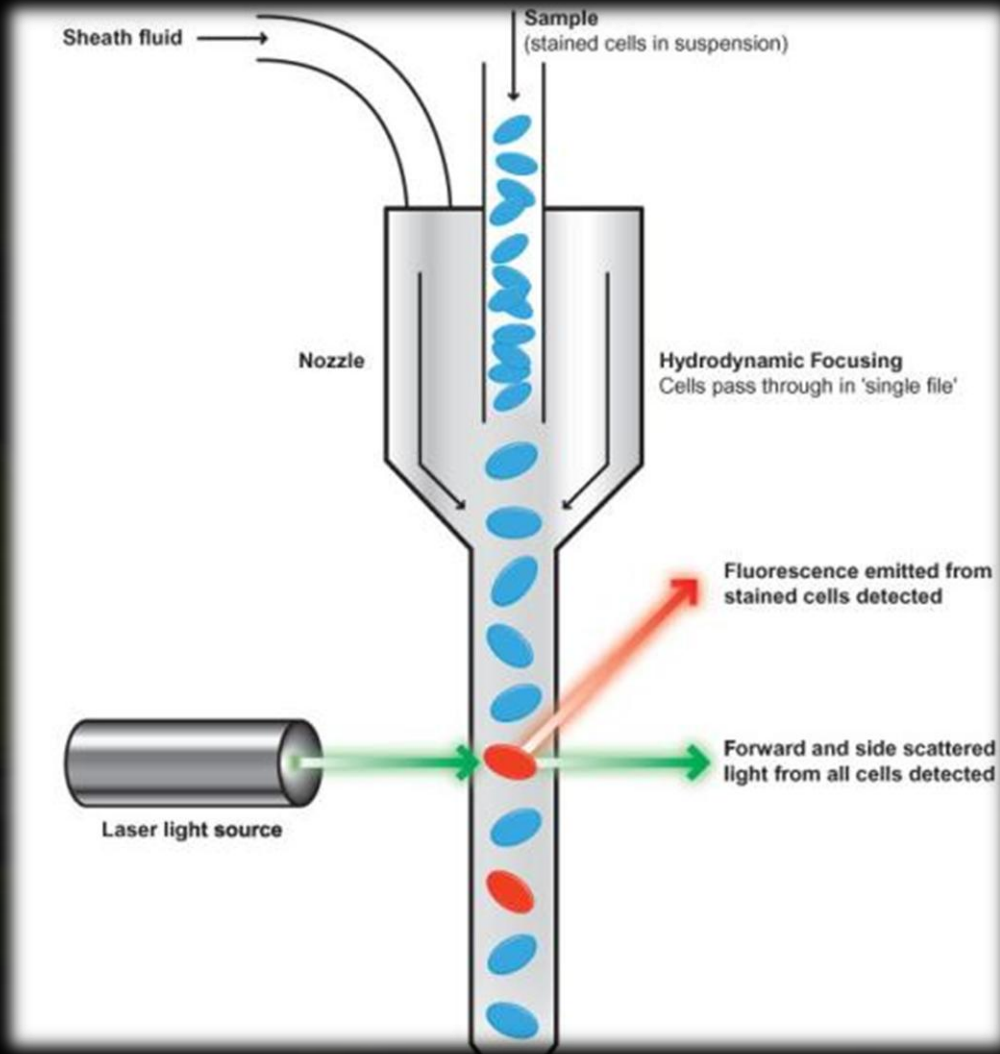
\$128 K

\$37 K

# Reducing Radiation in CT



# Flow Cytometry : Finding Cancer Cells



Flow Cytometer

# Post-Katrina Hurricane: What Google Showed



Biloxi-Ocean Springs Bridge

# What was really there : Digital Globe

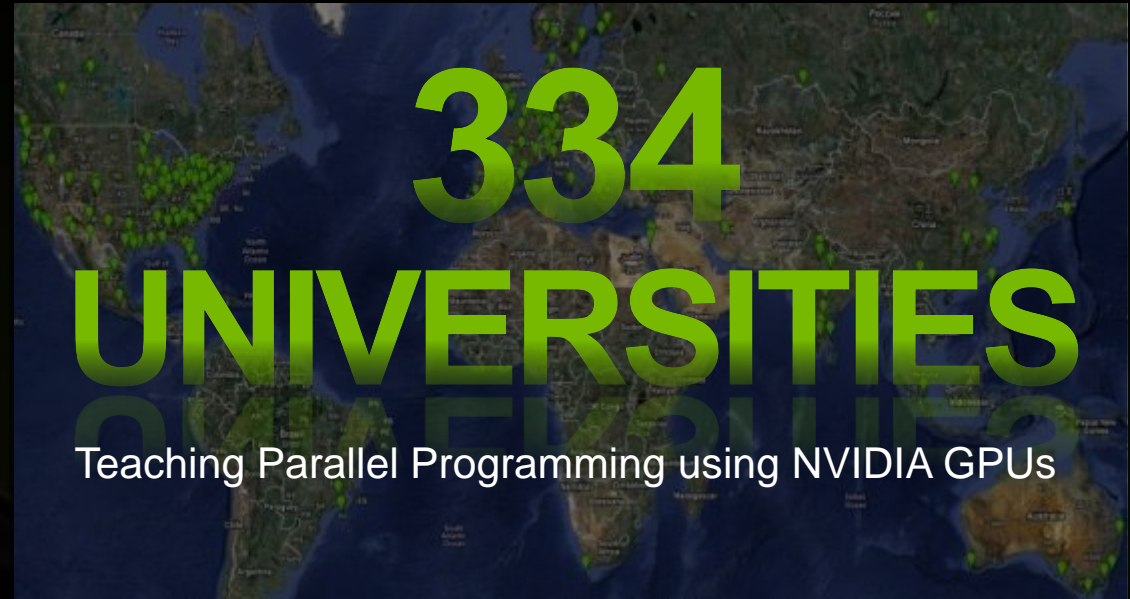
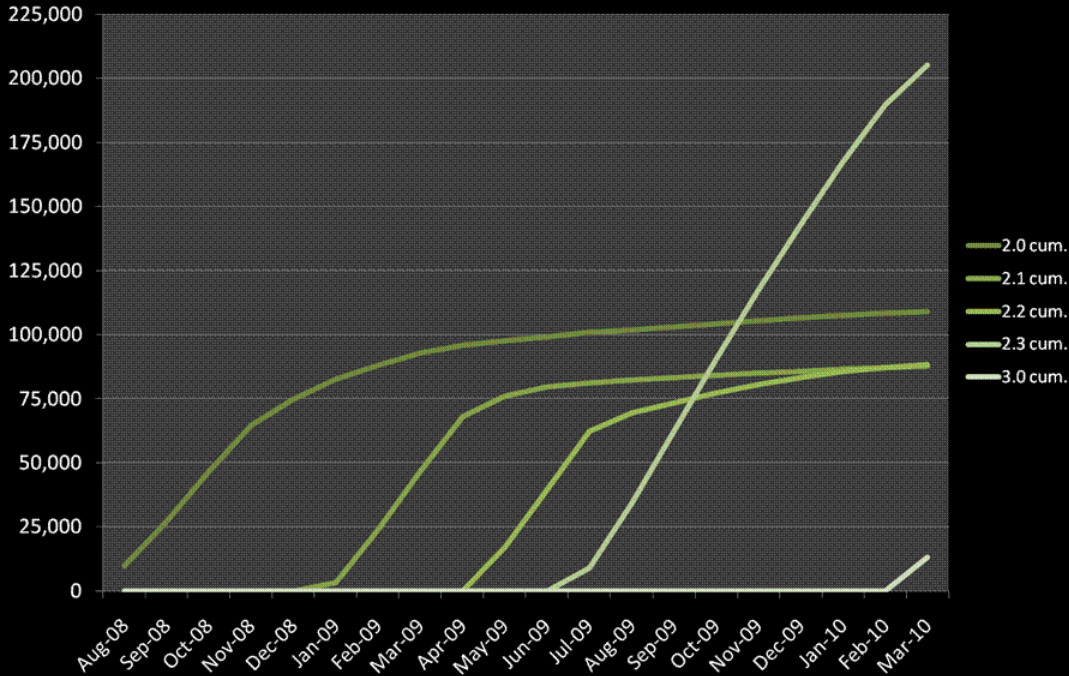


Biloxi-Ocean Springs Bridge

# Explosive Growth in CUDA Developers

200,000 CUDA Developer Downloads

Cumulative Toolkit Downloads by Version



# NVIDIA's Commitment to GPU Computing



1000+ Applications and Papers using NVIDIA GPUs



10 GPGPU Books using NVIDIA GPUs



# Programming GPUs

# Doing GPGPU Right: Combination of Hardware and Software

## GPU Computing Applications

### Libraries and Middleware

cuFFT	cuBLAS	CULA LAPACK	NPP & cuDPP	Video	PhysX Physics	OptiX Ray Tracing	mental ray iray Rendering	Reality Server 3D Web Services
-------	--------	----------------	----------------	-------	------------------	-------------------------	---------------------------------	---

C++

C

OpenCL™

Direct  
Compute

Fortran

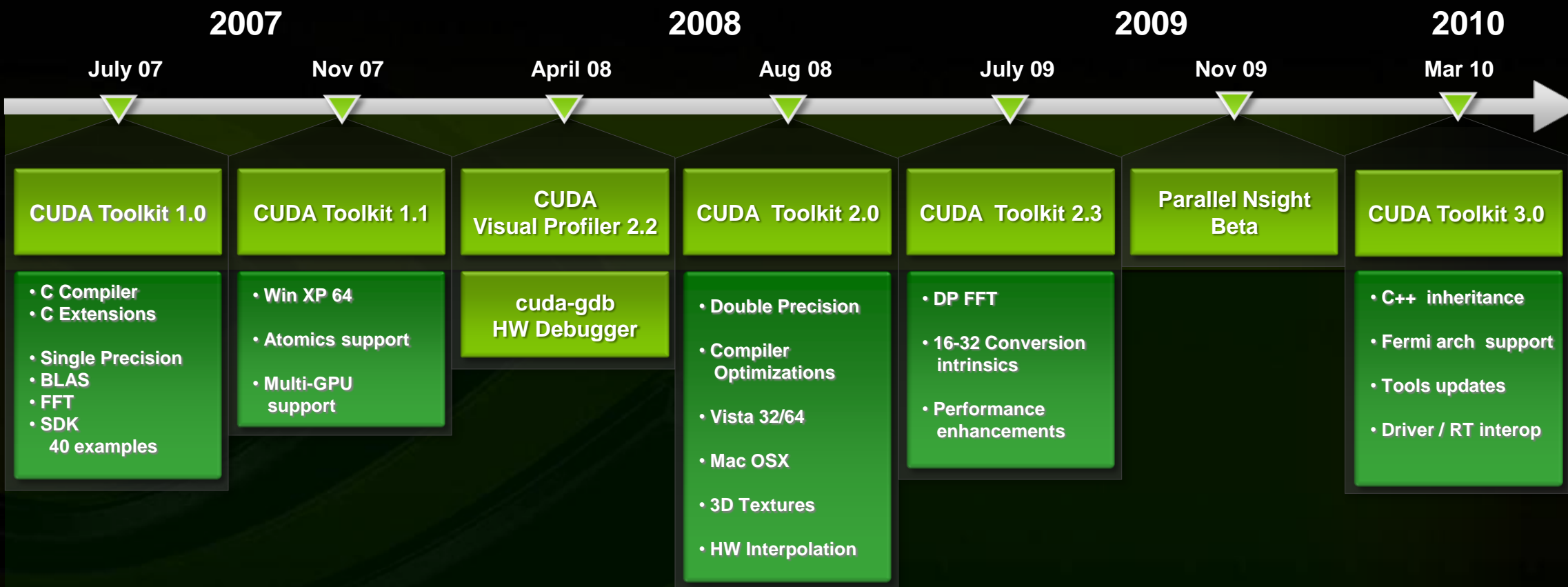
Java and  
Python



**NVIDIA GPU**

CUDA Parallel Computing Architecture

# CUDA C/C++ Continuous Innovation



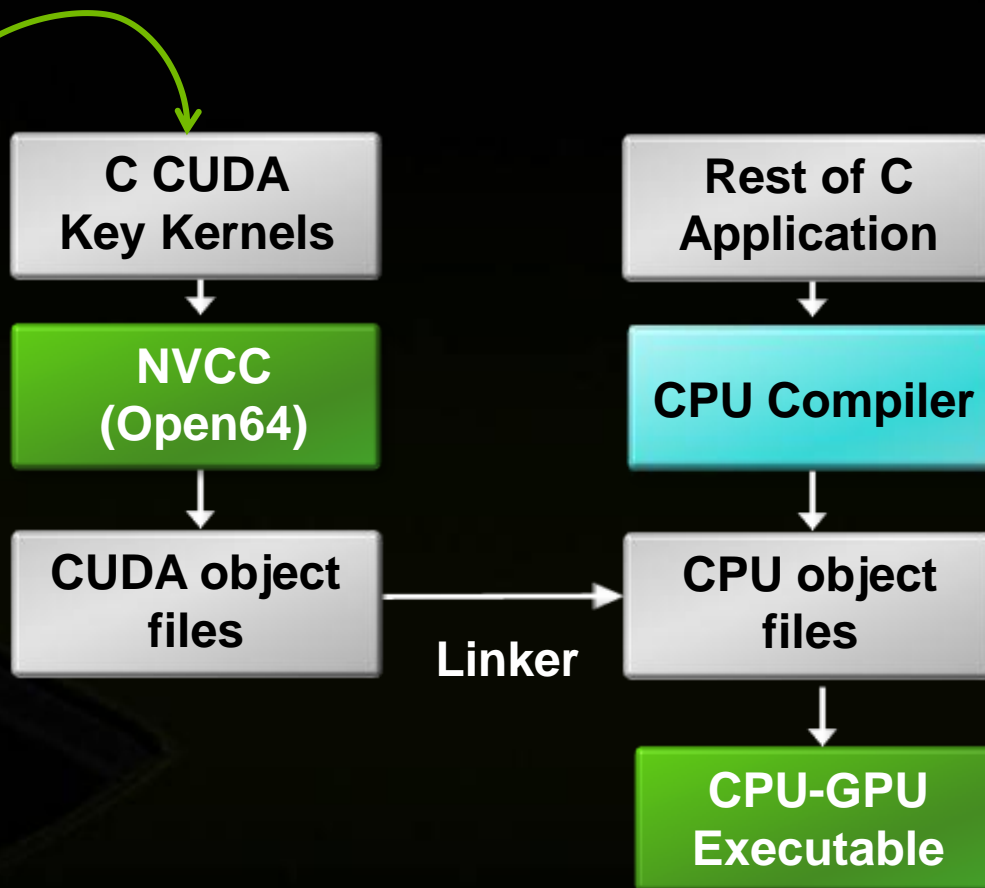
# Compiling C for CUDA Applications

```
void serial_function(... ) {  
    ...  
}  
void other_function(int ... ) {  
    ...  
}
```

```
void saxpy_serial(float ... ) {  
    for (int i = 0; i < n; ++i)  
        y[i] = a*x[i] + y[i];  
}
```

```
void main( ) {  
    float x;  
    saxpy_serial(..);  
    ...  
}
```

Modify into  
Parallel  
CUDA code



# C for CUDA : C with a few keywords

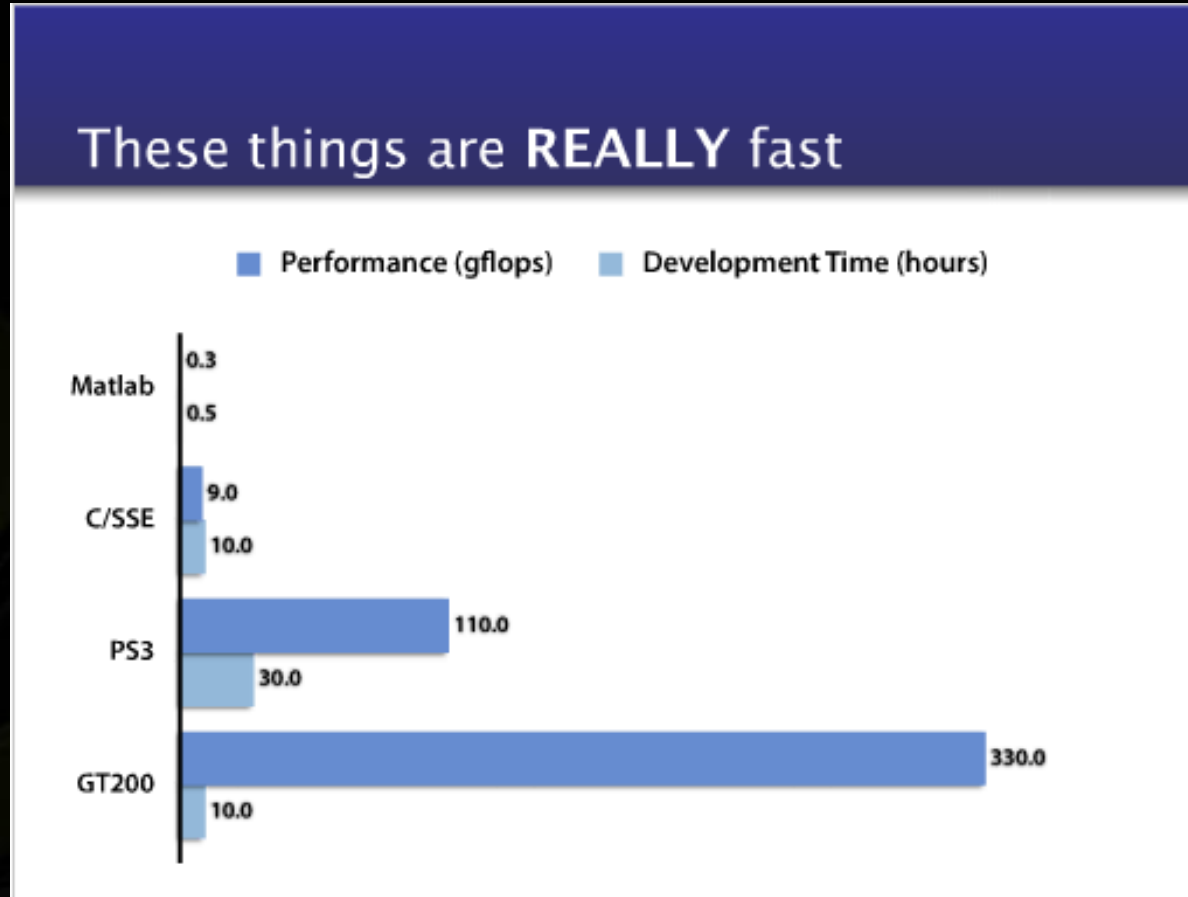
```
void saxpy_serial(int n, float a, float *x, float *y)
{
    for (int i = 0; i < n; ++i)
        y[i] = a*x[i] + y[i];
}
// Invoke serial SAXPY kernel
saxpy_serial(n, 2.0, x, y);
```

*Standard C Code*

```
__global__ void saxpy_parallel(int n, float a, float *x, float *y)
{
    int i = blockIdx.x*blockDim.x + threadIdx.x;
    if (i < n) y[i] = a*x[i] + y[i];
}
// Invoke parallel SAXPY kernel with 256 threads/block
int nblocks = (n + 255) / 256;
saxpy_parallel<<<nblocks, 256>>>(n, 2.0, x, y);
```

*Parallel C Code*

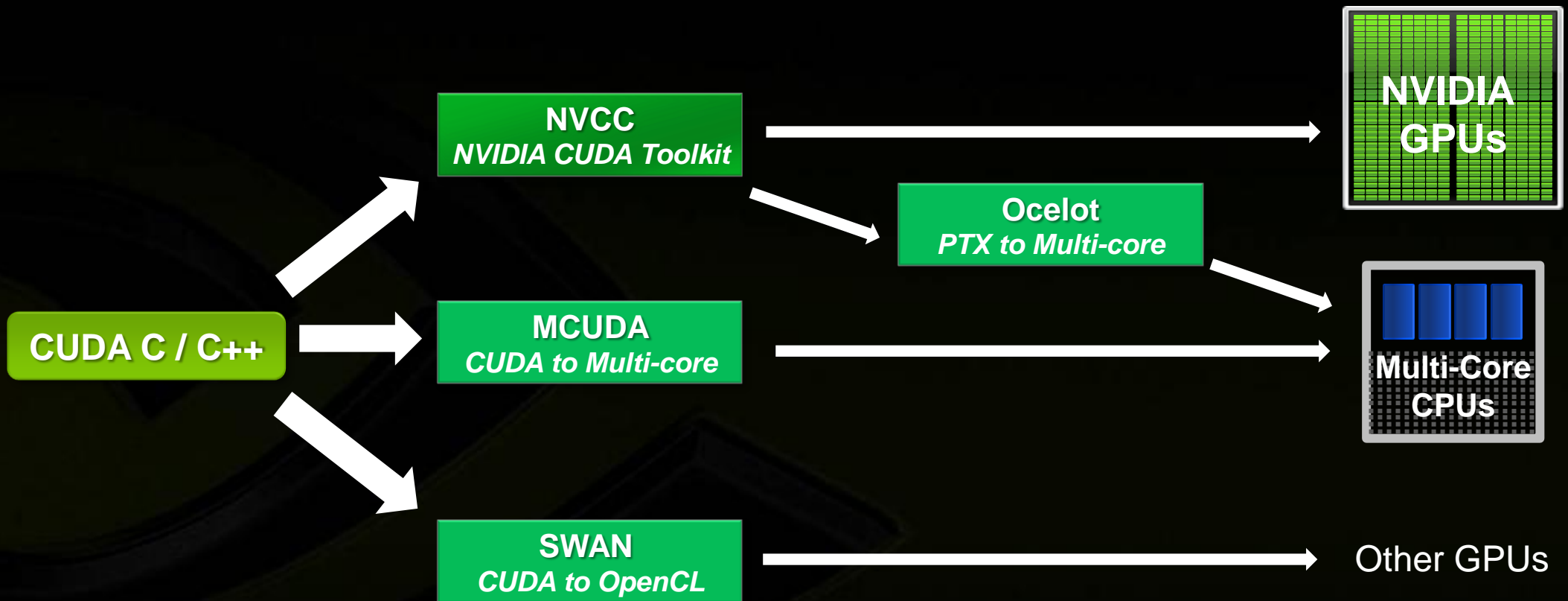
# CUDA Programming Effort / Performance



Source : MIT CUDA Course



# Targeting Multiple Platforms with CUDA



MCUDA: <http://impact.crhc.illinois.edu/mcuda.php>

Ocelot: <http://code.google.com/p/gpuocelot/>

Swan: <http://www.multiscalelab.org/swan>

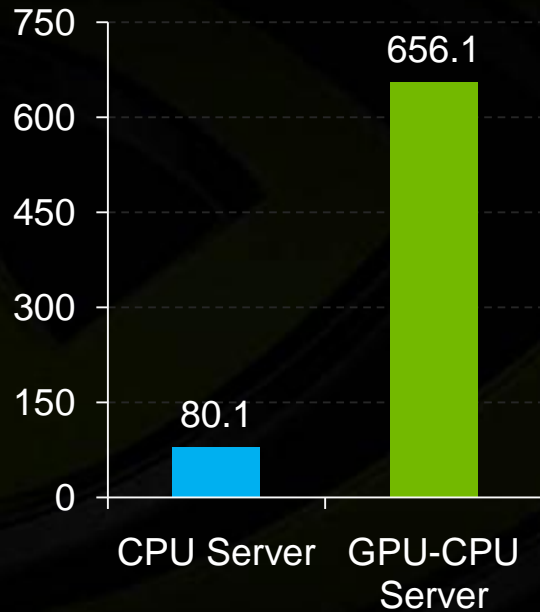


# Performance Benchmarks

# 8x Higher Linpack

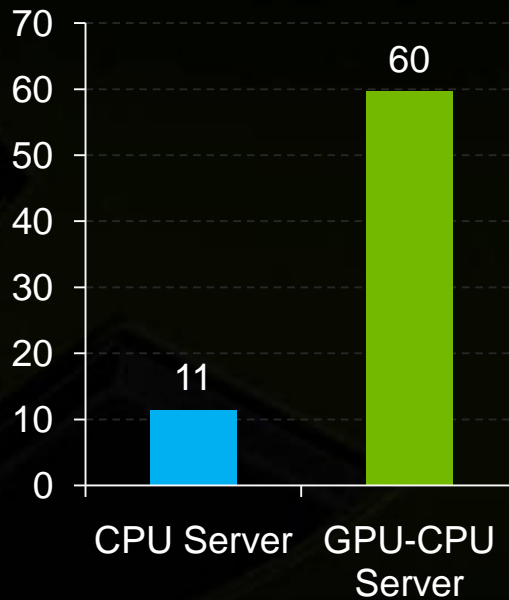
8x

Performance  
Gflops



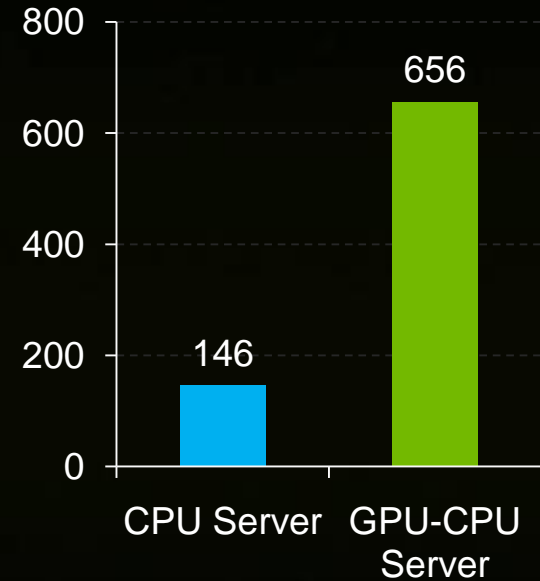
5x

Performance / \$  
Gflops / \$K



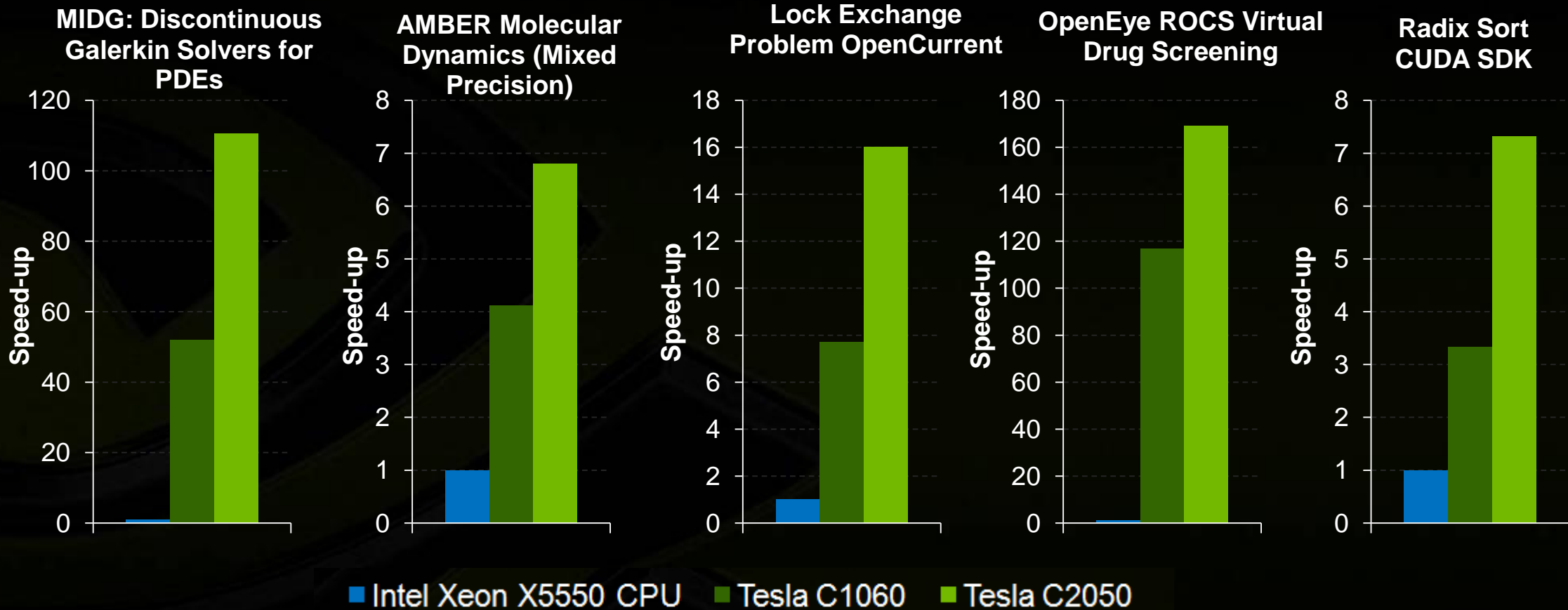
4.5x

Performance / watt  
Gflops / kwatt

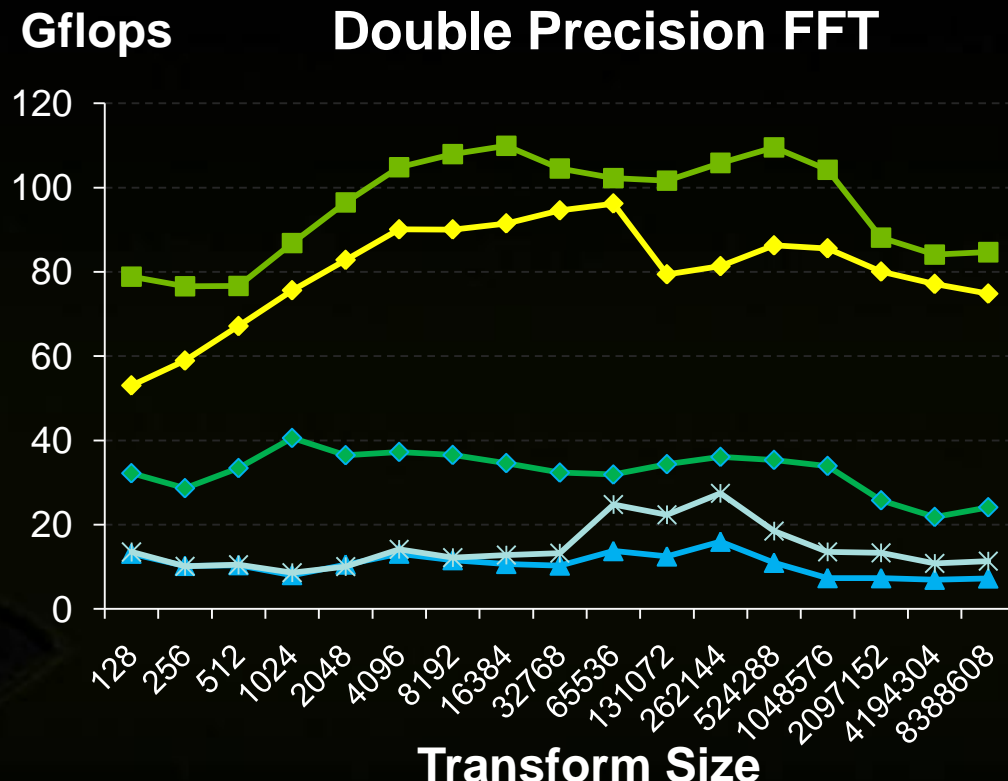
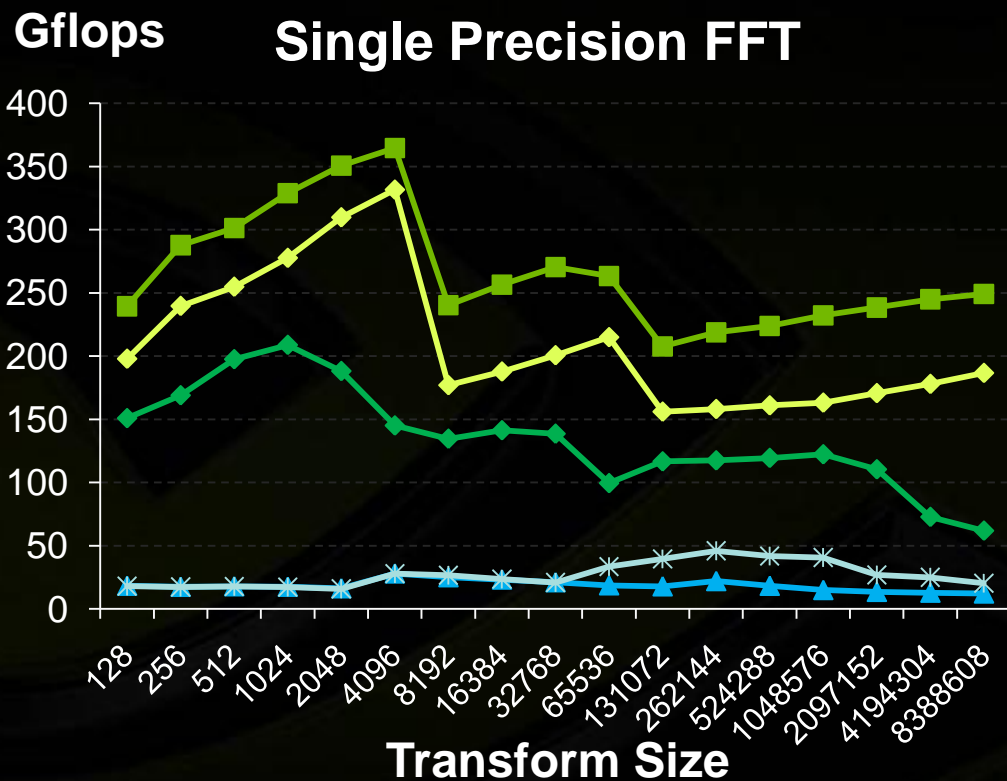


CPU 1U Server: 2x Intel Xeon X5550 (Nehalem) 2.66 GHz, 48 GB memory, \$7K, 0.55 kw  
GPU-CPU 1U Server: 2x Tesla C2050 + 2x Intel Xeon X5550, 48 GB memory, \$11K, 1.0 kw

# Performance Summary



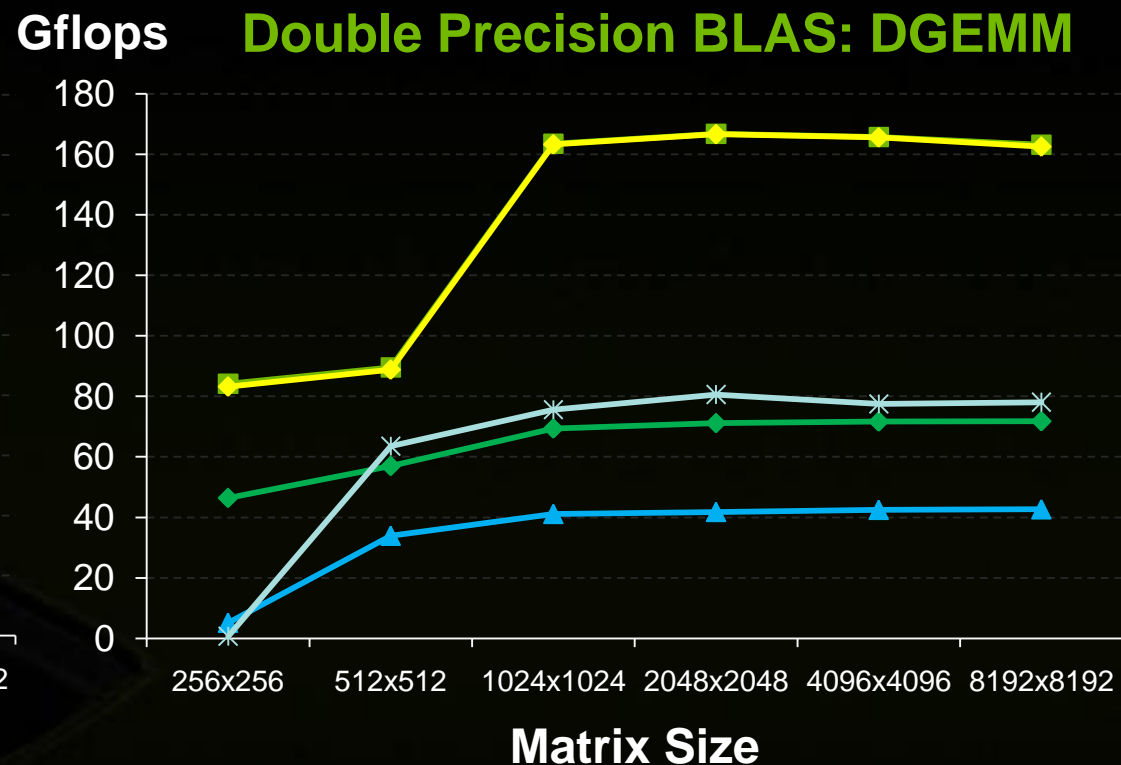
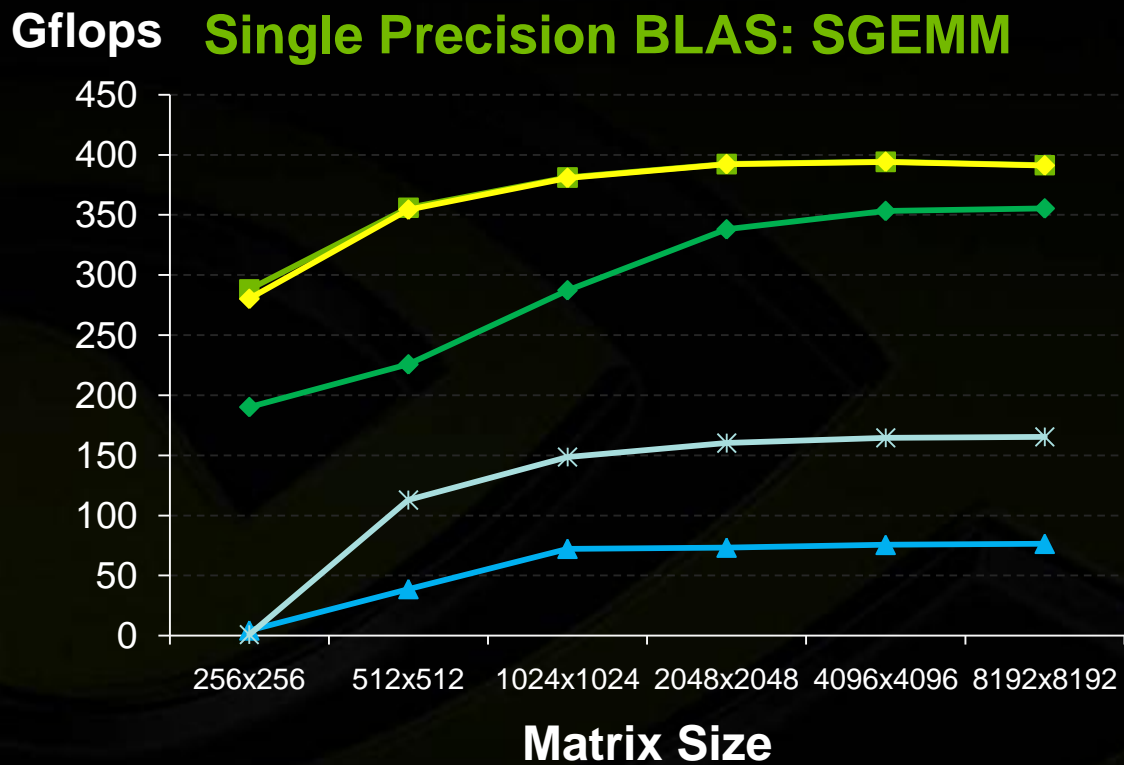
# Standard FFT Library: cuFFT 3.1



■ Tesla C2050 (ECC off)   
 ◆ Tesla C2050 (ECC on)   
 ◆ Tesla C1060   
 ▲ MKL 4 Threads   
 ✱ MKL 8 Threads

■ **cuFFT 3.1: NVIDIA Tesla C1060, Tesla C2050 (Fermi)**  
▲ **MKL 10.2.4.32: Quad-Core Intel Xeon 5550, 2.67 GHz**

# Standard BLAS Library: cuBLAS 3.1

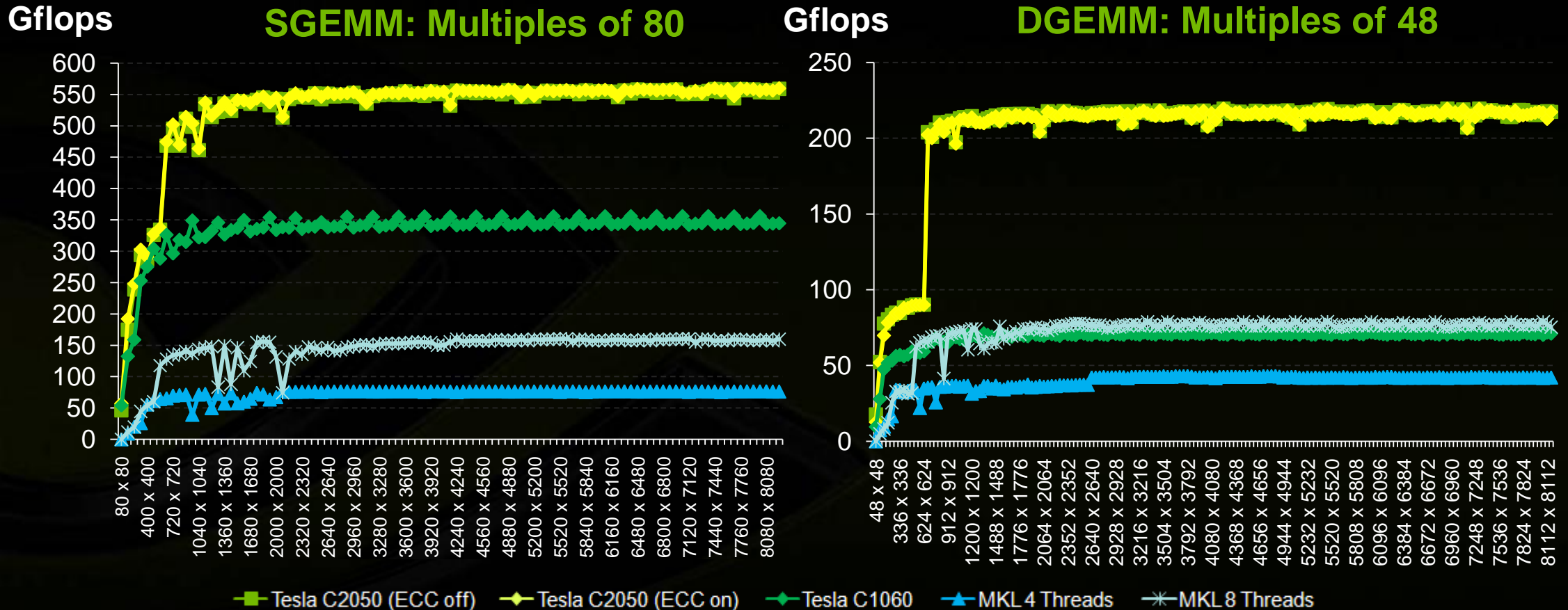


■ Tesla C2050 (ECC off) ■ Tesla C2050 (ECC on) ■ Tesla C1060 ■ MKL 4 Threads ■ MKL 8 Threads

cuBLAS 3.1: NVIDIA Tesla C1060, Tesla C2050 (Fermi)

MKL 10.2.4.32: Quad-Core Intel Xeon 5550, 2.67 GHz

# Matrix Size for Best cuBLAS Performance



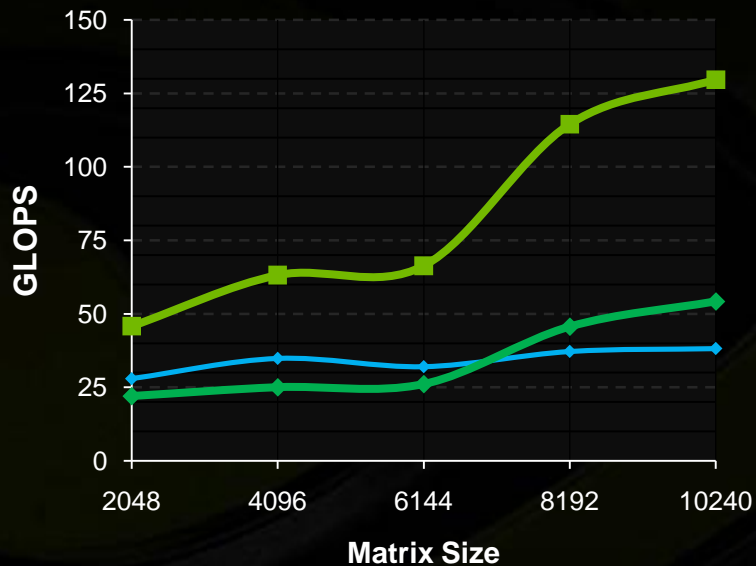
cuBLAS 3.1: NVIDIA Tesla C1060, Tesla C2050 (Fermi)

MKL 10.2.4.32: Quad-Core Intel Xeon 5550, 2.67 GHz

# CULA 1.3 LAPACK Library from EM Photonics

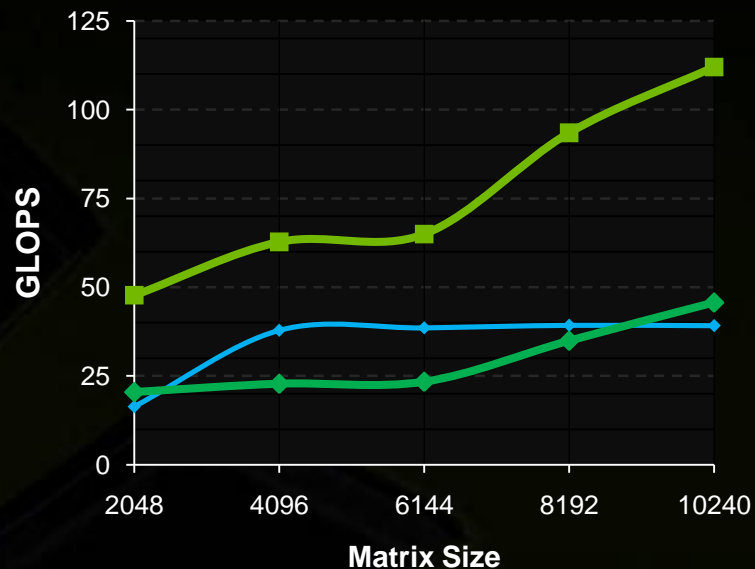
## QR Decomposition (DGEQRF)

Householder method; Operation count estimated as  $1.33N^3$



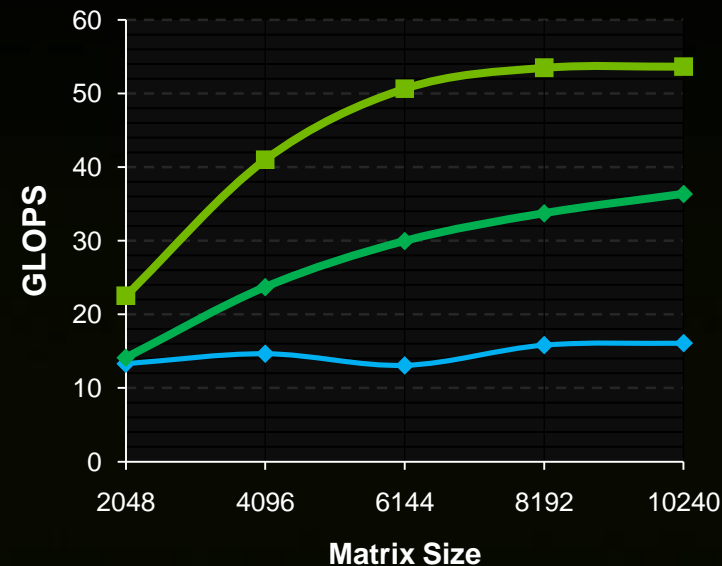
## LU Decomposition (DGETRF)

Partial pivoting; Operation count estimated as  $0.66N^3$



## Singular Value Decomposition (DGESVD)

Left & right singular vectors; Operation count estimated as  $21N^3$



Double Precision Results

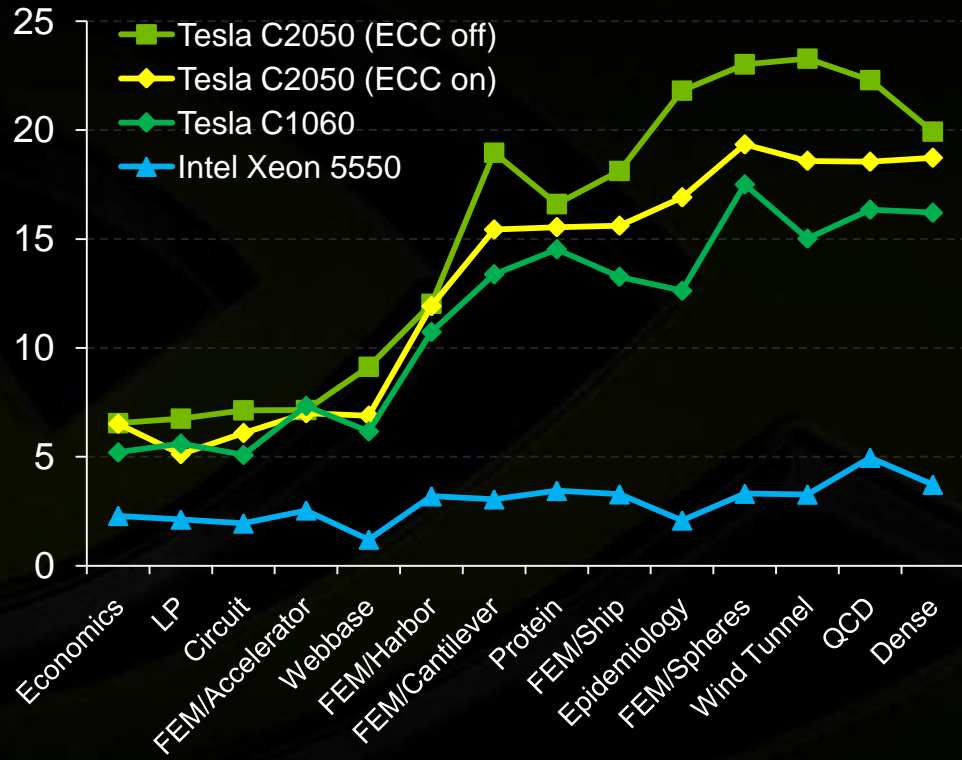
Data Courtesy: EM Photonics



# Sparse Matrix-Vector Multiplication (SpMV)

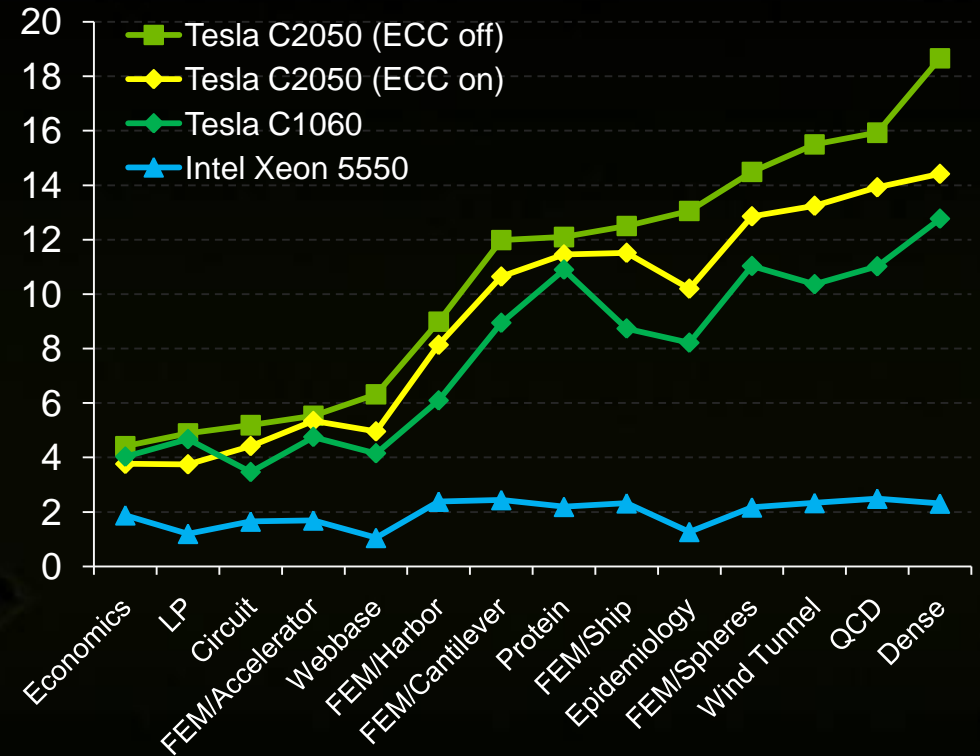
Gflops

Single Precision



Gflops

Double Precision



SpMv: CUDA 3.0, Tesla C1060 and Tesla C2050  
MKL 10.2: Intel Xeon 5550, 2.67 GHz



# Tesla GPU Computing Products

# Fermi: The Computational GPU

## Performance

- 7x Double Precision of CPUs
- IEEE 754-2008 SP & 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

# NVIDIA Tesla GPU Computing Products

## Server Module



Tesla M2070 /  
Tesla M2050



Tesla M1060

## 1U Systems



Tesla S2050



Tesla S1070

## Workstation Boards



Tesla C2070 /  
Tesla C2050



Tesla C1060

	Server Module		1U Systems		Workstation Boards	
GPUs	1 T20 GPU	1 T10 GPU	4 T20 GPUs	4 T10 GPUs	1 T20 GPU	1 T10 GPU
Single Precision	1030 GFlops	933 GFlops	4120 GFlops	4140 GFlops	1030 Gflops	933 GFlops
Double Precision	515 Gflops	78 GFlops	2060 GFlops	346 GFlops	515 Gflops	78 GFlops
Memory	6 GB / 3 GB	4 GB	12 GB (S2050)	16 GB 4 GB / GPU	6 GB / 3 GB	4 GB
Mem BW	148.4 GB/s	102 GB/s	148.4 GB/s	102 GB/s	144 GB/s	102 GB/s



“In testing our key applications, the Tesla GPUs delivered speed-ups that we had never seen before, sometimes even orders of magnitude”

**Satoshi Matsuoka**

Professor  
Tokyo Institute of Technology

“I believe history will record Fermi as a significant milestone.”

**Dave Patterson**

Director Parallel Computing Research Laboratory, U.C. Berkeley  
Co-Author of Computer Architecture: A Quantitative Approach



“Future computing architectures will be hybrid systems with parallel-core GPUs working in tandem with multi-core CPUs”



**Jack Dongarra**

Professor, University of Tennessee  
Author of Linpack



Mass market adoption

---

Applications

---

Scaling