



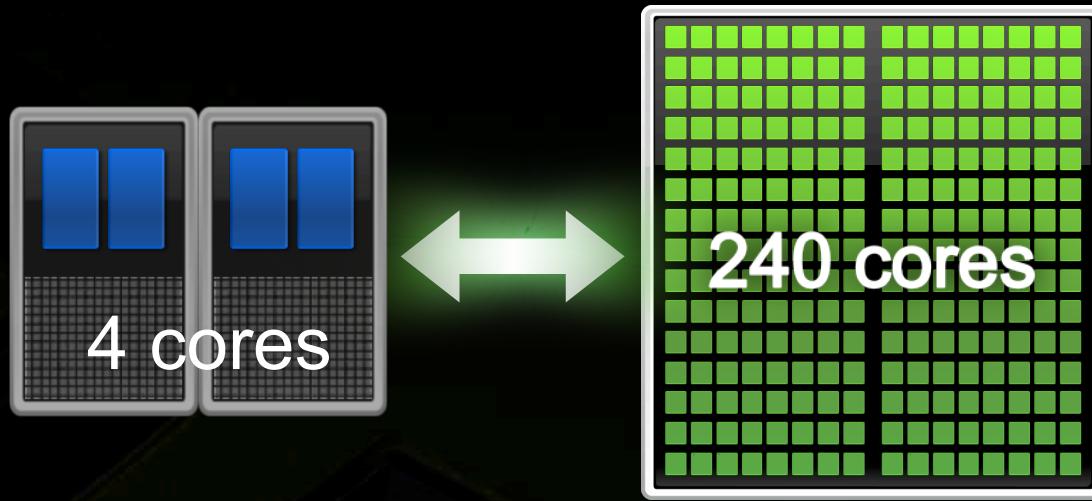
The Massively Parallel Computing Revolution

Ross Cunniff, Director of Desktop GPU Software

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

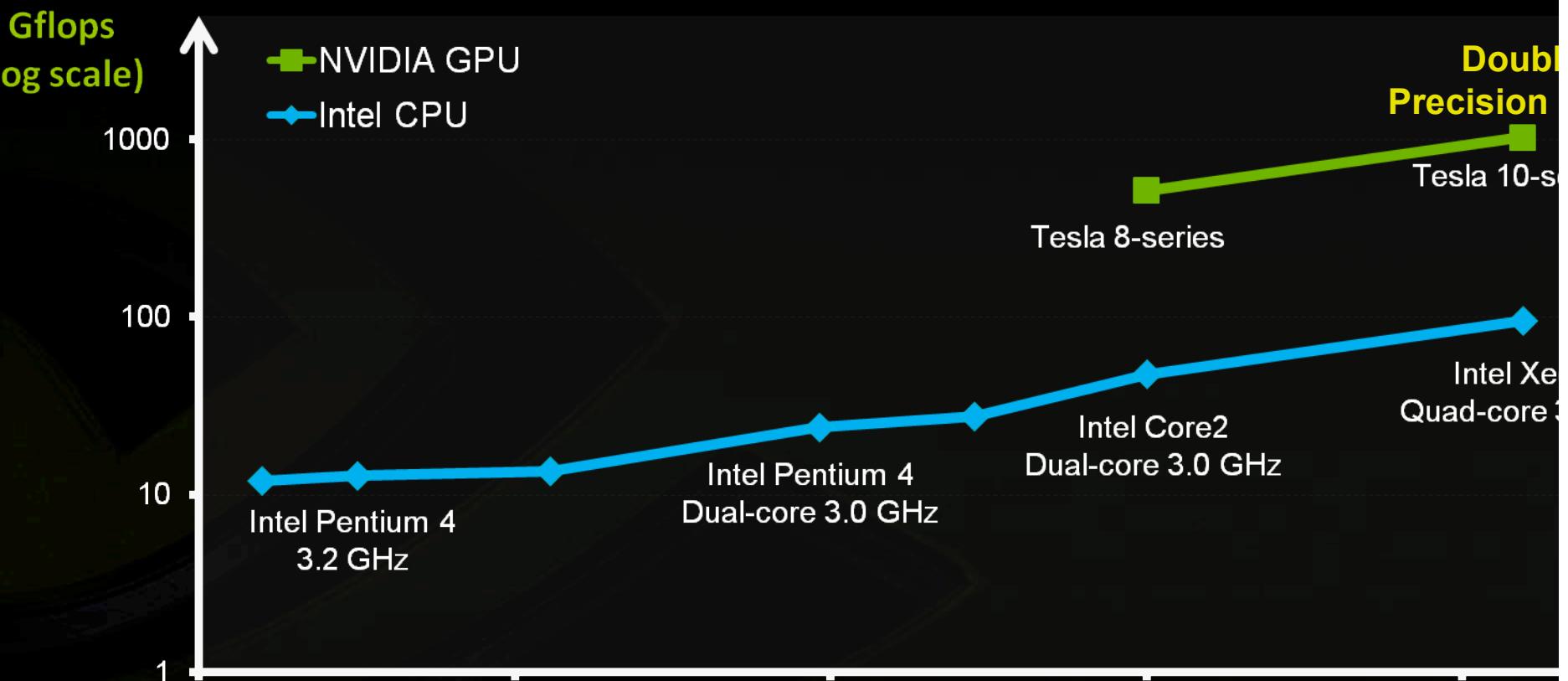
FRACTAL workshop, April 2009

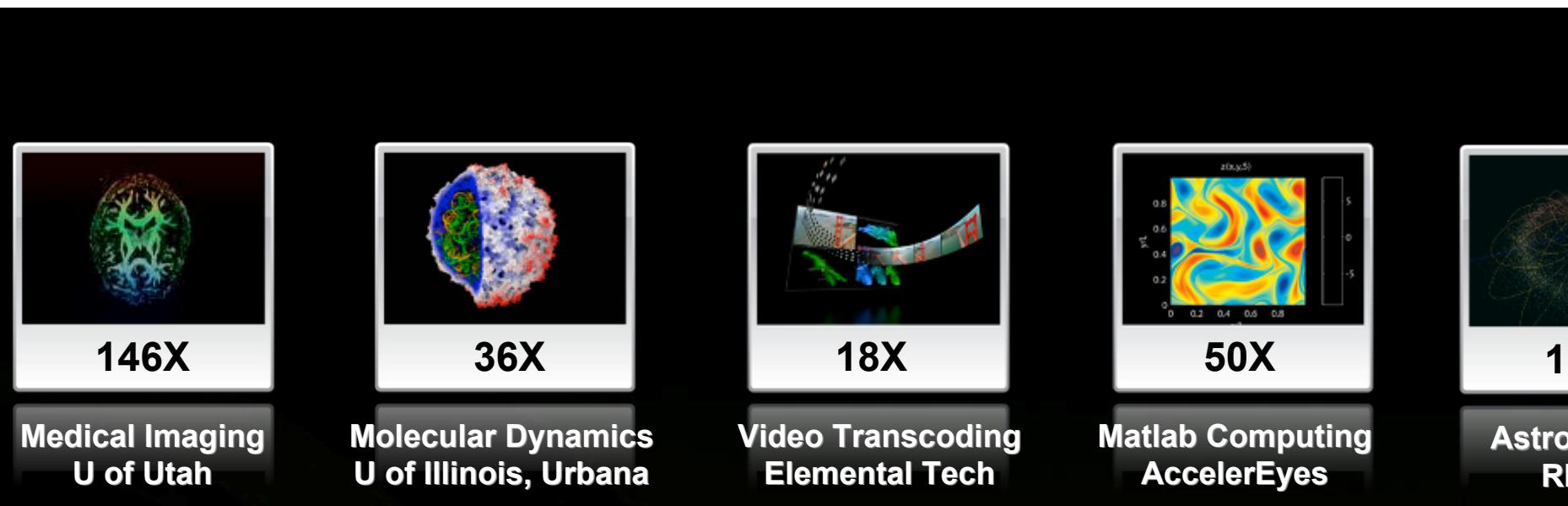
Move to Massive Parallelism



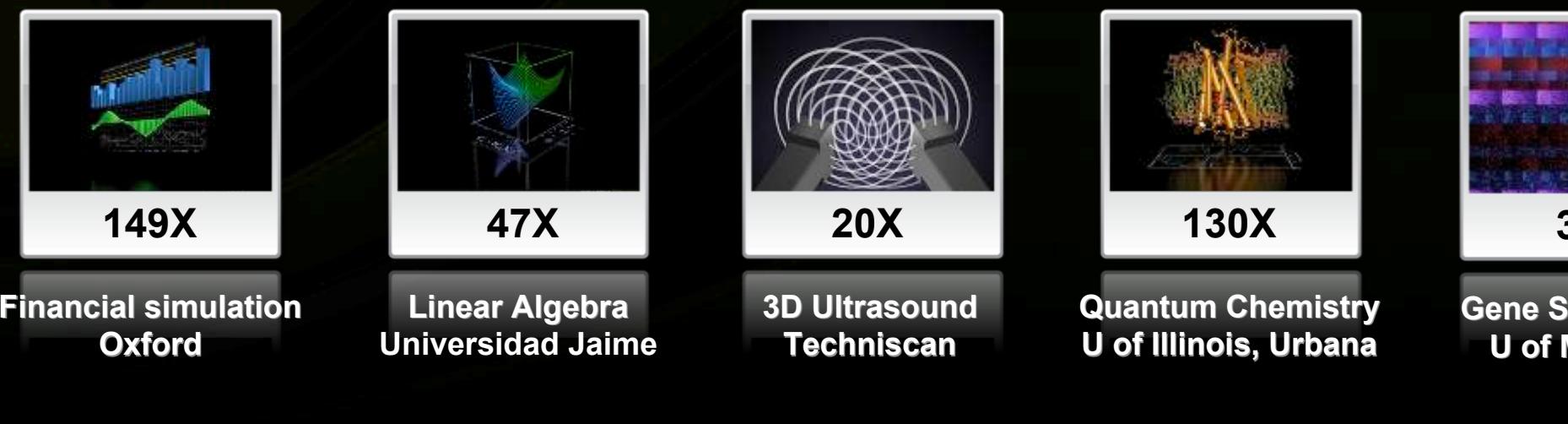
Computing with CPU + GPU
Heterogeneous Computing

Computation Discontinuity

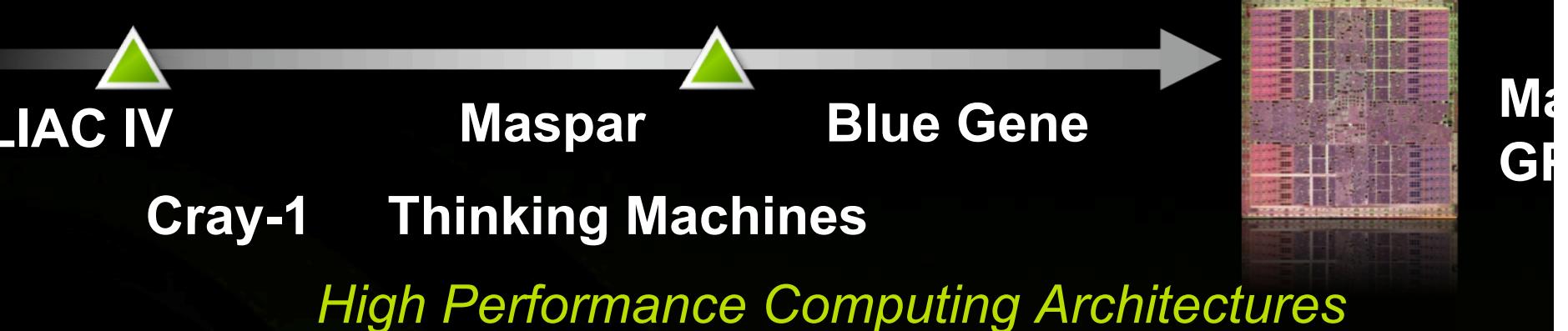


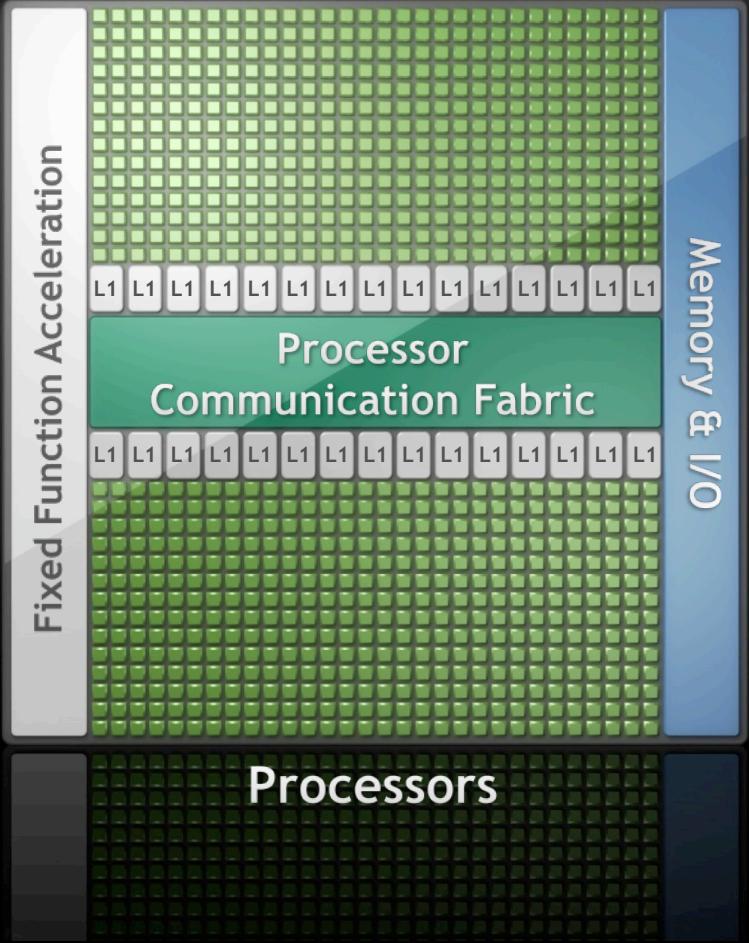


50x – 150x



Parallel vs Sequential Architecture Evolution





NVIDIA Tesla 10-Series GPU

Massively parallel, many core architecture

240 Processor Cores

1 Teraflops - 1,000 times Cray X-MP

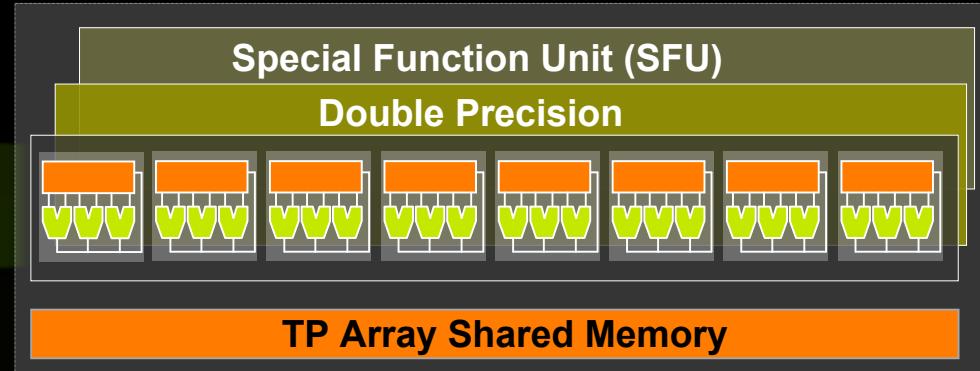
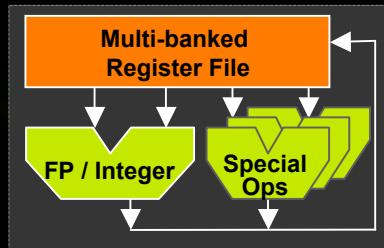
IEEE Compliant Double Precision Floating Point

Designed for Scientific Computing

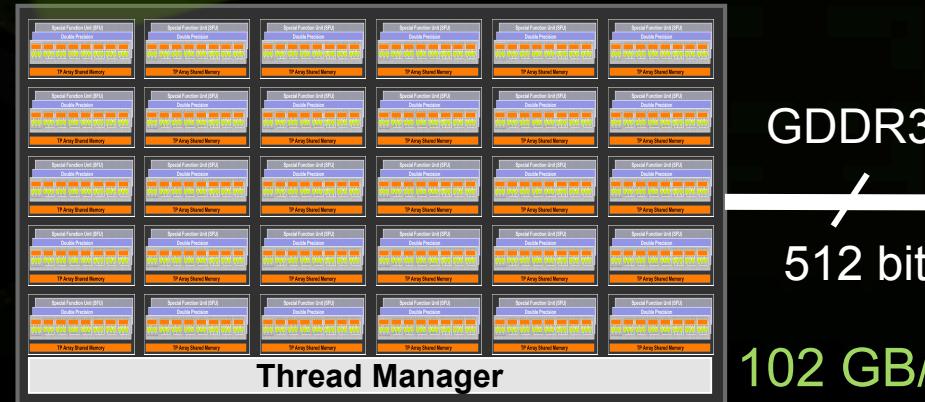
Tesla T10 GPU: 240 Processor Cores

Thread Processor Array (TPA)

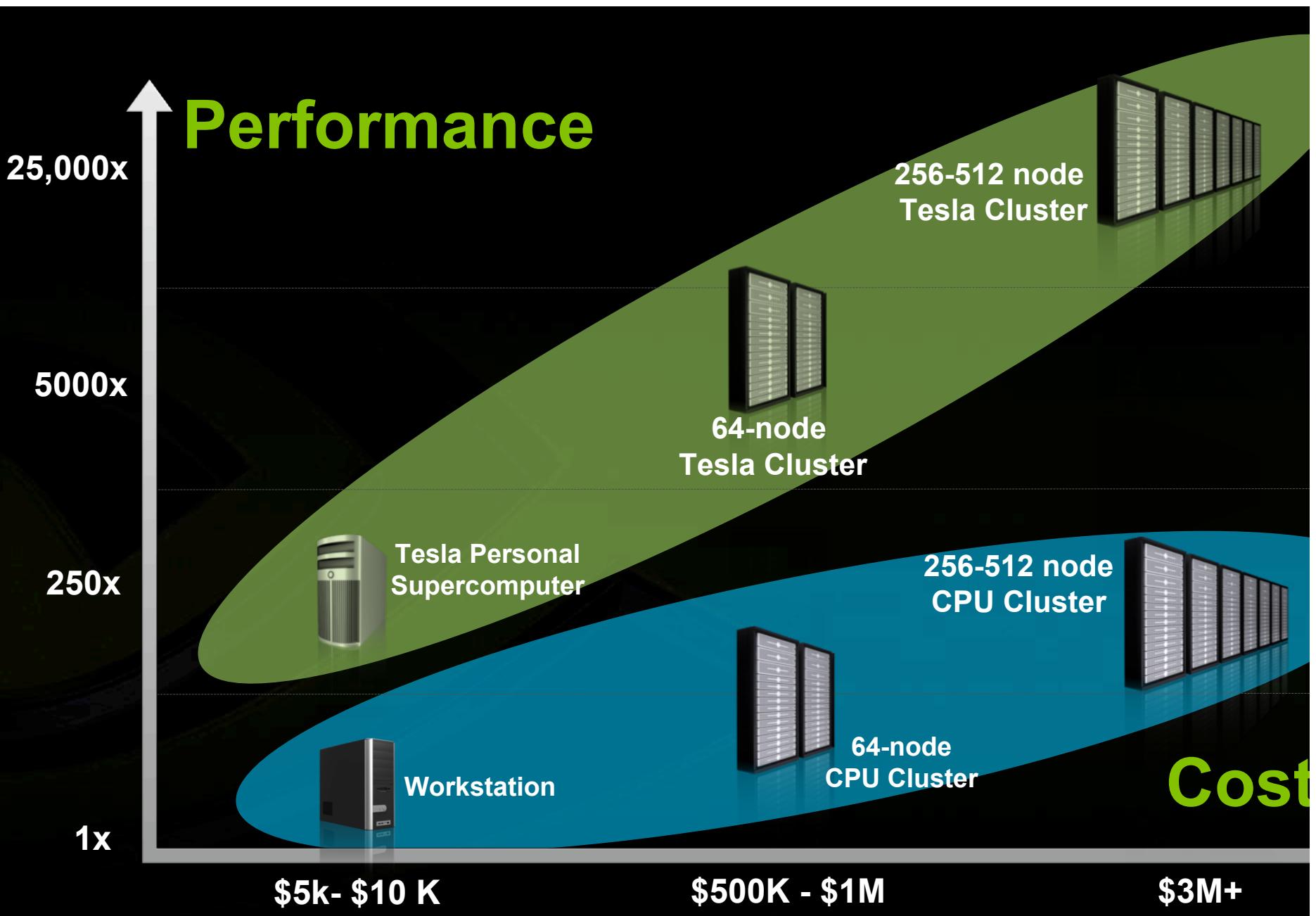
Thread Processor (TP)



- Processor core has
 - Floating point / Integer unit
 - Move, compare, logic, branch unit
- IEEE 754 floating point
 - Single and Double
- 102 GB/s high-speed interface to memory



30 TPAs = 240 Processors



5000+ Customers / ISVs

Life Sciences & Medical Equipment	Productivity / Misc	Oil and Gas	EDA	Finance	CAE / Mathematics	
Max Planck FDA barts Research Medtronic AGC olved machines mith-Waterman NA sequencing AutoDock NAMD/VMD olding@Home oward Hughes Medical RIBI Genomics	GE Healthcare Siemens Techniscan Boston Scientific Eli Lilly Silicon Informatics Stockholm Research Harvard Delaware Pittsburg ETH Zurich Institute Atomic Physics	CEA NCSA WRF Weather Modeling OptiTEx Tech-X Elemental Technologies Dimensional Imaging Manifold Digisens General Mills Rapidmind Rhythm & Hues xNormal Elcomsoft LINZIK	Hess TOTAL CGG/Veritas Chevron Headwave Acceleware Seismic City P-Wave Seismic Imaging Mercury Computer ffA	Synopsys Nascentric Gauda CST Agilent	Symcor Level 3 SciComp Hanweck Quant Catalyst RogueWave BNP Paribas	AccelerEyes MathWorks Wolfram National Instruments Ansys Access Analytics Tech-x RIKEN SOFA Renault Boeing

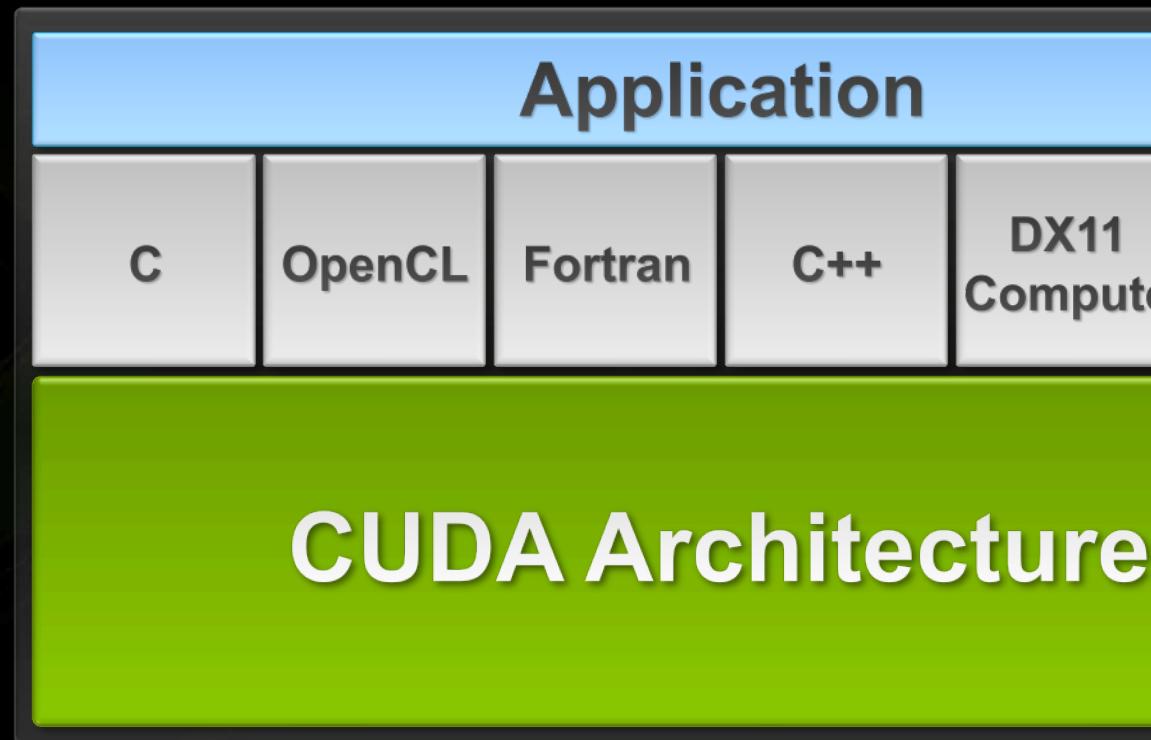
CUDA

CUDA Parallel Computing Architecture

Parallel computing architecture
and programming model

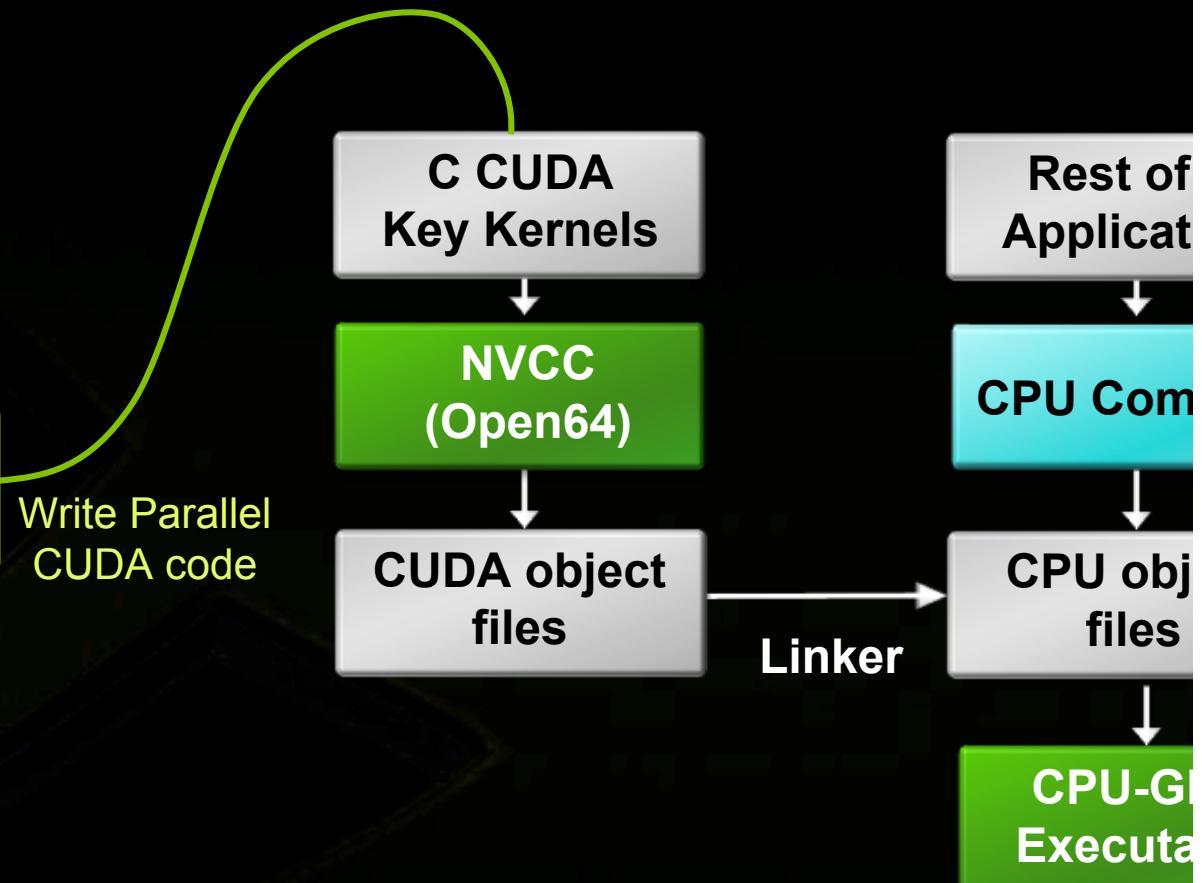
Includes a C compiler plus
support for OpenCL and
DX11 Compute

Architected to natively support
all computational interfaces
(standard languages and APIs)



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(..);  
    ...  
}
```



C for CUDA : C with 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

```
__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);
```

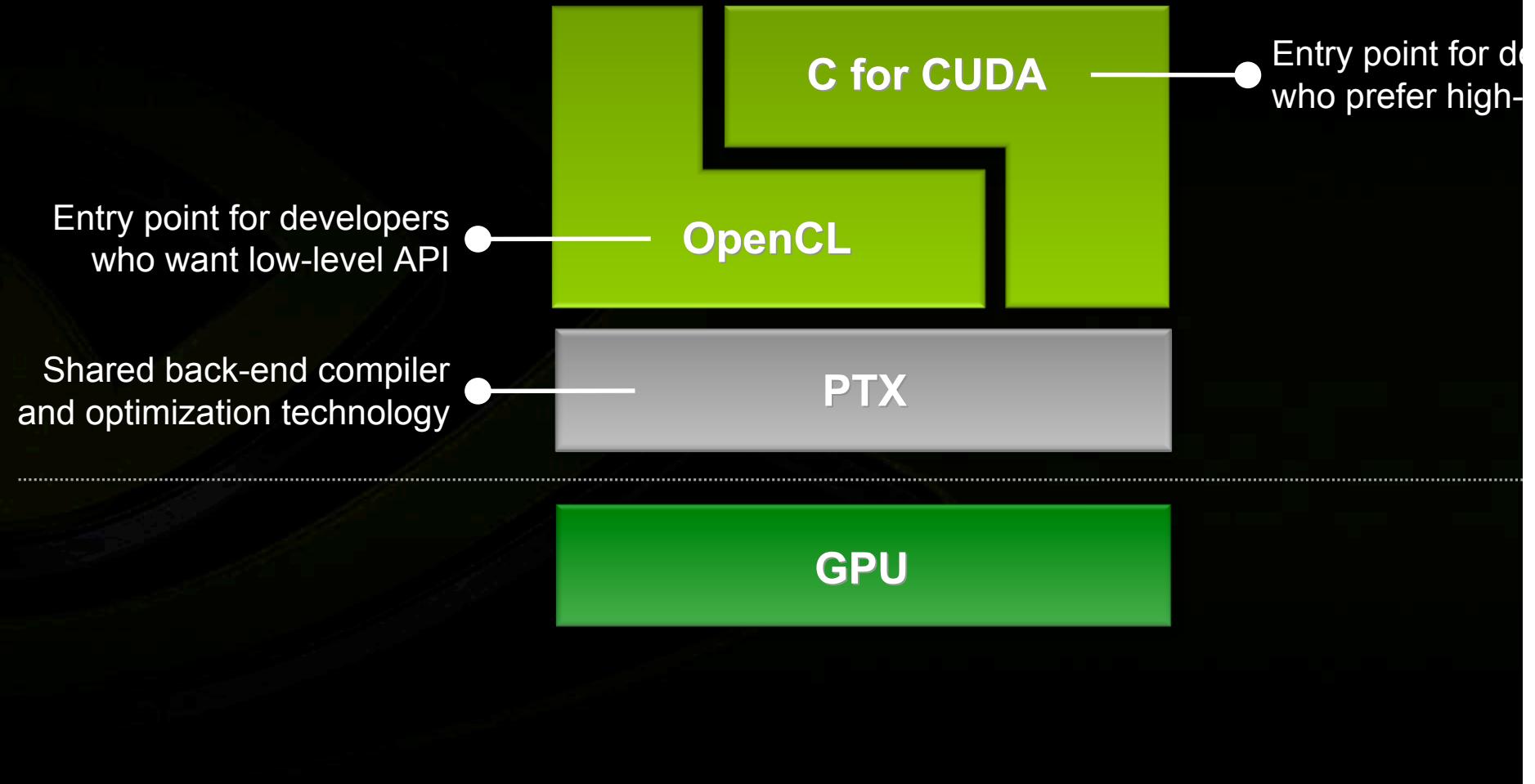
Parallel C

Support for Fortran, Python

- PGI Fortran to CUDA compiler (alpha release)
- Fortran Wrapper for CUDA
 - FLAGON : Library of numerics for Fortran 95
- Full Fortran support in Q4 2009
- Also Available
 - PyCUDA : Python Wrapper for CUDA
 - Java wrappers, .NET integration

Available from www.nvidia.com/tesla

NVIDIA C for CUDA and OpenCL



Different Programming Styles

- **C for CUDA**
 - C with parallel keywords
 - C runtime that abstracts driver API
 - Memory managed by C runtime
 - Generates PTX
- **OpenCL**
 - Hardware API - similar to OpenGL and CUDA driver API
 - Programmer has complete access to hardware device
 - Memory managed by programmer
 - Generates PTX

CUDA Facts

900+ Research Papers

115+ universities teaching CUDA



NVIDIA CUDA ZONE

USA -

DOWNLOAD CUDA WHAT IS CUDA CUDA U DEVELOPING WITH CUDA

LATEST CUDA NEWS NVIDIA Tesla Makes Personal SuperComputing A Reality

Dense Compressed/Hierarchical Linear System Solver 50 x

GPU-HMMER

Creation parallel dotplots for suite of protein sequences

GPU Accelerated Free Surface Flows Using Smoothed Particle Hydrodynamics 23 x

Multibody mechanical simulations on the GPU

GPU for Surveillance 20 x

3D Particle Boltzmann Solver 120 x

Accelerating Molecular Dynamic Simulations on GPUs Using OpenMM 100 x

Some Issues in Dense Linear Algebra for Multicore and Special Purpose Architectures 2 x

www.NVIDIA.com

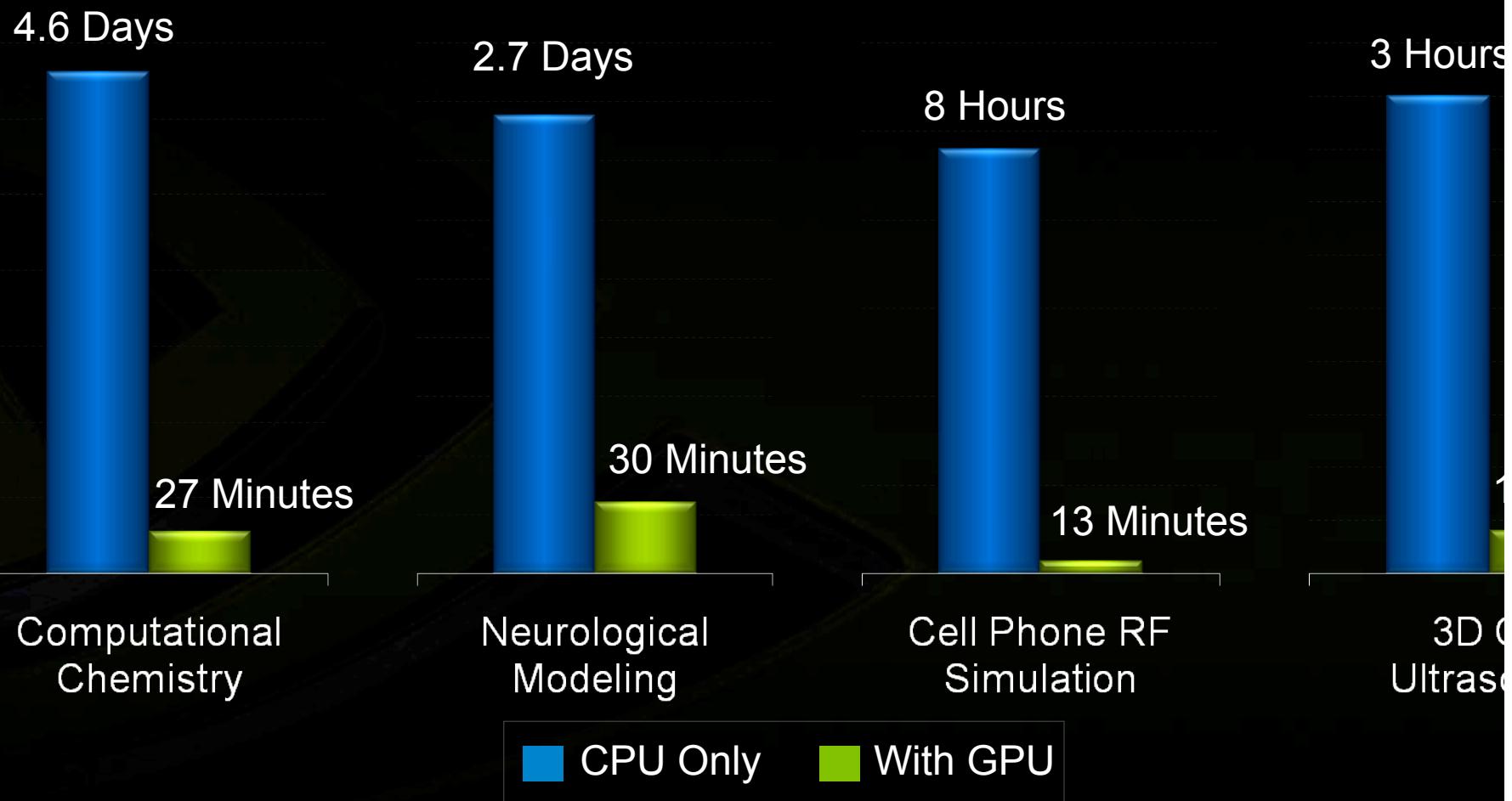
- 200+ papers and a
- 110 Million CUDA-E
- 60,000+ Active De

Application Domains

Details at

http://www.nvidia.com/object/vertical_solutions.html

Accelerating Time to Discovery



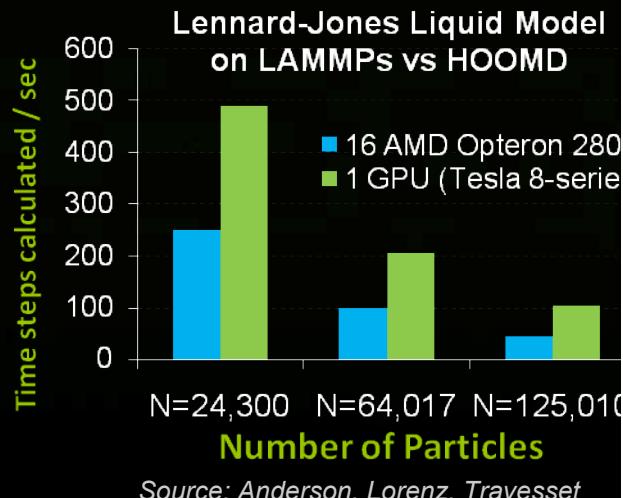
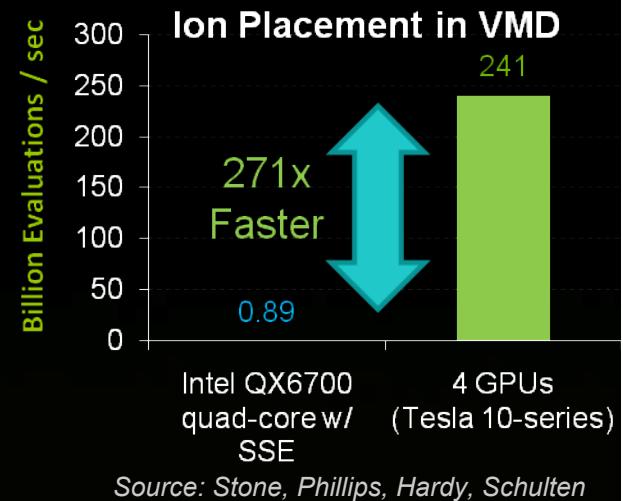
Molecular Dynamics

- **Available MD software**

- NAMD / VMD (alpha release)
- GROMACS (alpha release)
- HOOMD

- **OpenMM : Library for molecular modeling** <https://simtk.org/home/openmm>

- Jump start work on molecular dynamics



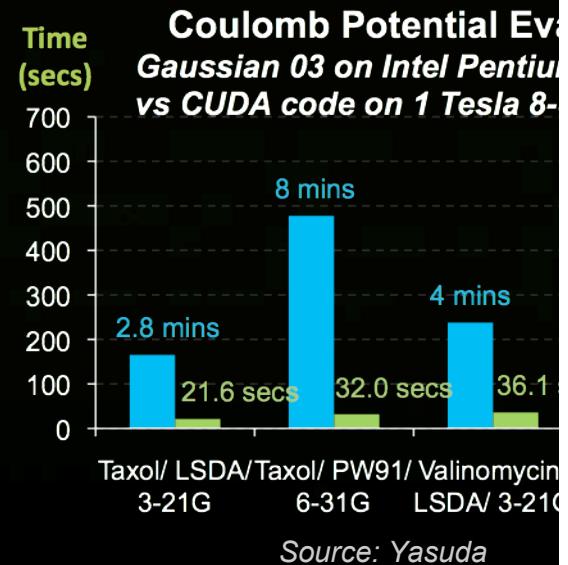
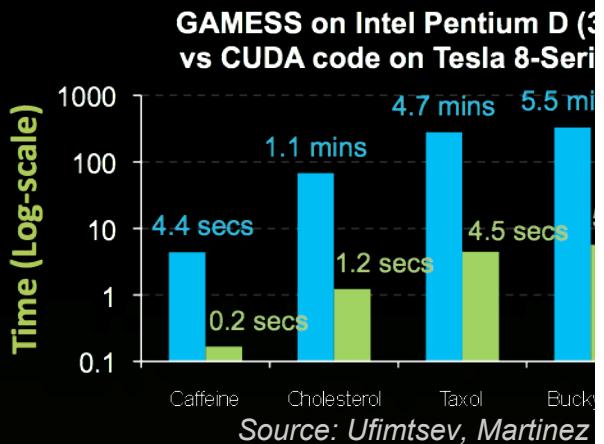
Quantum Chemistry

- Available MD software

- NAMD / VMD (alpha release)
- HOOMD
- ACE-MD
- MD-GPU

- Ongoing work

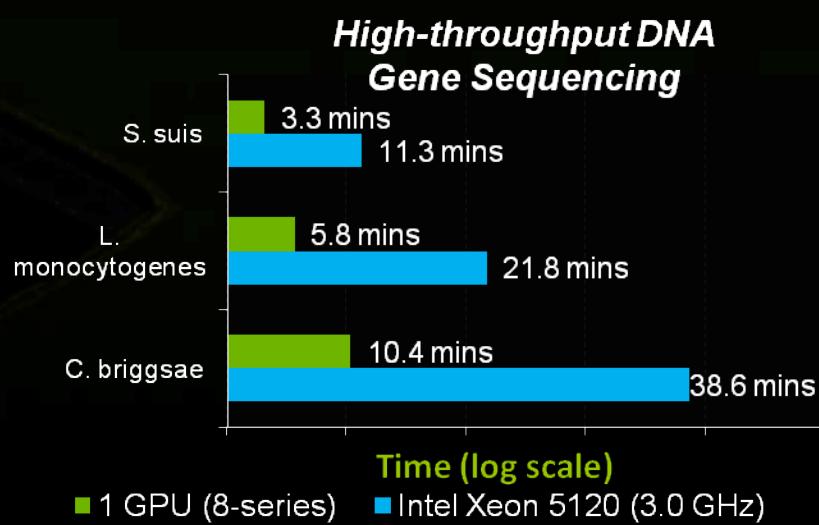
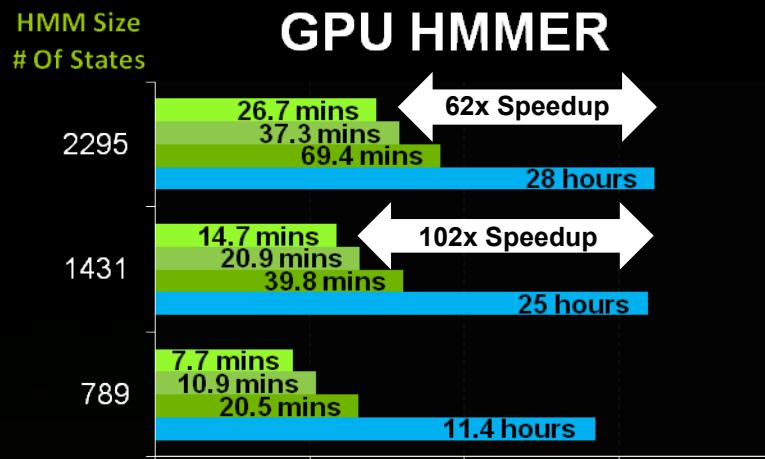
- LAMMPS
- CHARMM
- Q-Chem
- Gaussian



Bio-Informatics

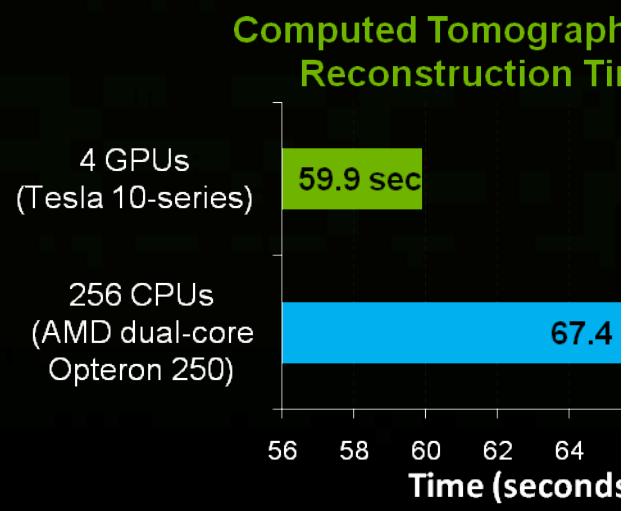
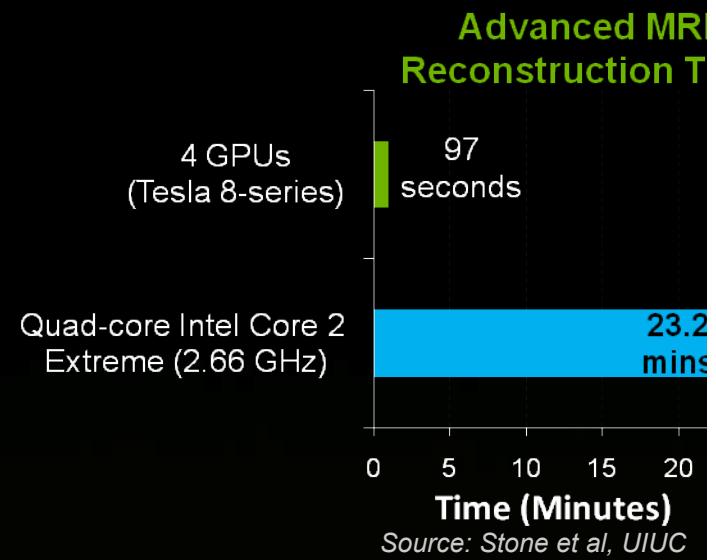
- Available applications

- GPU HMMER
- MUMmerGPU sequencing
- MATLAB acceleration
- Protein docking



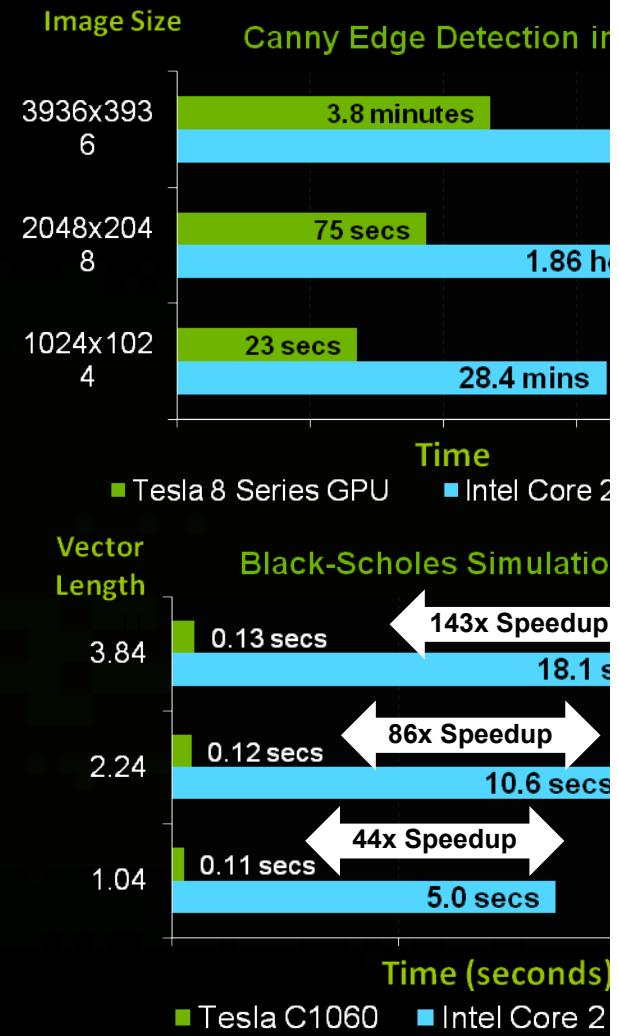
Medical Equipment

- GE Healthcare : CT
 - 40% increase in CT resolution
 - 2x increase in frame rate
- Techniscan: Ultra-sound
 - High resolution ultra-sound
 - 2x increase in acquisition
- Digisens : Tomography
 - Tomography reconstruction
- Several others on
 - X-Ray, Flow Cytometry, MRI, etc



MATLAB: Acceleration with GPUs

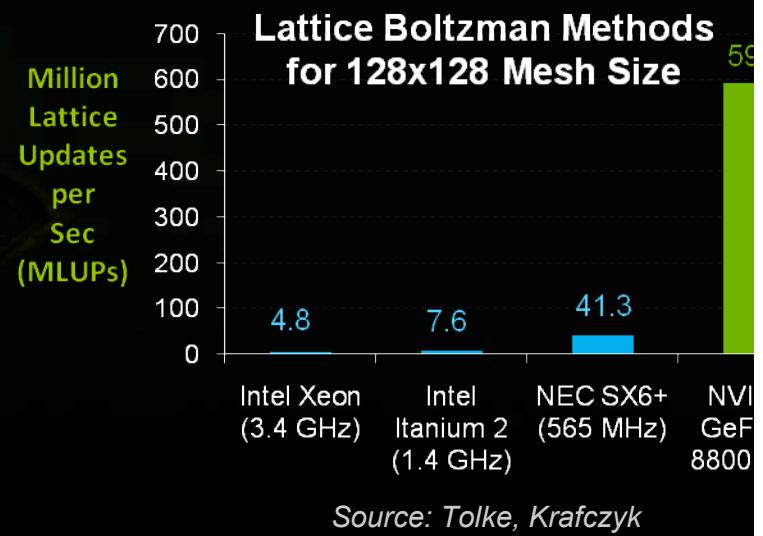
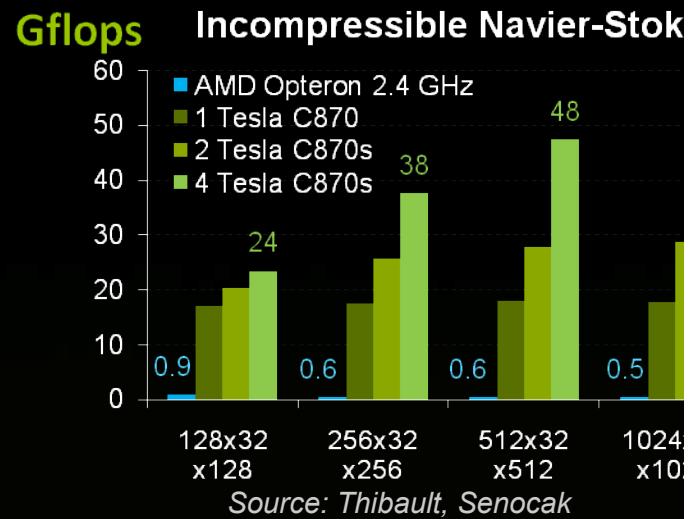
- **Jacket CUDA plugin from Accelereyes**
 - <http://www.accelereyes.com>
 - 15-day trial version available
- **Tesla GPU in a workstation**
 - For MATLAB and research



Computational Fluid Dynamics (CFD)

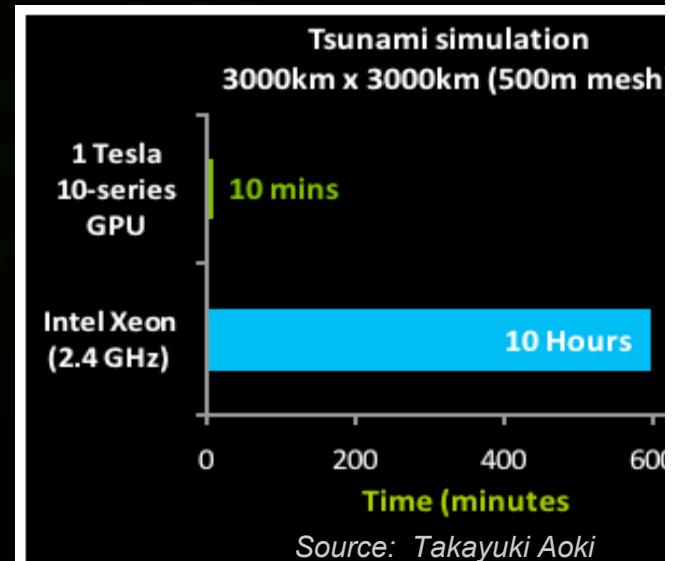
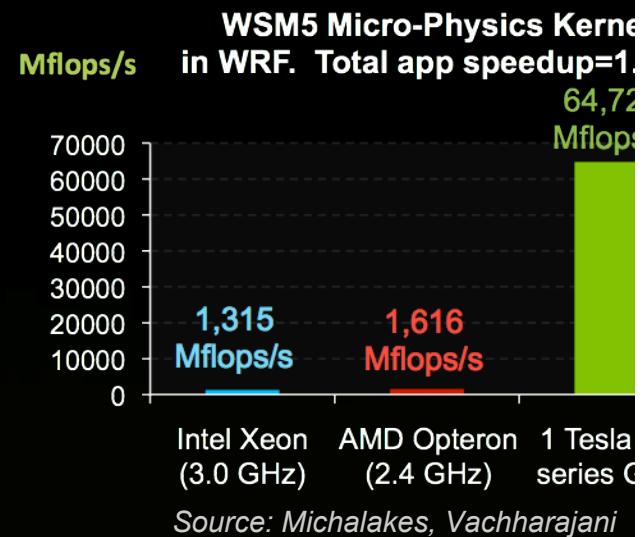
- **Ongoing work**

- **Navier-Stokes**
- **Lattice Boltzman**
- **3D Euler Solver**
- **Weather and ocean modeling**



Weather, Atmospheric, & Ocean Modeling

- CUDA-accelerated WRF available
 - 25-30% speedup in WRF so far
 - Other kernels in WRF being ported
- Ongoing work
 - Tsunami modeling
 - Ocean modeling
 - Several CFD codes

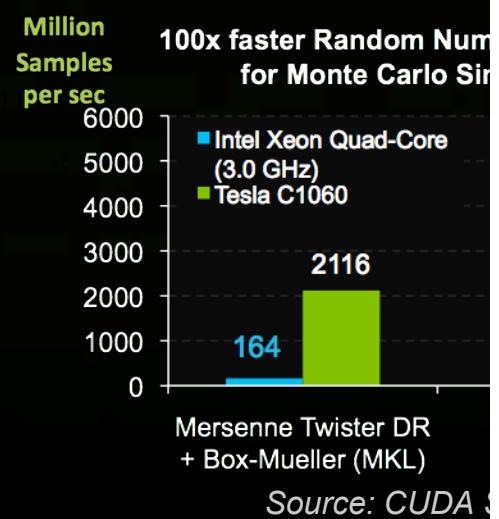


Electromagnetics / Electrodynamics

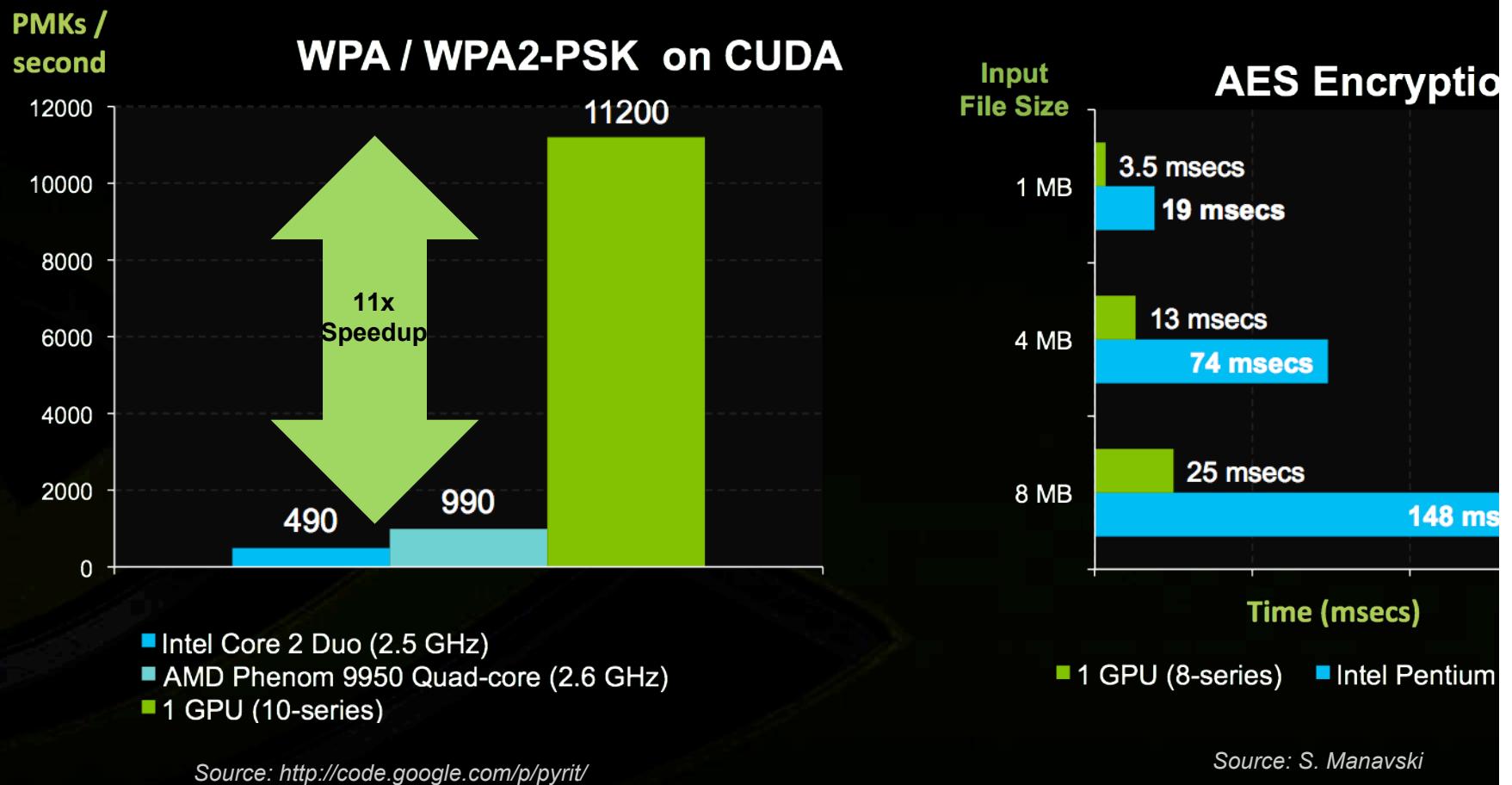
- **FDTD Solvers**
 - Acceleware
 - EM Photonics
 - CUDA Tutorial
 - **Ongoing work**
 - Maxwell equation solver
 - Ring Oscillator (FDTD)
 - Particle beam dynamics simulator
-
- The chart shows Speed (Mcells/s) on the y-axis (0 to 600) versus Time (T) on the x-axis. A single data point is plotted at approximately 9.9 Mcells/s, labeled with a red arrow. The chart title is 'Cell Phone Model Si' and the subtitle is 'Simulation size : 80'. The source is cited as 'Source: Acceleware'.
- Speed
Mcells/s
- Cell Phone Model Si
Simulation size : 80
- 9.9 Mcells/
s
- Intel Xeon (2.6 GHz)
- (T)
- FDTD Acceleration us
Source: Acceleware

Computational Finance

- **Financial Computing Software vendors**
 - SciComp : Derivatives pricing modeling
 - Hanweck: Options pricing & risk analysis
 - Aqumin: 3D visualization of market data
 - Exegy: High-volume Tickers & Risk Analysis
 - QuantCatalyst: Pricing & Hedging Engine
 - Oneye: Algorithmic Trading
 - Arbitragis Trading: Trinomial Options Pricing
- **Ongoing work**
 - LIBOR Monte Carlo market model
 - Callable Swaps and Continuous Time Finance



Encryption



Pattern Matching

String Matching in DNA Sequencing

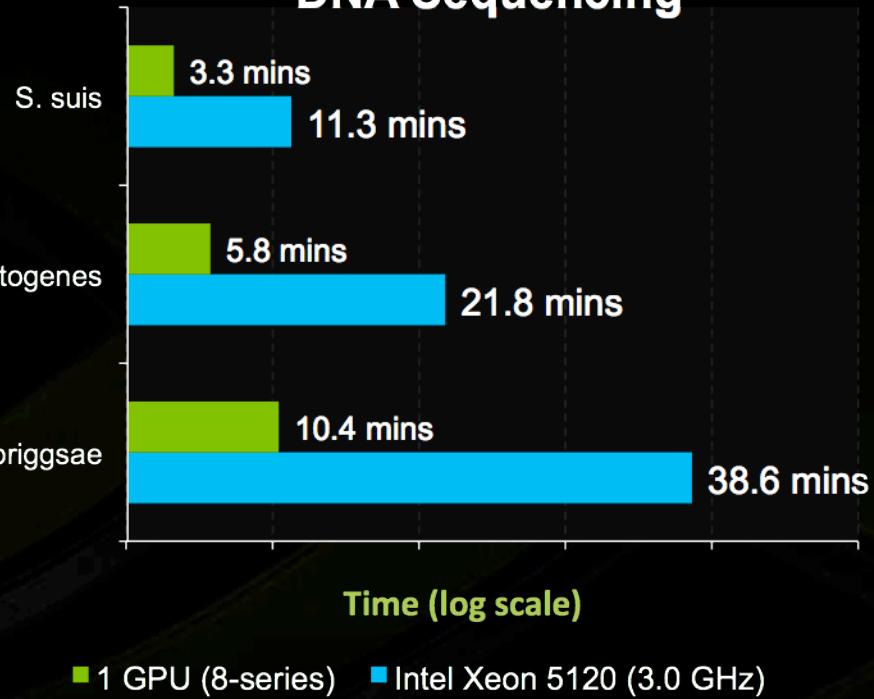
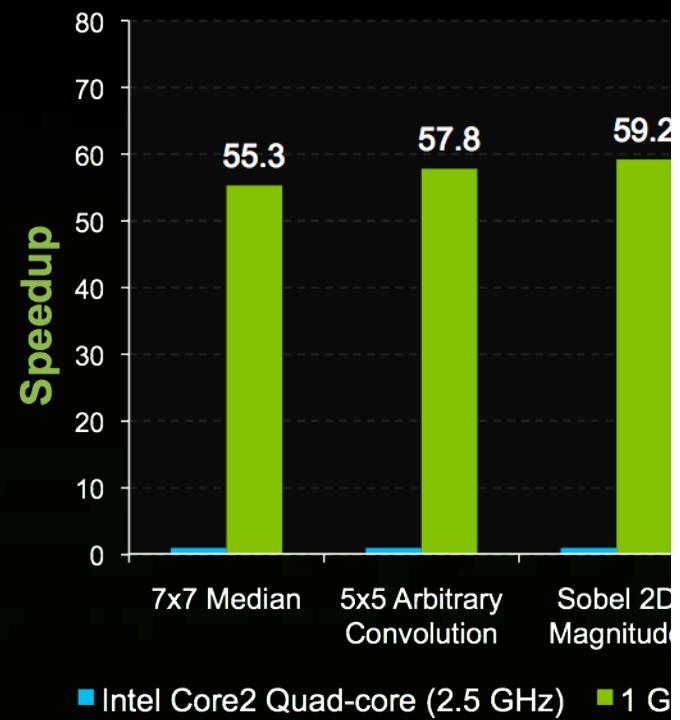


Image Process



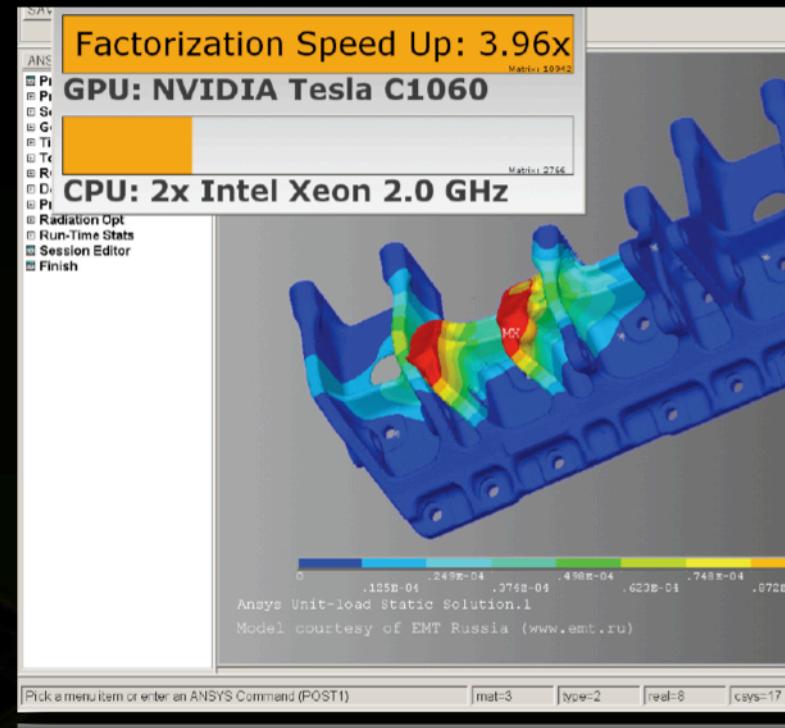
ANSYS Mechanical (Beta)

“CUDA allows us to leapfrog in a very compute intensive part of ANSYS.”

*Gene Poole
Principle SW Developer, ANSYS*



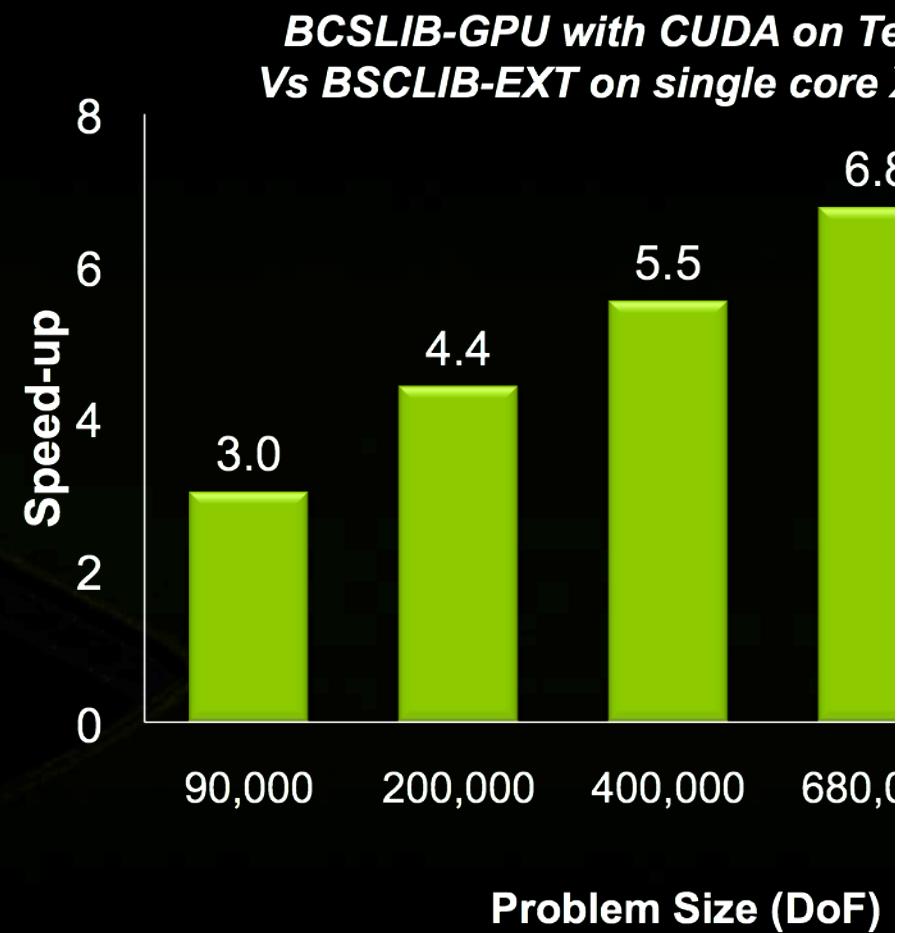
Hear Gene at <http://www.youtube.com/nvidiatesla>



Source: SuperComputing 2008, Ansys GPU acceleration of L

BCSLib-GPU : Sparse-Matrix Solver

- Used several solvers from Ansys, MSC
- Statics, Lanczos iterations, and Dynamics analyses
 - Not just LDL^T factorization
- Beta now, Release April 09



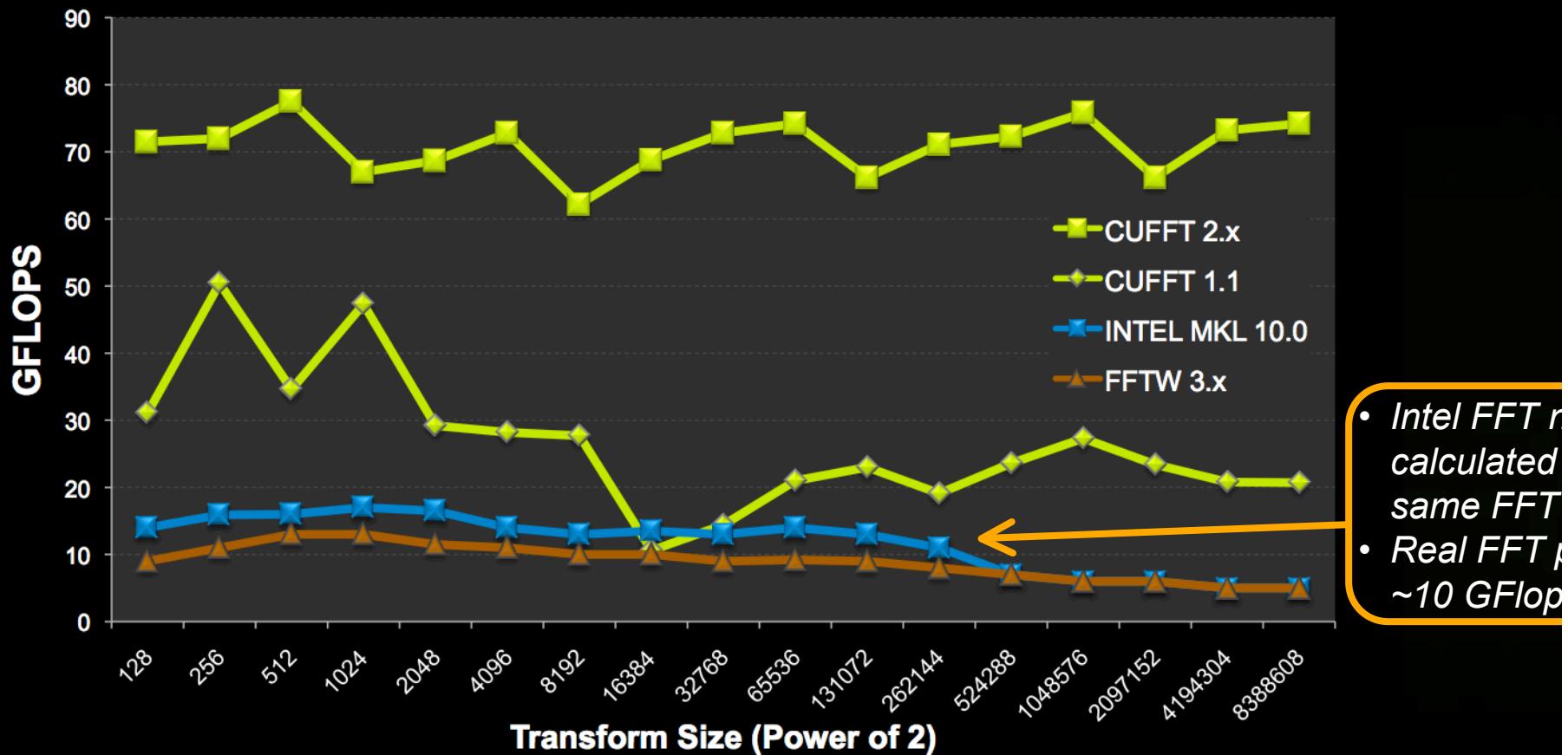
Libraries



FFT Performance: CPU vs GPU (8-Series)

1D Fast Fourier Transform On CUDA

NVIDIA Tesla C870 GPU (8-series GPU)
Quad-Core Intel Xeon CPU 5400 Series 3.0GHz,
In-place, complex, single precision

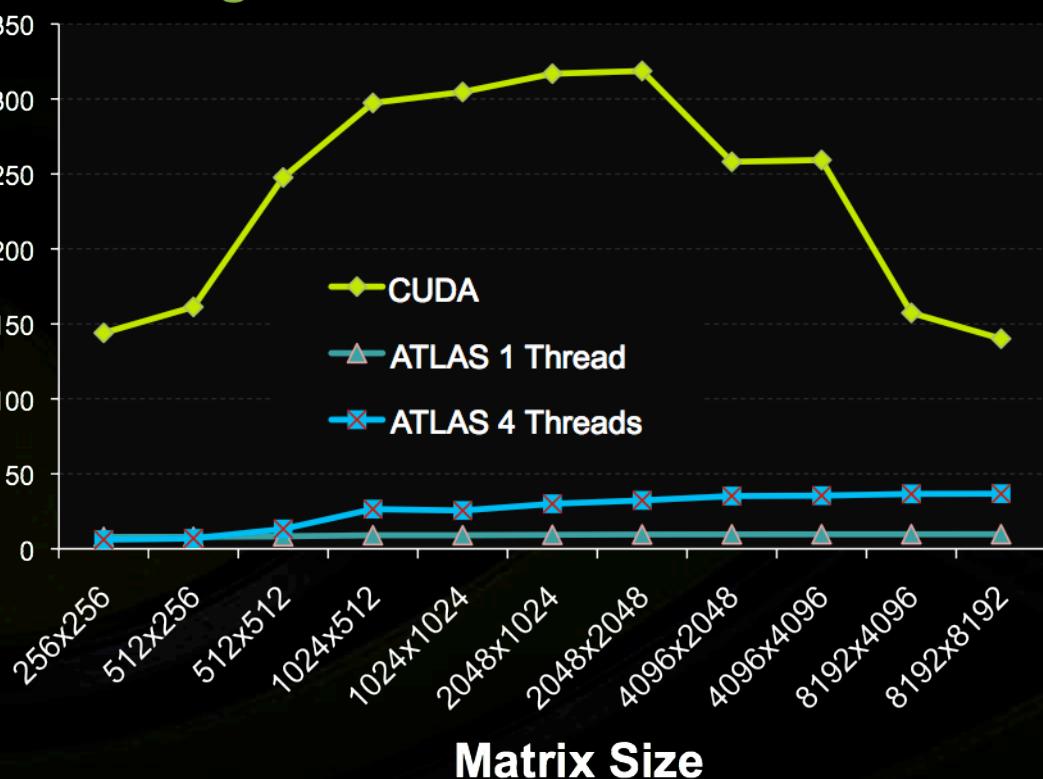


- Intel FFT results calculated same FFT
- Real FFT performance ~10 GFlop

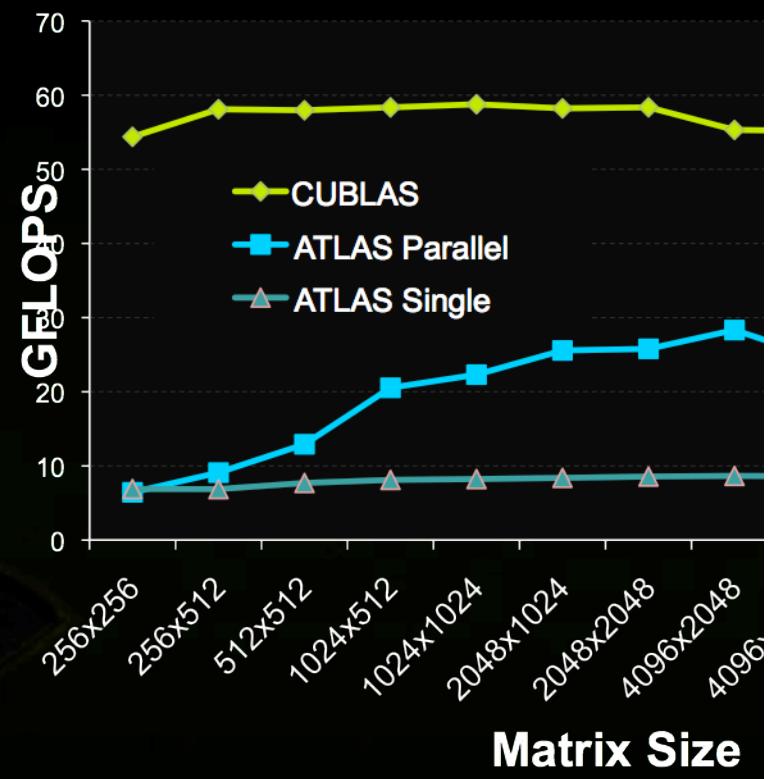
Source for Intel data : <http://www.intel.com/cd/software/products/asmo-na/eng/266852.htm>

BLAS: CPU vs GPU (10-series)

Single Precision BLAS: SGEMM



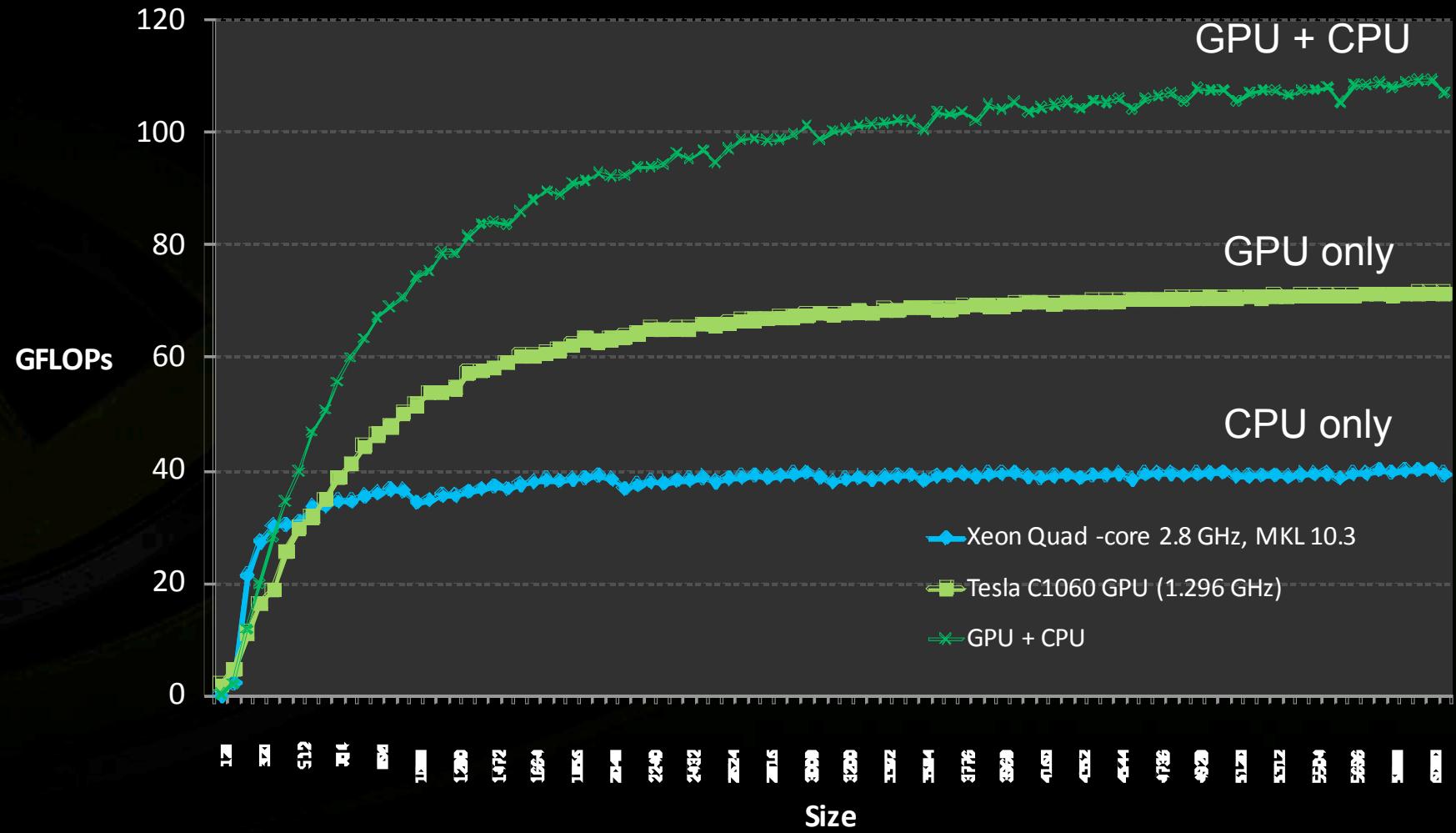
Double Precision BLAS:



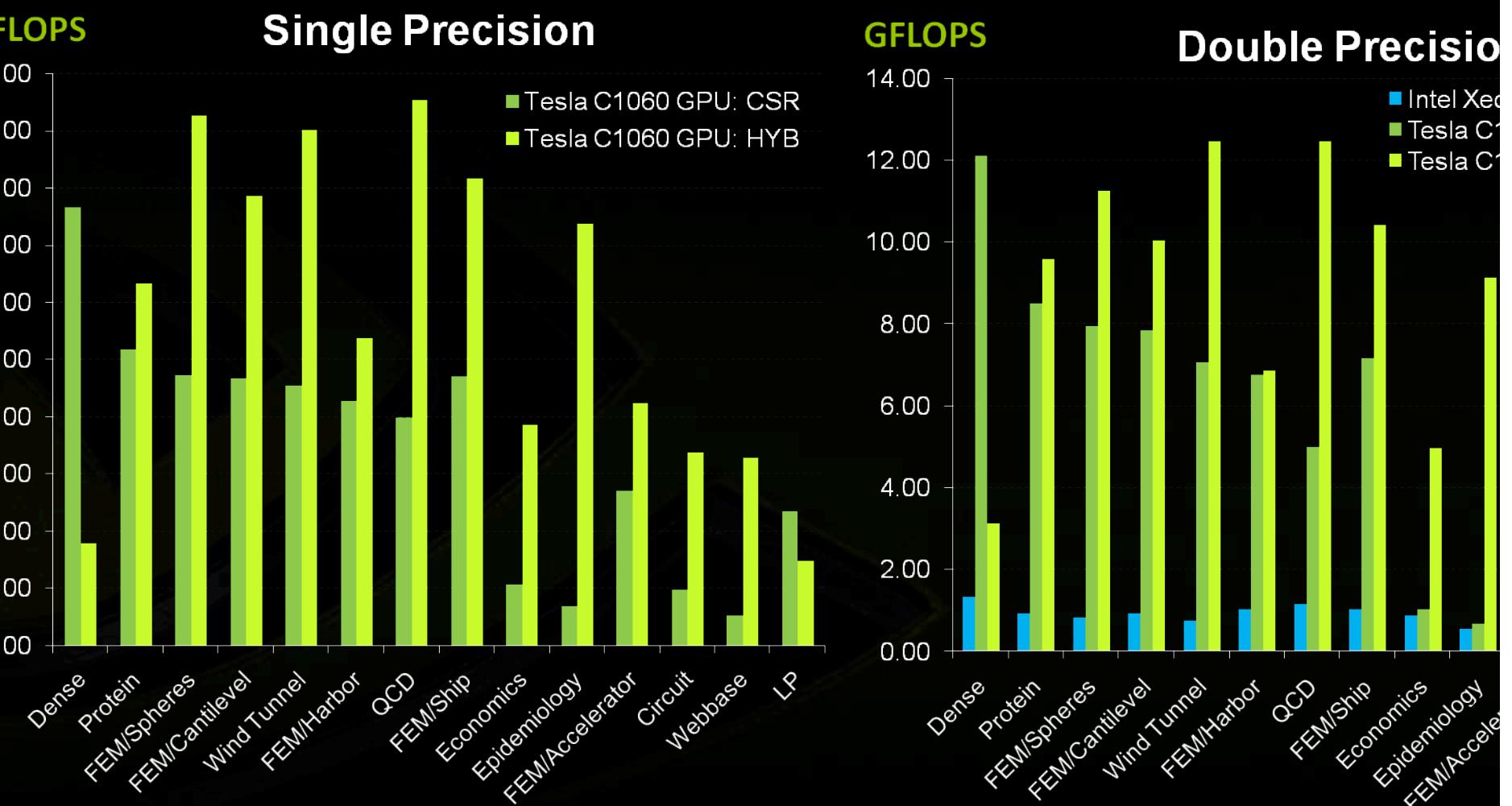
CUBLAS: CUDA 2.0, Tesla C1060 (**10-series** GPU)

ATLAS 3.81 on Dual 2.8GHz Opteron Dual-Core

GPU + CPU DGEMM Performance



Results: Sparse Matrix-Vector Multiplication (SpMV) on CUDA



CPU Results from "Optimization of Sparse Matrix-Vector Multiplication on Emerging Multicore Platforms", Williams et al, Supercomputing 2007

More Information

- **Tesla main page**
 - <http://www.nvidia.com/tesla>
 - Product Information
 - Software tools & libraries
- **Vertical Solutions**
 - http://www.nvidia.com/object/vertical_solutions.html
- **CUDA Zone**
 - <http://www.nvidia.com/cuda>
 - Applications, Papers,白皮书
 - Learn CUDA: self-paced tutorials
- **Hear from Developers**
 - <http://www.youtube.com/nvidiatechnology>

NVIDIA: Leadership in GPU computing

200+ Apps on CUDA Zone



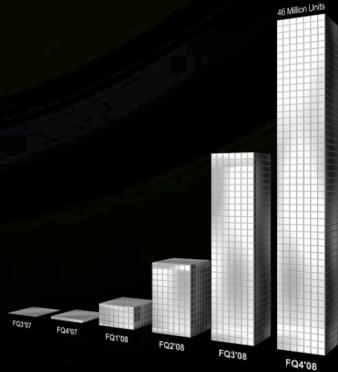
30+ CUDA GPU clusters



115+ Universities Teach
900+ research papers

Duke
Erlangen
ETH Zurich
Georgia Tech
Grove City College
Harvard
IISc Bangalore
IIIT Hyderabad
IIT
Illinois
INRIA
Iowa
ITESM
Johns Hopkins
Kent State
Kyoto
Lund
Maryland
McGill
MIT
North Carolina

110 M CUDA enabled GPUs
60,000+ active developers



150K CUDA compiler downloads

