

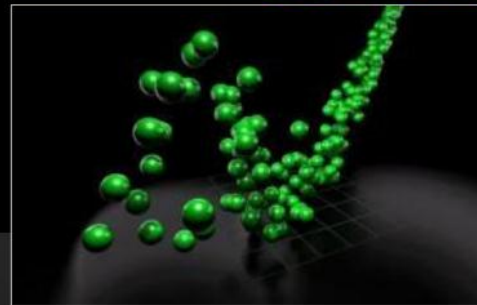
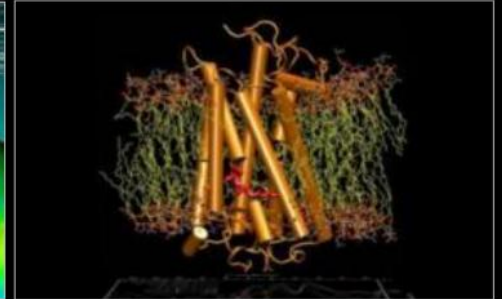
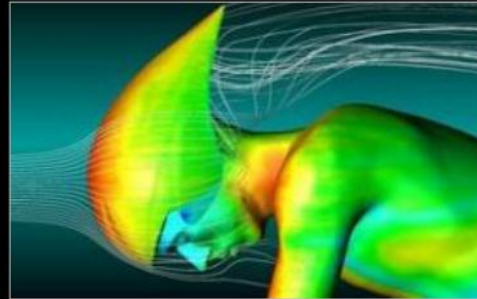
6th September 2011 – Belgrade

University of Belgrade, Mathematics faculty



NVIDIA HARDWARE FOR HIGH PERFORMANCE COMPUTING

Piero Altoè – HPC Sales Manager (EMEA)



TESLA

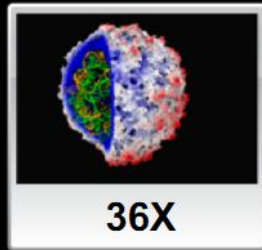
GPU Computing

Accelerating High Performance Computing

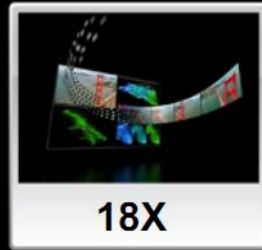
<http://www.nvidia.com/tesla>



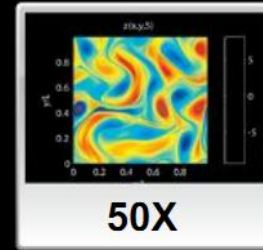
Medical Imaging
U of Utah



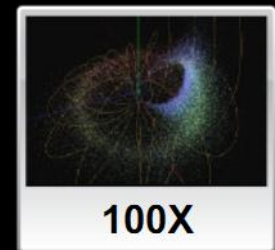
Molecular Dynamics
U of Illinois, Urbana



Video Transcoding
Elemental Tech

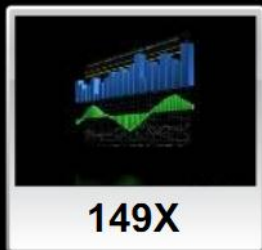


Matlab Computing
AccelerEyes

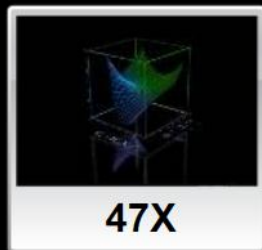


Astrophysics
RIKEN

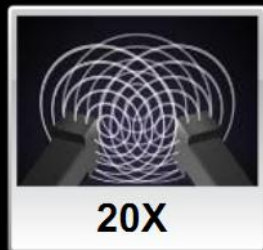
GPUs Accelerate Science



Financial Simulation
Oxford



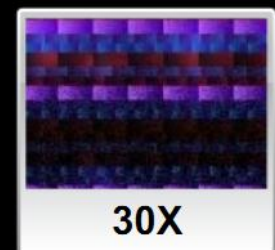
Linear Algebra
Universidad Jaime



3D Ultrasound
Techniscan

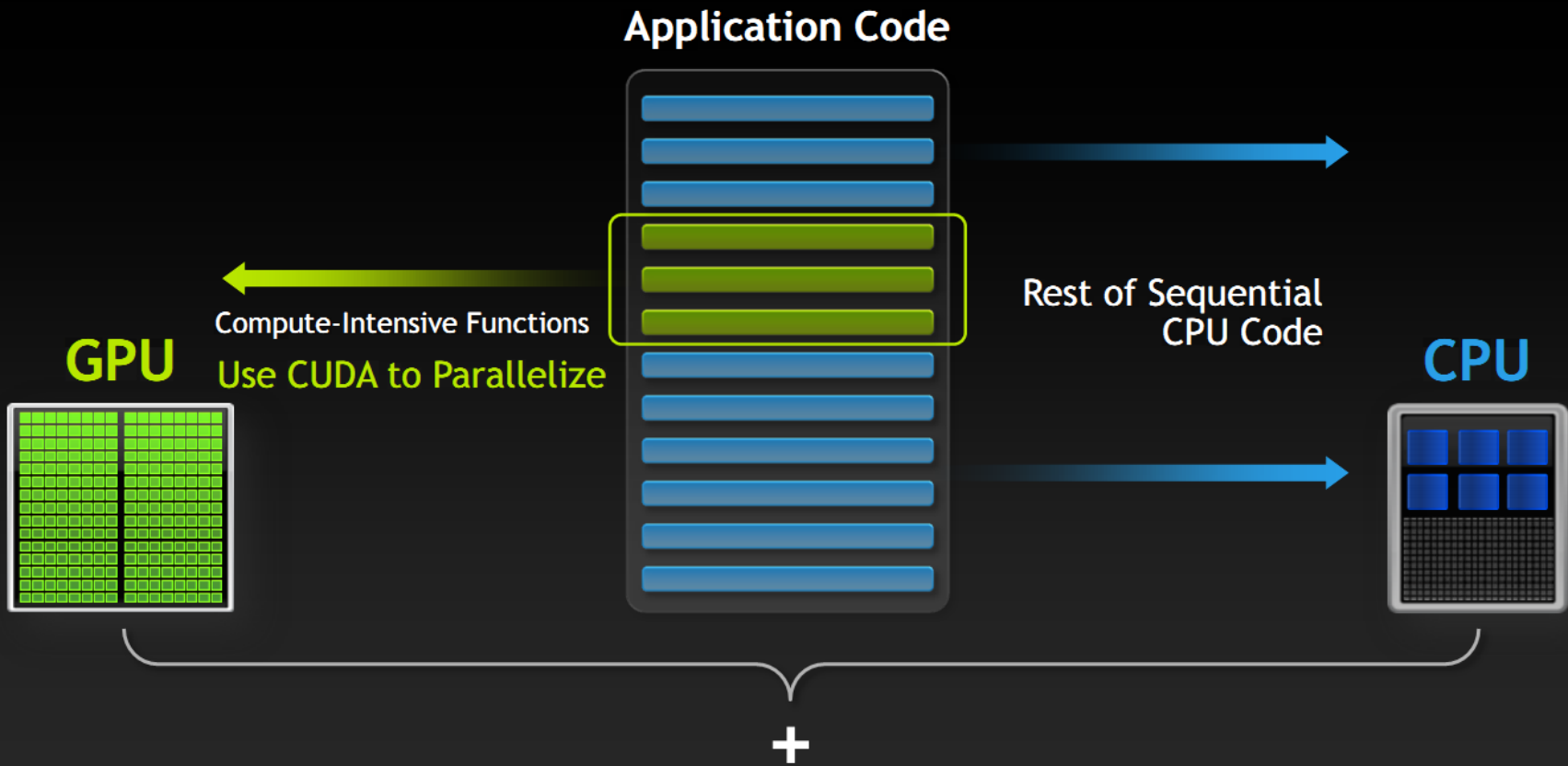


Quantum Chemistry
U of Illinois, Urbana



Gene Sequencing
U of Maryland

Minimum Change, Big Speed-up



Tesla GPUs Power 3 of Top 5 Supercomputers

#1 : Tianhe-1A

7168 Tesla GPU's 2.5 PFLOPS



#3 : Nebulae

4650 Tesla GPU's 1.2 PFLOPS



#4 : Tsubame 2.0

4224 Tesla GPU's 1.194 PFLOPS



“

We not only created the world's fastest computer, but also implemented a heterogeneous computing architecture incorporating CPU and GPU, this is a new innovation. ”

Premier Wen Jiabao

Public comments acknowledging Tianhe-1A

World's Fastest MD Simulation

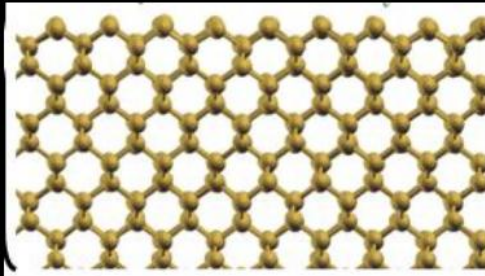
Sustained Performance of 1.87 Petaflops/s

Institute of Process Engineering (IPE)

Chinese Academy of Sciences (CAS)

**Used all 7168 Tesla GPUs on
Tianhe-1A GPU Supercomputer**

MD Simulation for Crystalline Silicon



World's Greenest Petaflop Supercomputer

Tsubame 2.0

Tokyo Institute of Technology

- 1.19 Petaflops
- 4,224 Tesla M2050 GPUs



Commercial Apps Accelerated by GPUs

▶ Molecular Dynamics	AMBER	CHARMM	DL_POLY	GROMACS	LAMMPS	NAMD
▶ Fluid Dynamics	Altair Acusolve Turbostream	Autodesk Moldflow	OpenFOAM	Prometech Particlework		
▶ Earth Sciences	ASUCA	HOMME	NASA GEOS-5	NOAA NIM	WRF	
▶ Engineering Simulation	Agilent EMPro Impetus AFEA	ANSYS Mechanical Remcom XFDTD	ANSYS Nexxim	SIMULIA Abaqus	CST Microwave Studio	
▶ Others	GADGET2 PARATEC	MATLAB Schlumberger Petrel	Mathematica	NBODY	Paradigm VoxelGeo	

Widespread Adoption of GPUs

Oil and gas

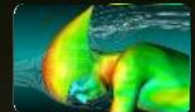
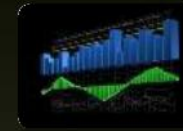
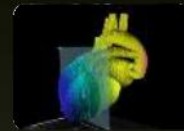
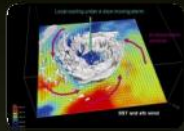
Edu/Research

Government

Life Sciences

Finance

Manufacturing



Seismic Processing
Reservoir Sim

Astrophysics
Molecular
Dynamics
Weather / Climate

Signal Processing
Satellite Imaging
Video Analytics

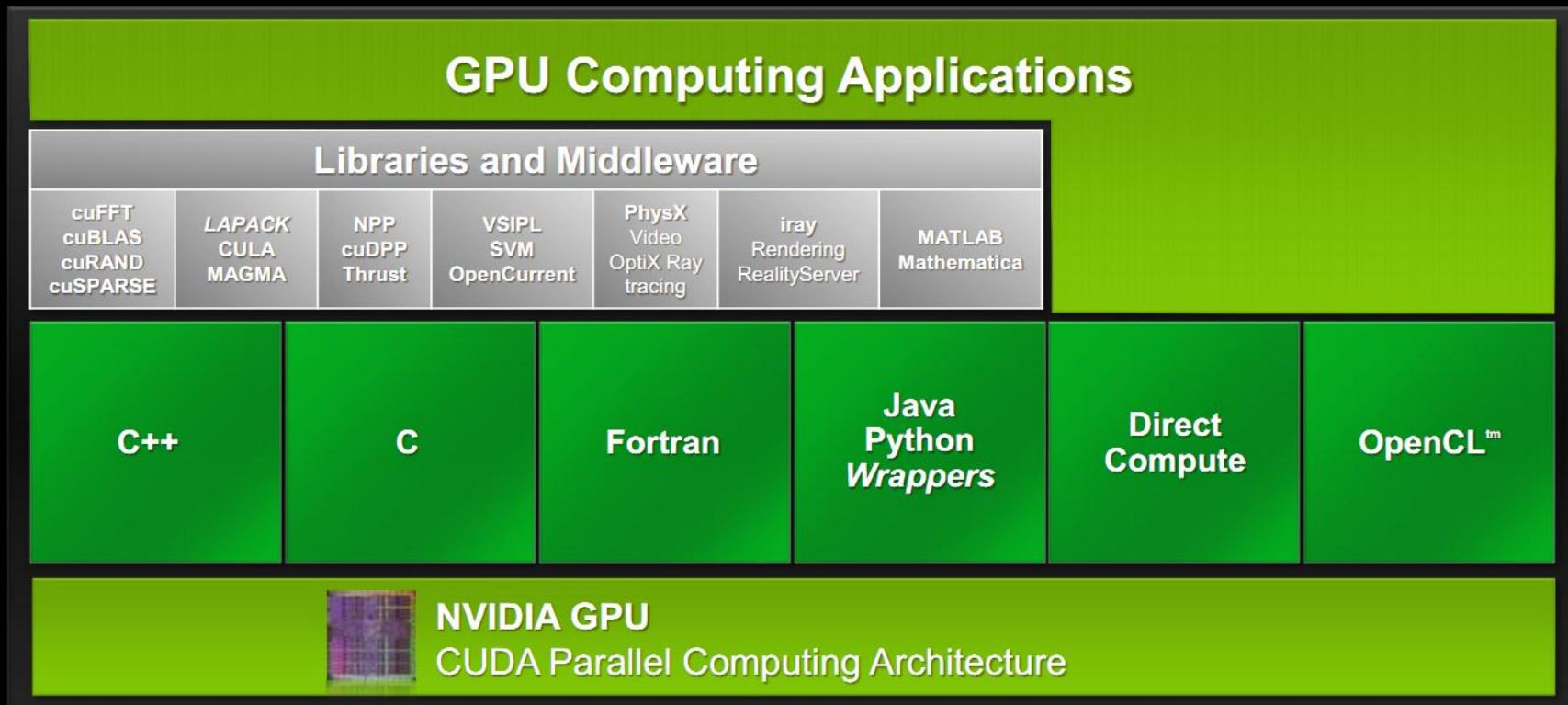
Bio-chemistry
Bio-informatics
Material Science
Genomics

Risk Analytics
Monte Carlo
Options Pricing
Insurance

Structural
Mechanics
Computational
Fluid Dynamics
Electromagnetics



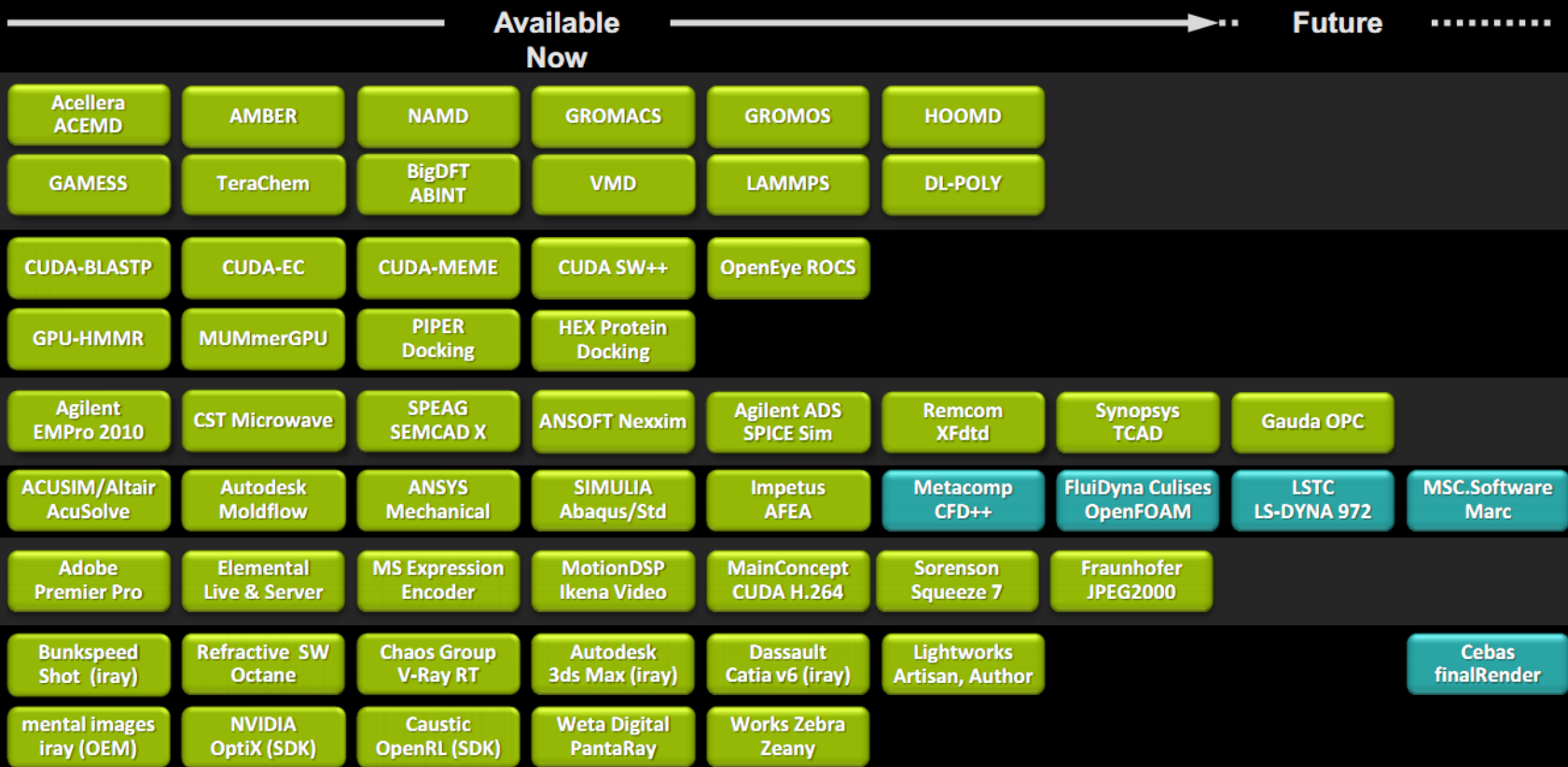
CUDA: Easy to Use Parallel Programming Model



Increasing Number of Professional CUDA Applications

	Available Now							Future
Tools & Libraries	CUDA C/C++	Parallel Nsight Vis Studio IDE	NVIDIA Video Libraries	ParaTools VampirTrace	PGI Accelerators	EMPhotonics CULAPACK	Allinea DDT Debugger	TauCUDA Perf Tools	PGI CUDA-X86	
	NVIDIA NPP Perf Primitives	PGI Fortran	Thrust C++ Template Lib	Bright Cluster Manager	CAPS HMPP	MAGMA	GPU Packages For R Stats Pkg	Platform LSF Cluster Mgr	GPU.net	
	pyCUDA	R-Stream Reservoir Labs	PBSWorks	MOAB Adaptive Comp	Torque Adaptive Comp	TotalView Debugger	IMSL			
Oil & Gas	Headwave Suite	OpenGeo Solns OpenSEIS	GeoStar Seismic	Acceleware RTM Solver	StoneRidge RTM	Seismic City RTM	Tsunami RTM		Schlumberger Petrel	
	ffa SVI Pro	Paradigm SKUA	VSG Open Inventor	Paradigm GeoDepth RTM	VSG Avizo	SVI Pro	SEA 3D Pro 2010	Schlumberger Omega	Paradigm VoxelGeo	
Numerical Analytics	LabVIEW Libraries	AccelerEyes Jacket: MATLAB	MATLAB	Mathematica						
Finance	NAG RNG	Numerix CounterpartyRisk	SciComp SciFinance	Aquimin AlphaVision	Hanweck Volera Options Anlysi	Murex MACS				
Other	Siemens 4D Ultrasound	Digisens CT	Schrodinger Core Hopping	Useful Prog Medical Imag	ASUCA Weather Model					
	Manifold GIS	MVTech Mach Vision	Dalsa Mach Vision	WRF Weather						

Increasing Number of Professional CUDA Applications



CUDA 4.0: Highlights

Easier Parallel Application Porting

- Share GPUs across multiple threads
- Single thread access to all GPUs
- No-copy pinning of system memory
- New CUDA C/C++ features
- Thrust templated primitives library
- NPP image/video processing library
- Layered Textures

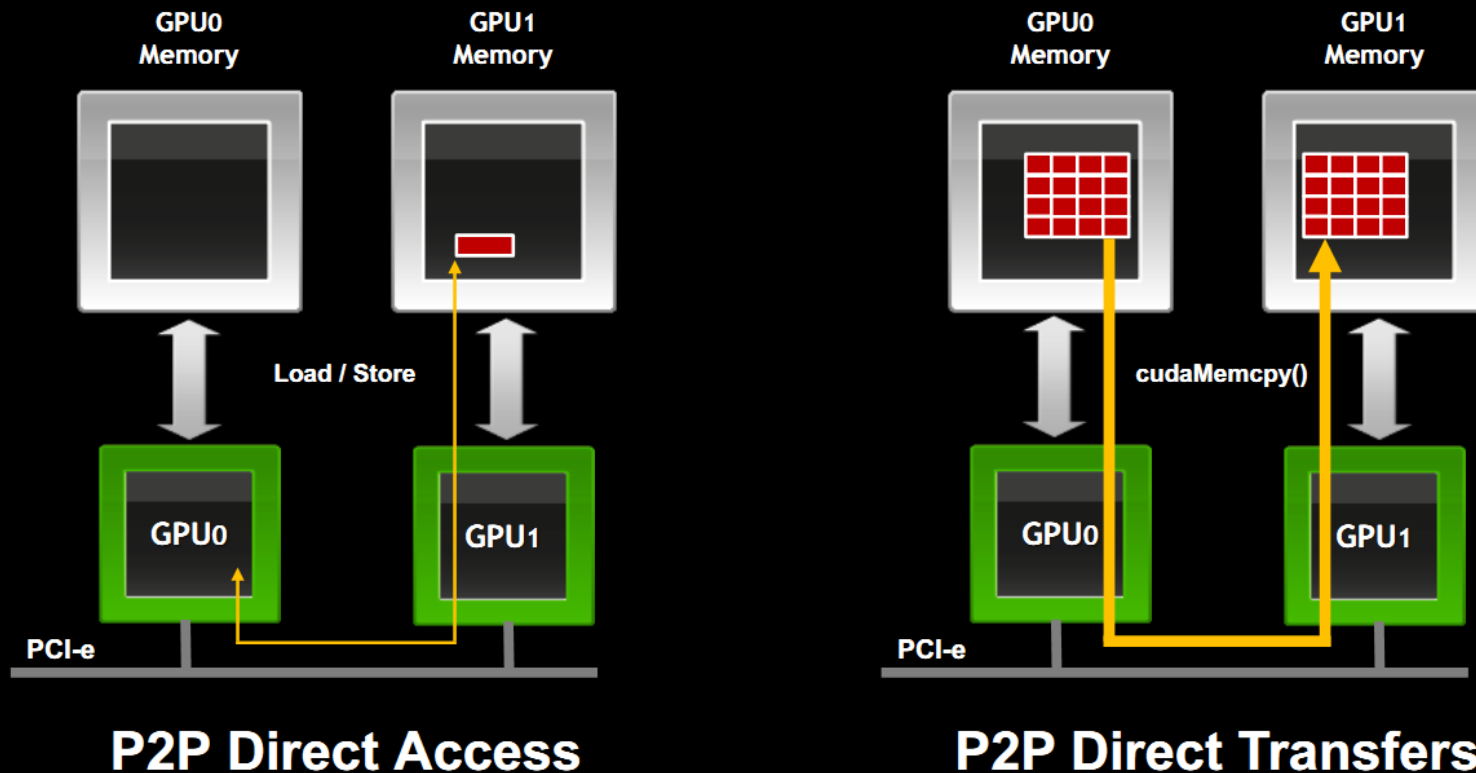
Faster Multi-GPU Programming

- Unified Virtual Addressing
- NVIDIA GPUDirect™ v2.0
 - Peer-to-Peer Access
 - Peer-to-Peer Transfers
 - GPU-accelerated MPI

New & Improved Developer Tools

- Auto Performance Analysis
- C++ Debugging
- GPU Binary Disassembler
- cuda-gdb for MacOS

GPUDirect v2.0: Peer-to-Peer Communication



Tesla Data Center & Workstation GPU Solutions



Tesla M-series GPUs
M2090 | M2070 | M2050

Servers & Blades



Tesla C-series GPUs
C2070 | C2050

Workstations

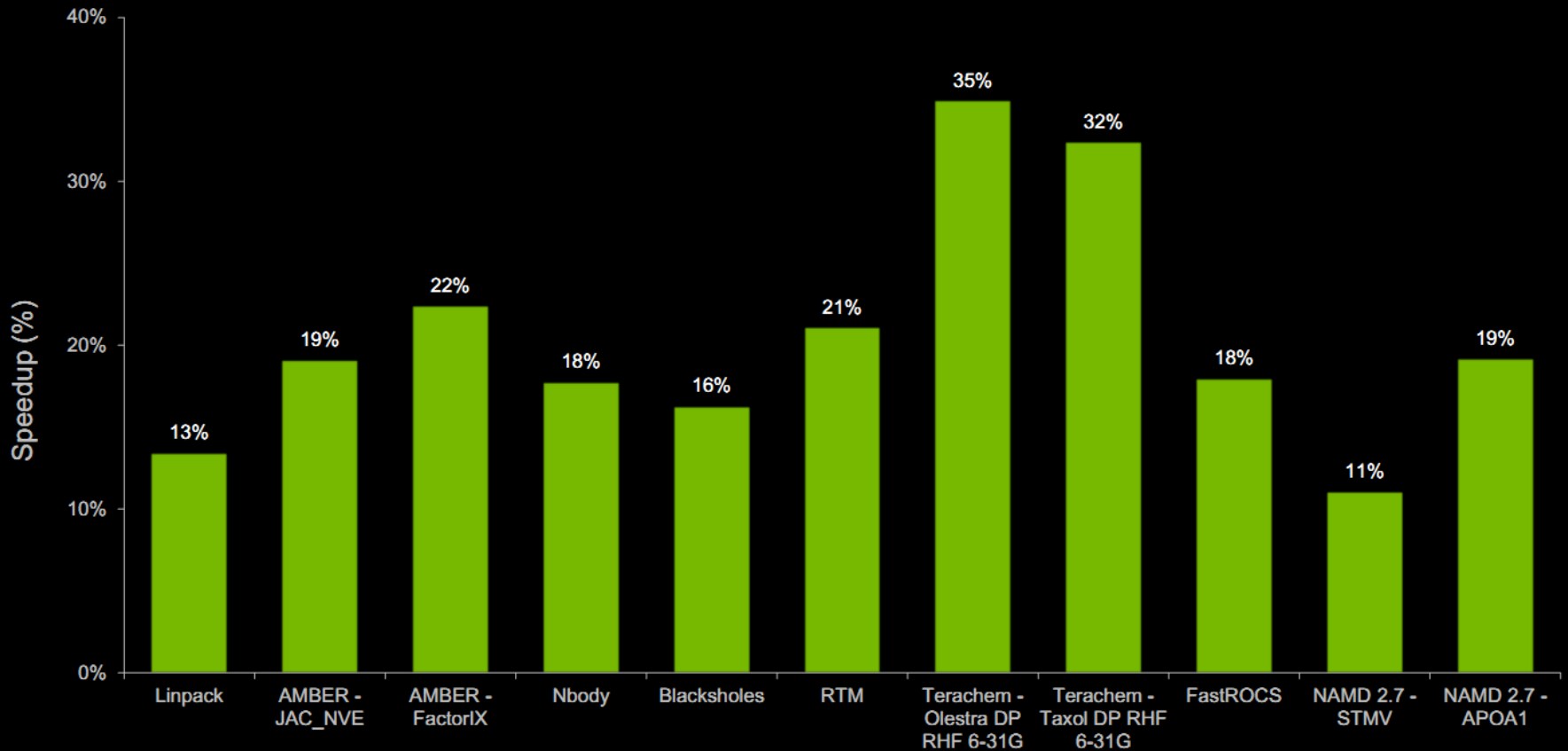
		M2090	M2070	M2050
Cores		512	448	448
Memory		6 GB	6 GB	3 GB
Memory bandwidth (ECC off)		177.6 GB/s	150 GB/s	148.8 GB/s
Peak Perf Gflops	Single Precision	1331	1030	1030
	Double Precision	665	515	515

C2070	C2050
448	448
6 GB	3 GB
148.8 GB/s	148.8 GB/s
1030	1030
515	515

GPU-Based OEM Systems

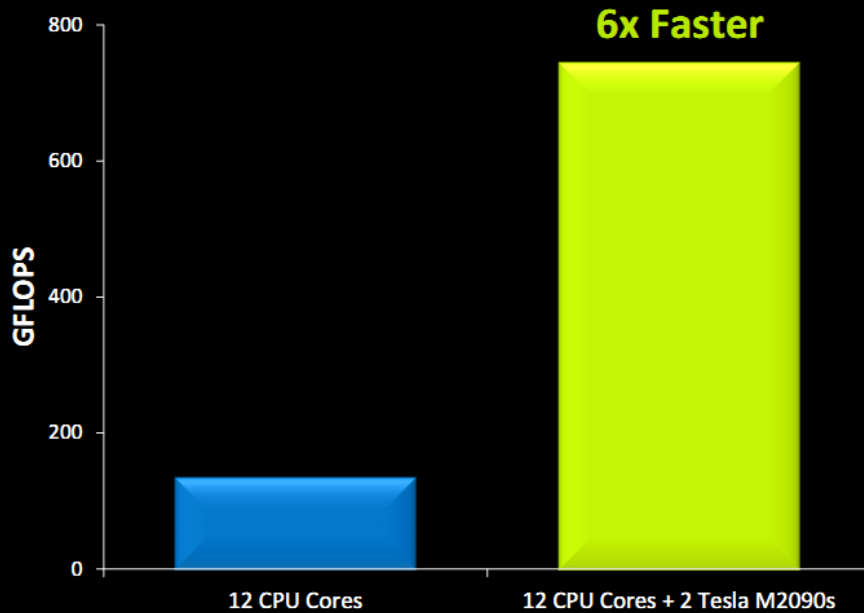
OEM	System	# of GPUs
Appro	Tetra	4GPUs in 1U
Asus	ESC4000	4GPUs in 2U
BULL	Bullx	18 GPUs in 7U
Dell	C61x	16 GPUs in 4U
HP	SL390G7	3 GPUs in 1U 8 GPUs in 2U
IBM	iDataplex	2 GPUs in 2U
NextIO	vCORE Express	4 GPUs in 1U
SGI	Prism XL	1 GPU in stick
Supermicro	6016GT-1U, TwinBlade-7126	2 GPUs in 1U 2 GPUs/blade
Tyan	FT72-B7015, GN70	8GPUs in 4U 3GPU/2U

Tesla M2090 vs M2070

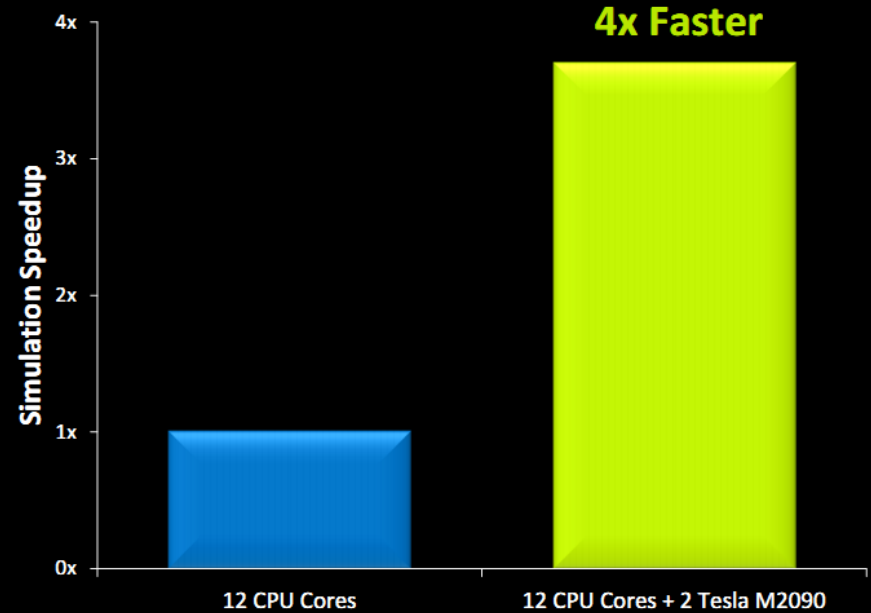


Maximize Compute Perf with Tesla M2090

LINPACK



NAMD 2.7



Benchmark

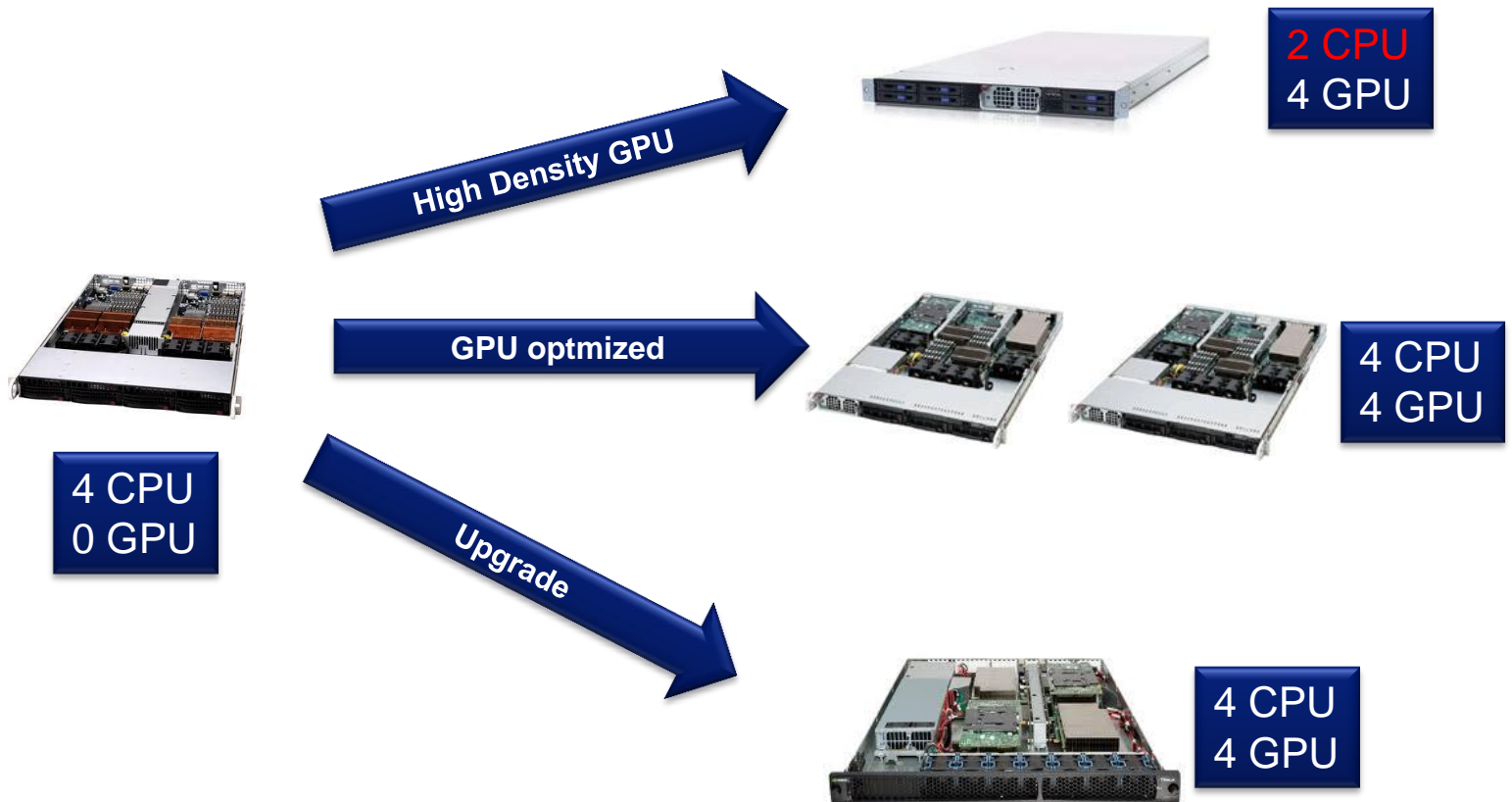
CPU: Dual socket Intel Xeon X5670, 2.93 GHz (12 cores)

Memory: 48 GB DDR3

CUDA 3.2

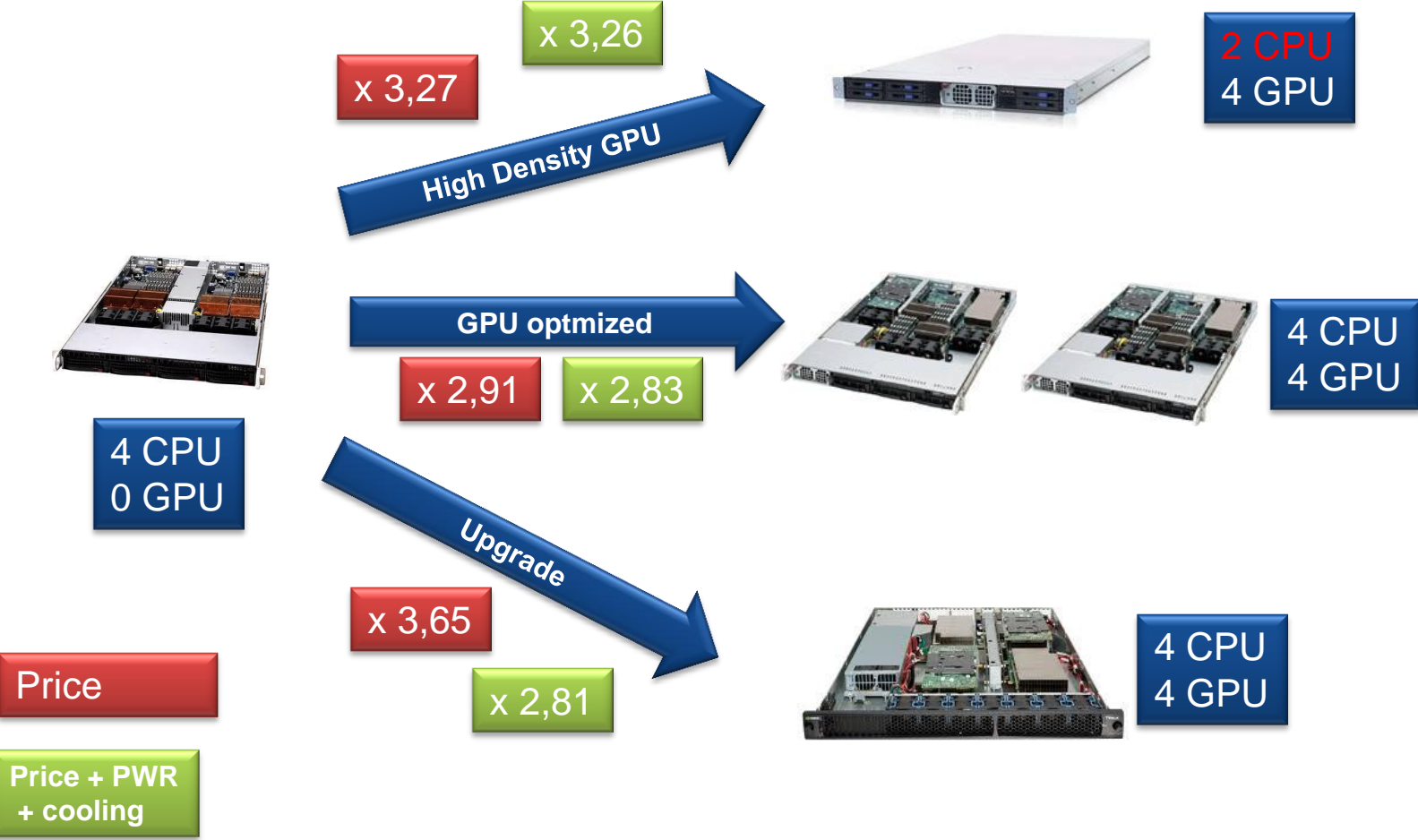


GPU: when & why?

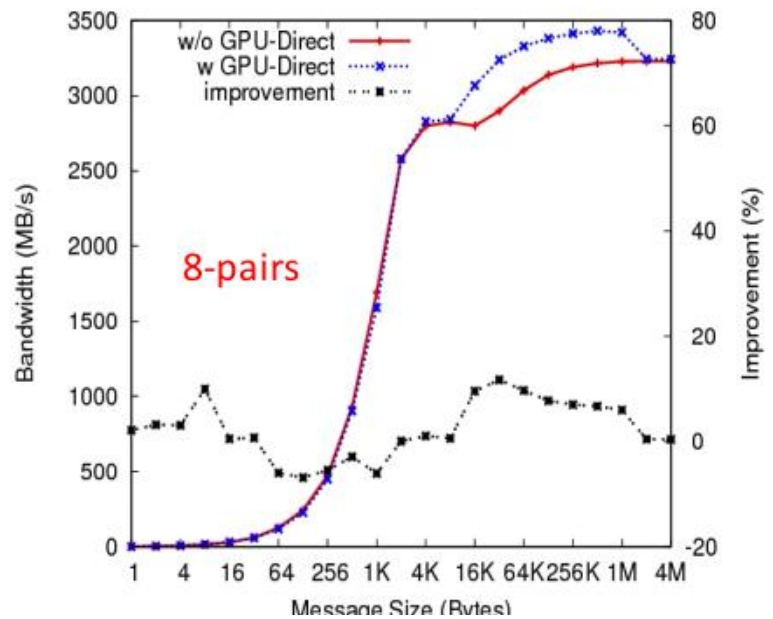
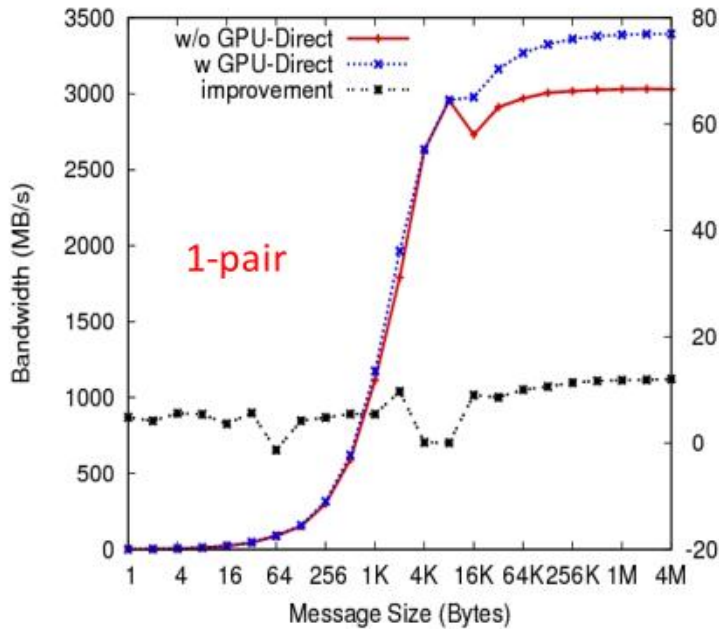




GPU: cost/performance

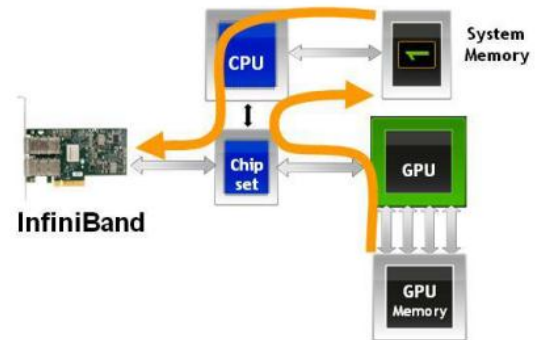


Communication GPU to GPU



GPU direct: shared host buffer between GPU and HCA IB

Better performance up to 12%





How to compare CPU and GPU calculations:

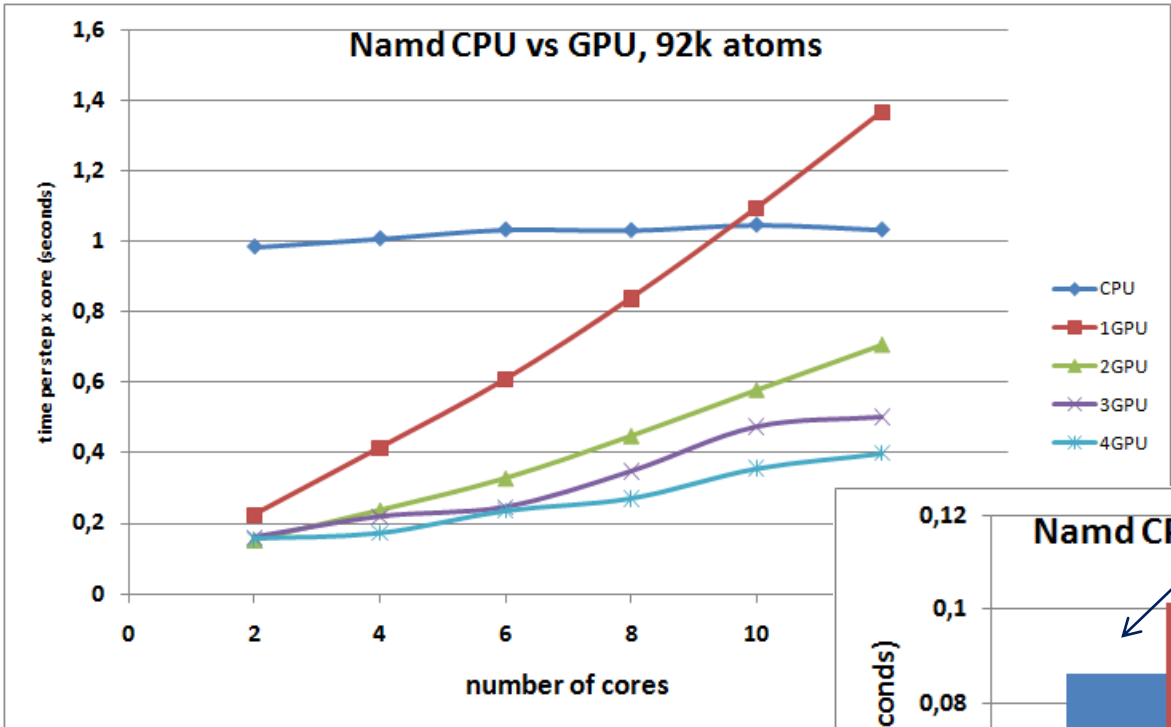
1. Never compare single core vs. single GPU
2. Take a look in the hardware configuration of the node
3. Comparison has to be done between the node with GPUs and without GPUs



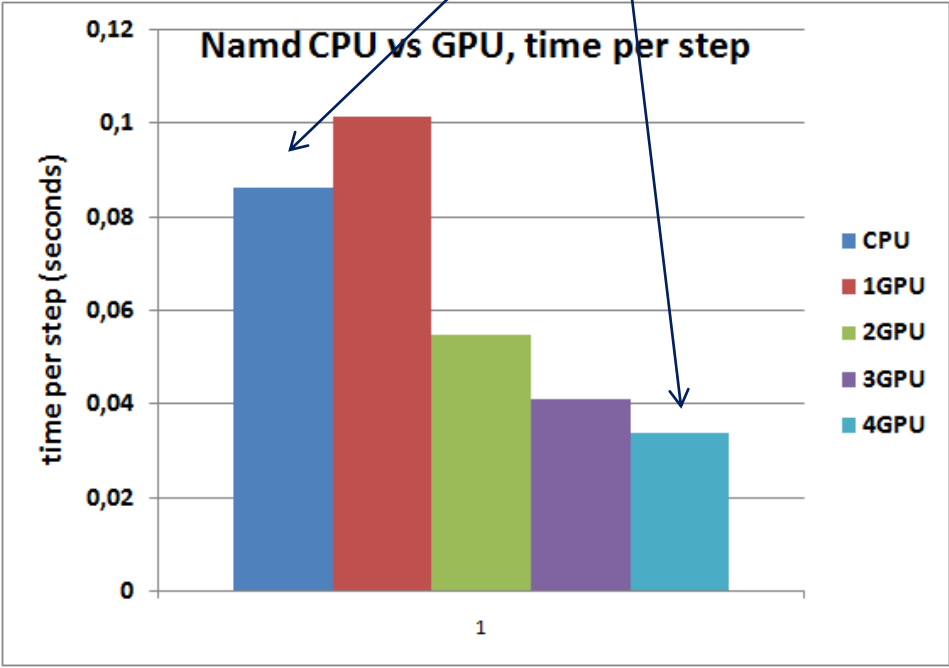
Why NAMD?

Because you can define independently
the number of CPU cores & number of GPUs

```
>: Charmrun ++local +p 8 namd2 +idlepoll +devices 0,1 input.namd
```

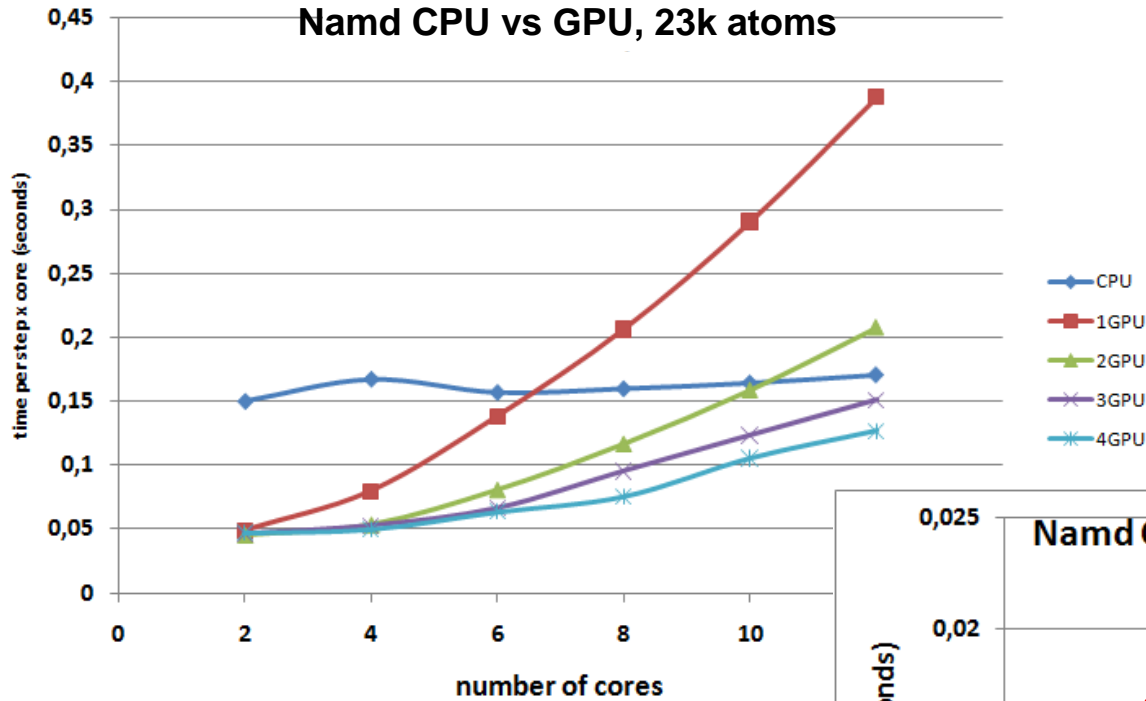


Speedup 2.6x



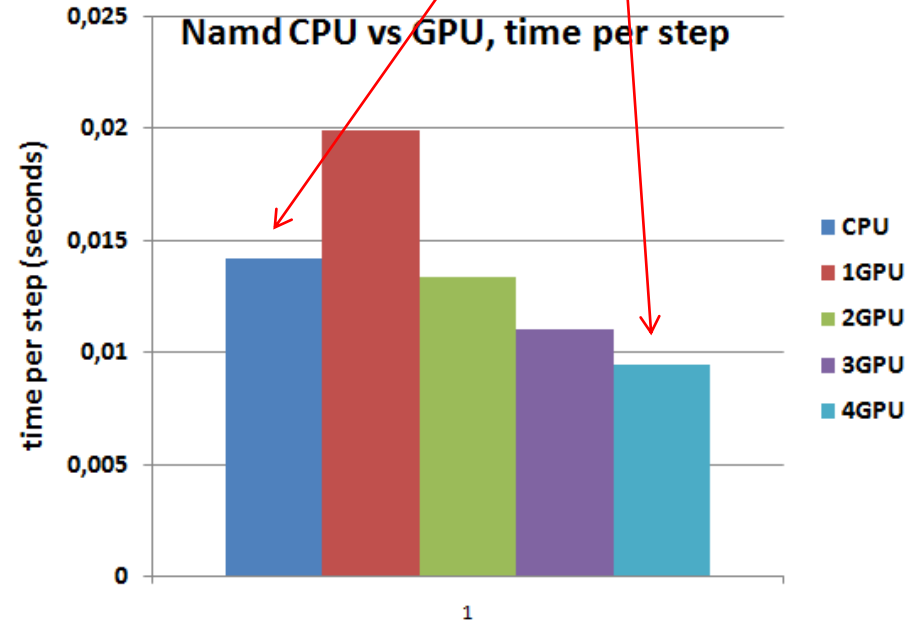


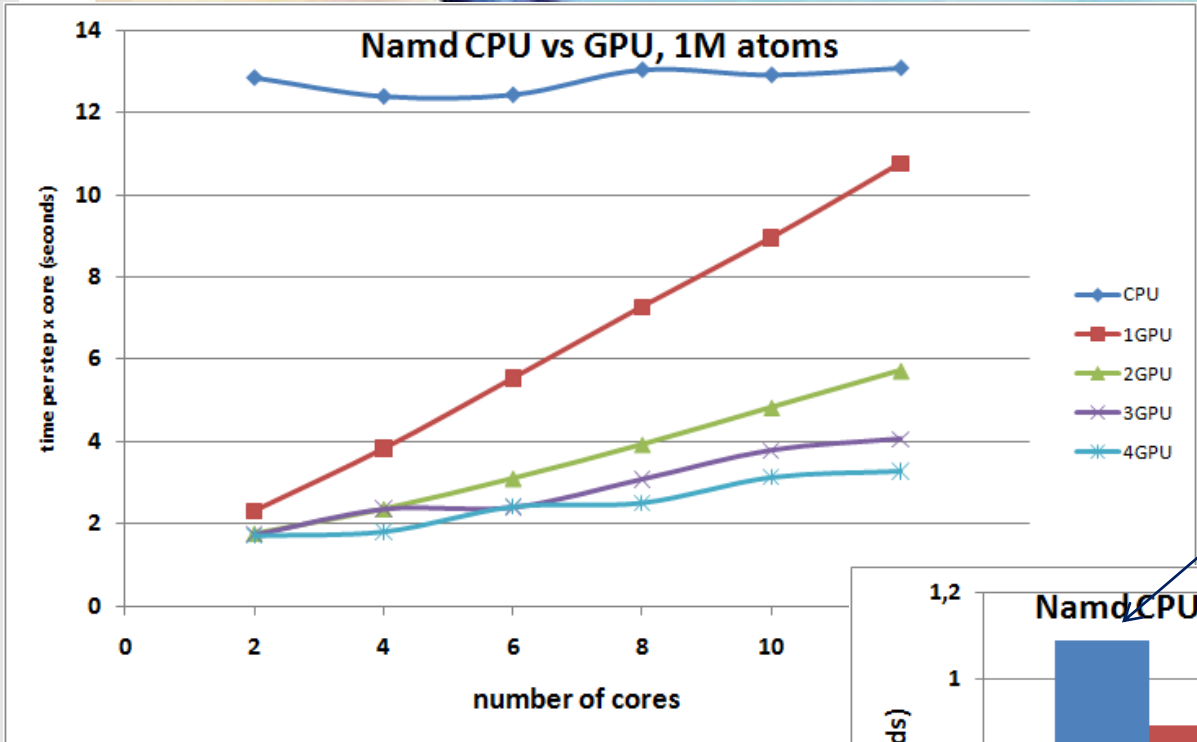
Namd CPU vs GPU, 23k atoms



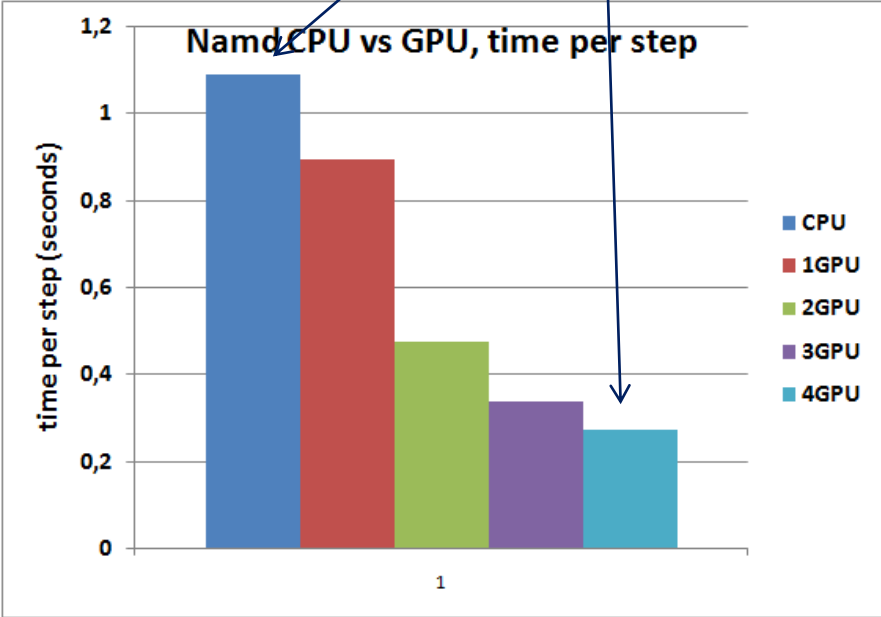
NO SPEEDUP

Namd CPU vs GPU, time per step



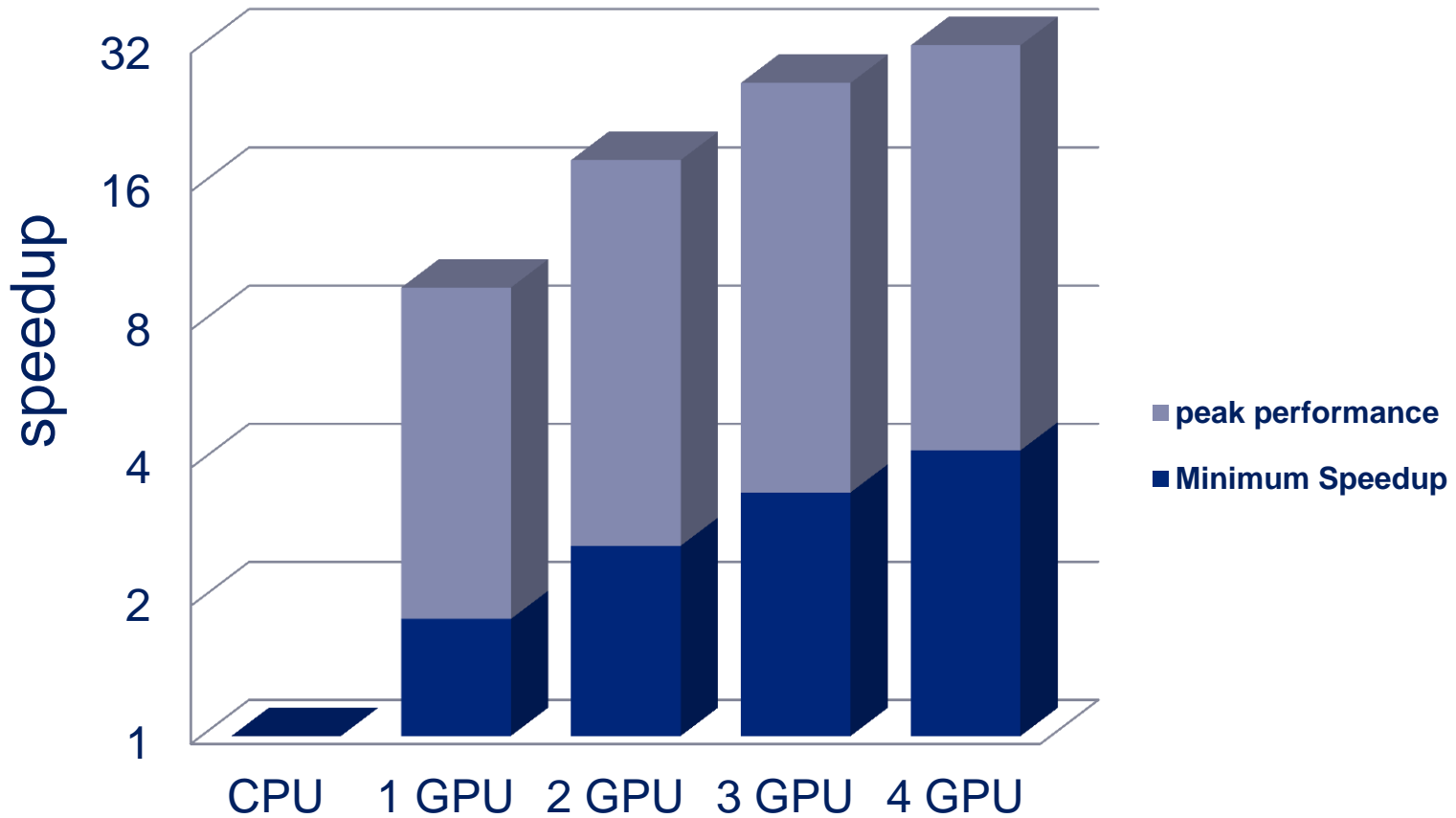


Speedup 4x





Minimum & Maximum Speedup





Features and Benefits

- Featuring up to four NVIDIA Tesla 20-Series GPUs per 1U server
- Supports up to two Quad/Six-Core Intel Xeon processor 5600 series per server
- Up to 96GB of DDR3 memory in 12 DIMM sockets
- Up to 3.0TB storage per server
- Two PCIe 2.0 x16 slots with riser card for GPU cards and one PCIe 2.0 x4 slot with riser
- Tool-less access to chassis, memory, HDDs, PCI card, blowers, and power supply
- Choice of Linux or Windows operating systems
- Offers a flexible, reliable and scalable compute platform
- Easy to maintain and service – hot-swappable drives and fans
- Energy-efficient - 1400W high-efficiency power supply and twelve cooling fans
- Best TCO with improved system density, computing capability while keeping the datacenter cool



server E7228+Tesla™

The new
supercomputing
technology

Host Server



2 x Tesla S1070
2 x Tesla S2050

TECHNICAL SPECIFICS TESLA™ S2050

2 x Tesla S2050 each one of which with the following features:

- 4 x GPU Tesla™
- 4 Teraflops in Single Precision
- 2 Teraflops in Double Precision
- GPU processor Clock 1.55 GHz peak clock
- Memory: up to 24 GB DDR5
- System I/O: two PCIe connection
Each connection leads to two of 4 GPUs

TECHNICAL SPECIFICS E7228

4 Motherboards Dual Socket Intel® Xeon® Six-Core serie 5600, each one of them with the following features:

- 1 x chipset Intel® 5520, 1333/1066/800 MHz
- 1 x 96GB REG ECC DDR3 1333/1066/800MHz (12 DIMMs)
- 1 x PCI-Express 16x 2.0 (Low Profile)
- 3 x SAS/SATA drive bays (2 bays)
- 1 x Intel® ICH10R – SATA II SW Raid - 6 ports (integrated)
- 1 x LAN Intel® Dual Gigabit 82576 controllers
(2 LAN ports total)
- 1 x Matrox G200eW
- 1 x sets of rear I/O ports including 2 USB 2.0, VGA, COM
- 1 x sets of IPMI 2.0 with dedicated Realtek 10/100Mb/s LAN port, KVM-over-LAN

E7228 is also equipped with:

- 2 x 1400W PFC power supply
Independent power control each node has its own power management

workstation E7095



PERSONAL SUPERCOMPUTER
NVIDIA TESLA™

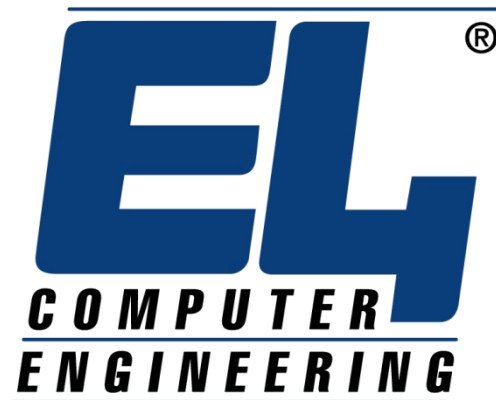


E7095 TECHNICAL FEATURES

- Workstation tower/rackmount low noise with optimized cooling and hetsink
- 8 bays SAS/SATA hot swap
- 1400W redundant power supply
- M/B dual Intel® Xeon®
- Chipset Intel® 5500
- CPU Intel® Xeon® serie 5600
- Up to 96 GB memory DDR3 reg ECC
- DVD writer dual layer
- Up to 4 x Tesla C 1060/C2050 or 2x S1070/S2050
- Graphic adapter NVIDIA Quadro NVS o FX



Thank you



E4 Computer Engineering SpA

Via Martiri della Liberta' 66

42019 - Scandiano (RE), Italy

www.e4company.com

Switchboard: +39.0522.991811





QUESTIONS?

