

# MSc CSTE, Fall 2018

## High Performance Computing

I. Moulitsas

Cranfield University

# Objective

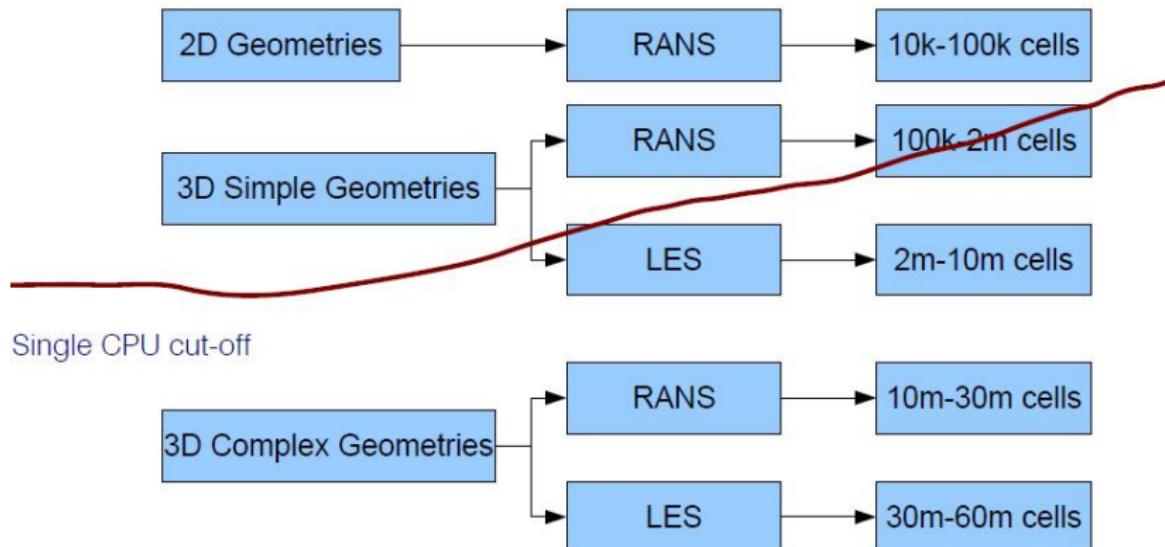
In this session we will discuss:

- Motivation for HPC
- Architecture and requirements of a typical small-scale HPC
- Scalability
- Hardware/Software choice & state-of-the-art
- Commercial codes benchmarks on HPC clusters
- MPI

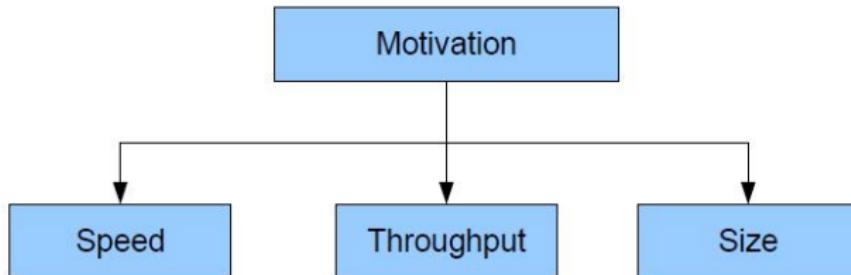
# Recommended Literature

- P.S. Pacheco, "Parallel Programming With MPI", Morgan Kaufmann, 1997
- W. Gropp et al, "Using MPI. Portable Parallel Programming with Message-Passing Interface", MIT Press, 1999
- R. Chandra et al, "Parallel Programming in OpenMP", Morgan Kaufmann, 2001
- "Parallel Methods in Numerical Analysis", Stanford University lecture notes, <http://www.stanford.edu/cla>
- S. Booth, J. Fishera, N. MacDonald, P. Maccallum, J. Malard, A. Ewing, E. Minty, A. Simpson, S. Paton, and S. Breuer. Introduction to the T3D. A one day course. Edinburgh Parallel Computing Centre, The University of Edinburgh, 2000.  
<http://www.es.embnet.org/Doc/Computing/index.html>

# Motivation for HPC in Scientific Computing Applications

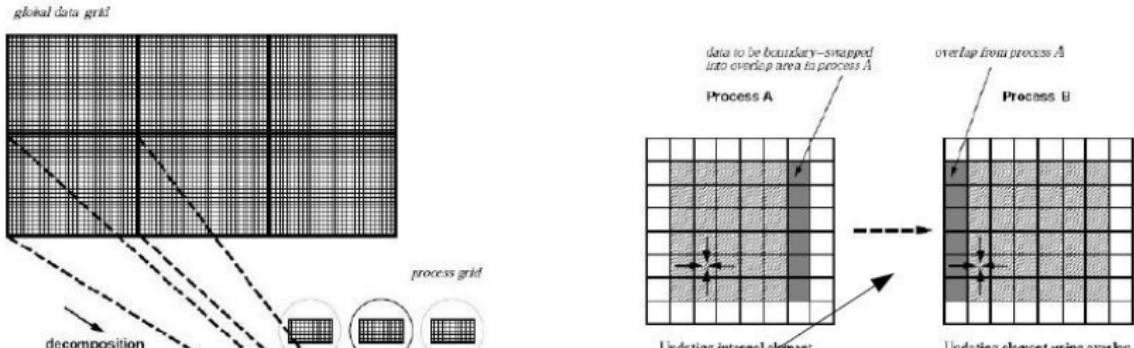


# Motivation for HPC in Scientific Computing Applications



- Enable multiple analysis in a fixed amount of time
- Decrease time necessary to complete one solution
- Enable higher grid resolution than possible in single processor machines
- Increase the accuracy of modelling of our physical system
- Introduce additional physical models that were impossible to tackle before

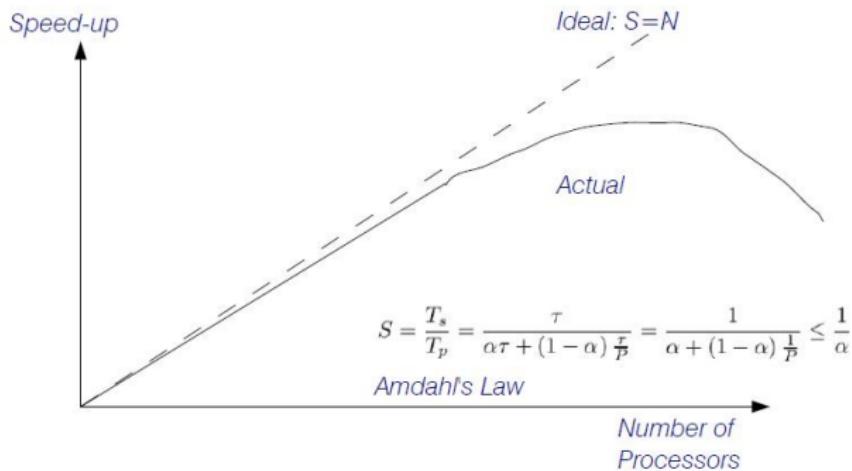
# Structured Grid Problem Breakdown



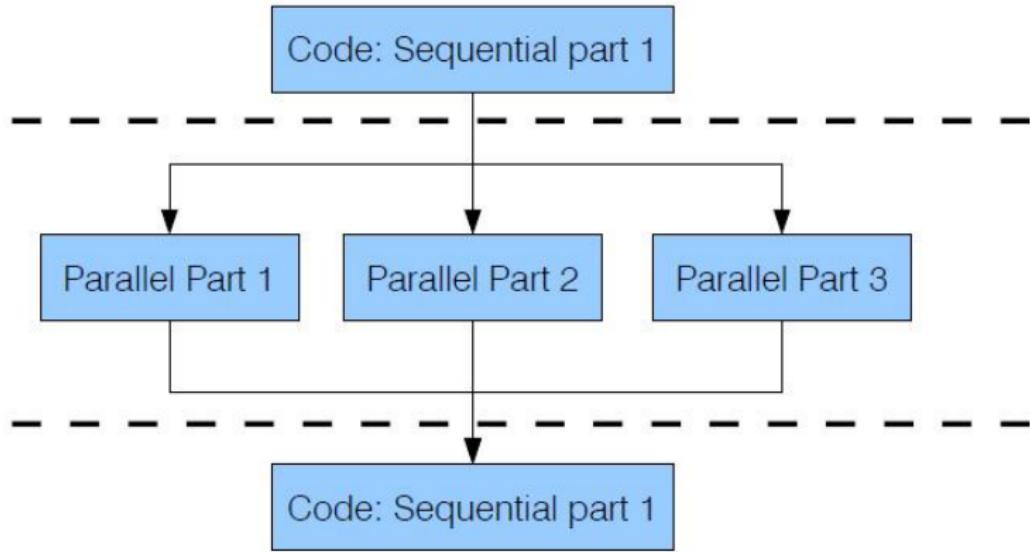
*Challenge 1: Communication. Mostly: Message Passing Interface (MPI)*

# Scalability

- Parallel speed-up : the ratio of the execution time of the parallel algorithm on a single processor to the execution time of the parallel algorithm on P processors.
- Parallel efficiency - Speed-up divided by P.

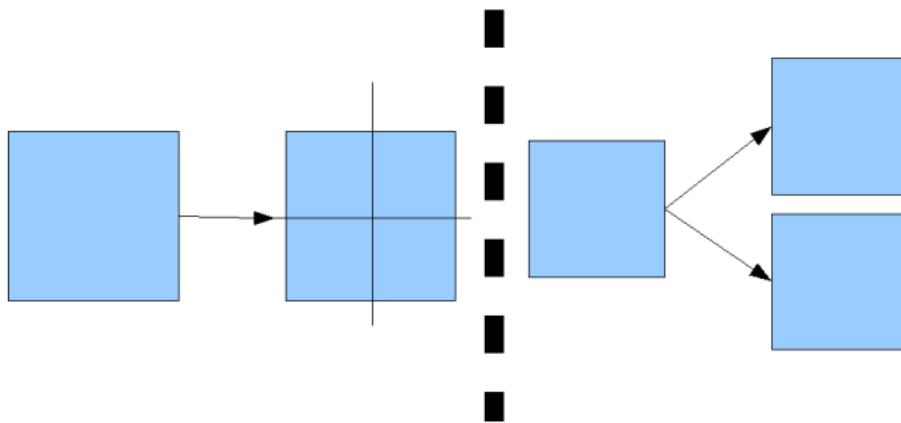


# Speed-Up



# Scalability: Weak vs Strong

- *Strong scalability* - fixed total problem size
- *Weak scalability* - based on fixed problem size per core



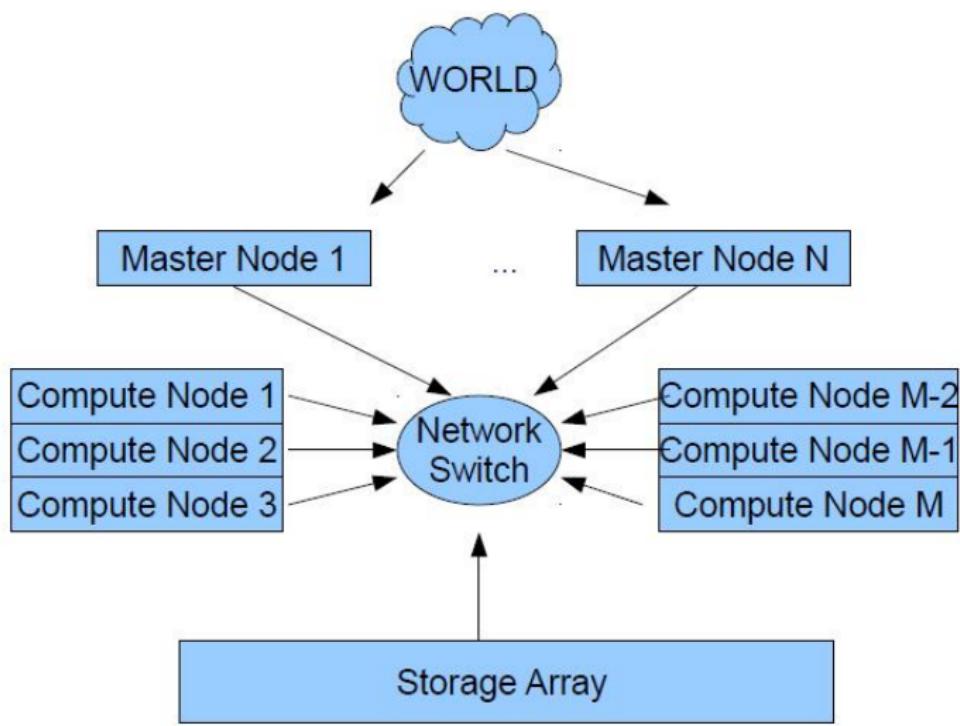
# Architecture of a Small-Scale Cluster

A compute cluster includes of the following hardware:

- Worker nodes
- Master node
- Storage array
- Network switches
- Power supply/protection
- Rack cabinet(s)

# Communications Scheme

Omitting power and rack, the simplified topology of a cluster is as follows:



# Initial Requirements

Initial considerations to take into account:

- Problem-based
  - ▶ Estimate size
  - ▶ Estimate speed
  - ▶ Estimate scalability
  - ▶ Postprocessing?
- Budget-based
  - ▶ Estimate typical usage pattern.
  - ▶ Limiting factors: CPU/Memory/Interconnect

# Costs

Considerations to take into account in the costs:

- Procurement:

- ▶ Hardware
- ▶ Software
- ▶ Space
- ▶ Installation

- Running costs:

- ▶ Power
- ▶ Air conditioning
- ▶ Support

- Lifetime

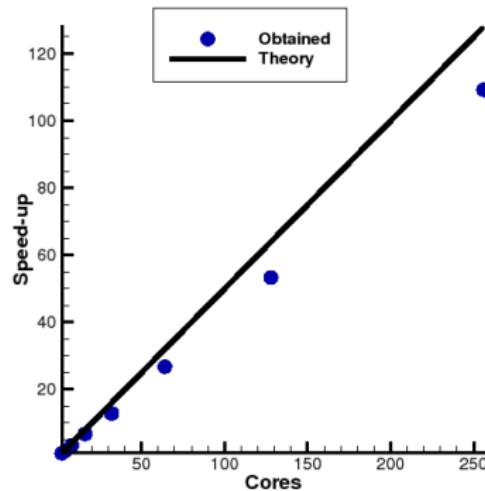
# Scalability: Developer's Perspective

## Question

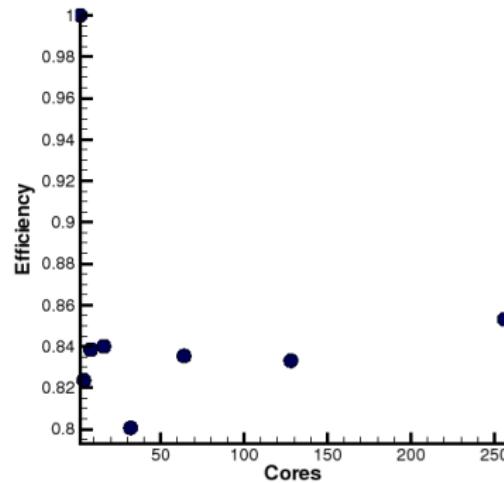
What are the key elements of an explicit compressible (or incompressible artificial compressibility based) solver and what factors can limit strong scalability of such solver?

# Scalability Example: In-house Incompressible

Astral HPC - an HP DL140 G3 cluster utilising a Xeon 51x0 (Woodcrest) 3GHz processors with Infiniband interconnect

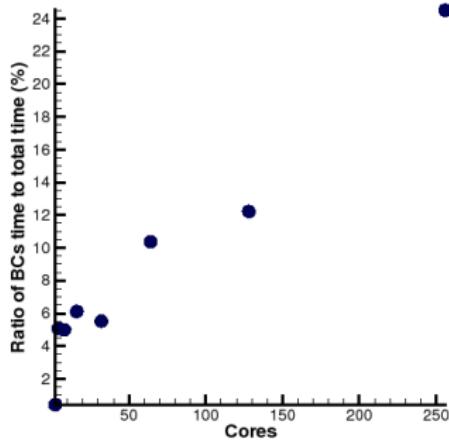


(a) Speed-up

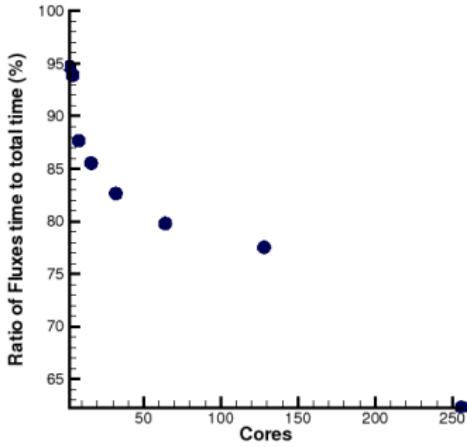


(b) Efficiency

# Limitations on Scalability



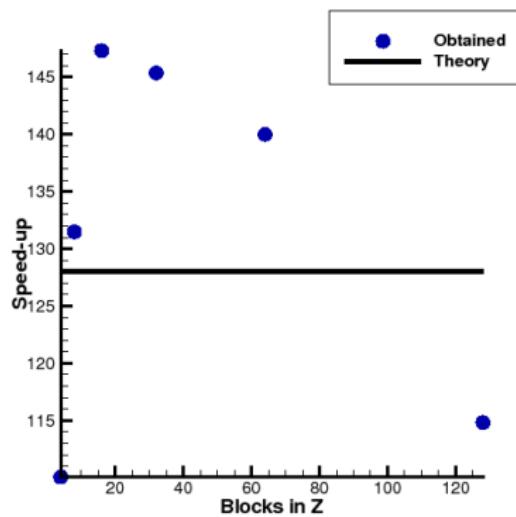
(a) Time spent in BCs



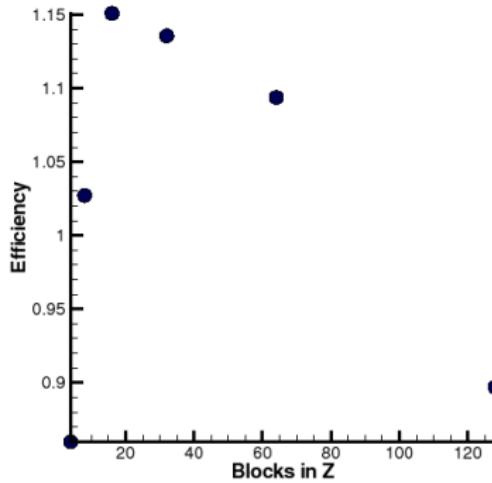
(b) Time spent in fluxes

# Cache Effects

Consider a fixed problem size and fixed number of cores (128)



(a) Speed-up



(b) Efficiency

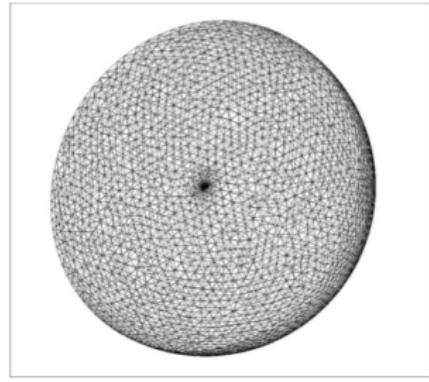
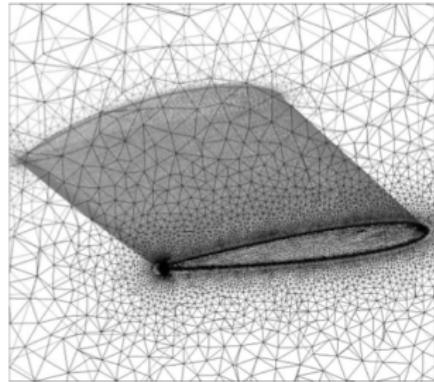
# Scalability Example: Flow over a UAV morphing wing

- In-house massively scalable solver for unstructured grids
- Hybrid Unstructured 3D grid

## Question

Find the optimal split for the computational work

**Limitation:** Data is too big to handle directly. Problem is NP-hard

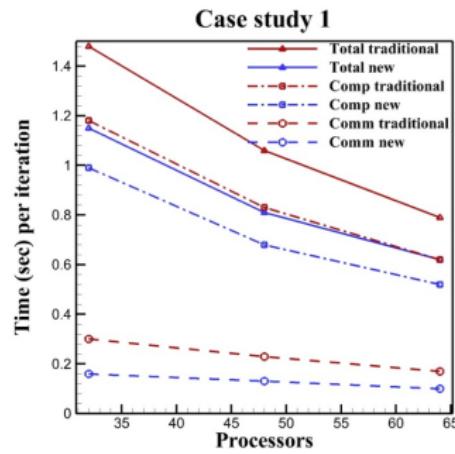


# Scalability Example and Load Balancing

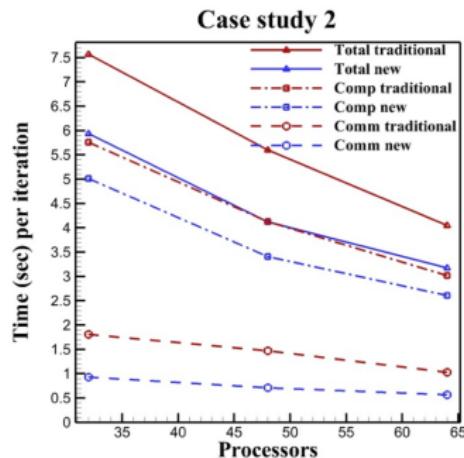
- 50% Reduction in communication cost
- 20% Reduction in computational cost
- Solution is obtained quickly
- Efficient use of computational resources

Times shown are computational cost and communication cost per iteration. Massive performance improvement when considering that several thousands of iterations are needed to achieve convergence.

WENO 3<sup>rd</sup> order



WENO 5<sup>th</sup> order



## Worker Nodes/Master Node

Worker nodes are bought for number-crunching. Therefore the properties you are interested in most are:

- CPU
- Memory
- Overall performance
- Local disk
- The total power output

Most of the servers at present time will include 2-8 CPUs with 2-4 cores per CPU. Master node usually does not participate in the parallel computations and is used for debugging, compilation, post-processing.

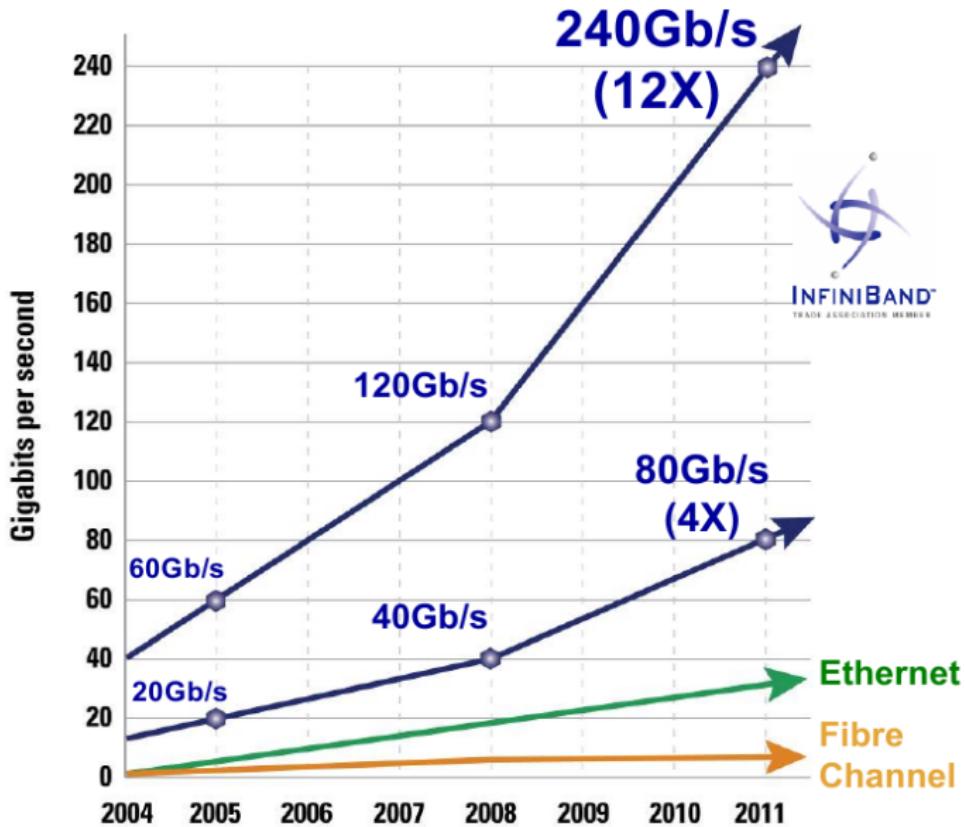
# Storage

- Storage can be either located on the master node itself, however this approach limits the expandability of the cluster. For this reason, the usual setup includes a separate Network-Attached Storage (NAS) - a dedicated data server.
- NAS server will have a separate processor and will support Redundant Array of Independent Disks (RAID) technology.

# Storage & RAID

- RAID arrays write data to multiple disks “stripes” keeping data highly distributed and, if needed writing redundant data (parity). Different levels e.g.
- RAID 0: Striped Set (2 disks minimum) without parity: improved performance but no fault tolerance .
- RAID 1: Mirrored Set (2 disks minimum) without parity: provides fault tolerance, may increase performance on read but decreases on writes.
- RAID 5: Striped Set (3 disk minimum) with Distributed Parity: requires all but one drive to be present to operate. If a drive fails, information restored from the distributed parity.
- Galileo’s storage is: ***Infortrend External RAID unit with 8 x 400GB discs configured in RAID 5 mode***

# Networking: Speed



# Power

- A power surge or the loss of power may result in the loss of data and, what is worse, - in corruption of the storage.
- Uninterruptable Power Supply (UPS) provides:
  - ▶ Power filtering
  - ▶ Notifications to your nodes/storage in the event of power fault
  - ▶ Backup power for the time sufficient in order for the system to shut down gracefully.

It can be too expensive and not quite necessary to shield all client nodes, but it is necessary to shield master node and storage.

# Software: Building Blocks

## Building blocks

- System
  - ▶ OS
  - ▶ Cluster management system
  - ▶ Job scheduling system
- User
  - ▶ Compilers
  - ▶ MPI libraries
  - ▶ Application software

# Software: OS Family

Top 500:

<b>Operating system Family</b>	<b>Count</b>	<b>Share %</b>
Linux	446	89.20%
Unix	25	5.00%
Mixed	23	4.60%
Windows	5	1.00%
BSD Based	1	0.20%

# Software: OS Type

Top 500:

Operating System	Count	Share %
Linux	391	78.20%
AIX	22	4.40%
CNK/SLES 9	20	4.00%
SLES10 + SGI ProPack 5	16	3.20%
CNL	12	2.40%
<u>CentOS</u>	8	1.60%
<u>SuSE Linux Enterprise Server 9</u>	5	1.00%
Windows HPC 2008	5	1.00%
<u>Redhat Linux</u>	4	0.80%
<u>SUSE Linux Enterprise Server 10</u>	4	0.80%
<u>RedHat Enterprise 4</u>	3	0.60%
<u>RedHat Enterprise 5</u>	2	0.40%
<b>Total</b>		<b>98.40%</b>

# Software: CMS & Job Schedulers

Cluster management system (CMS):

- Administration/installation/upgrade
- Source - Proprietary (e.g. IBM CMS:  
[www.ibm.com/systems/clusters/software/csm/index.html](http://www.ibm.com/systems/clusters/software/csm/index.html))
- Source - Open Source (e.g. Open Source Cluster Application Resources - [svn.oscar.openclustergroup.org](http://svn.oscar.openclustergroup.org), Rocks CMS - [www.rocksclusters.org](http://www.rocksclusters.org))

Job scheduler:

- SGE Gridengine - Sun, [www.sun.com/software/sge/](http://www.sun.com/software/sge/)
- Load Scheduling Facility (LSF) - Platform Computing,  
[www.platform.com](http://www.platform.com)

# Software: MPI

A number of implementations:

- Proprietary (e.g. HP-MPI, Intel MPI)
- OpenMPI ([www.open-mpi.org](http://www.open-mpi.org))
- MPICH2 ([www.mcs.anl.gov/research/projects/mpich2](http://www.mcs.anl.gov/research/projects/mpich2))
- LAM-MPI ([www.lam-mpi.org](http://www.lam-mpi.org))

# State of the Art

Top500, Nov 2012. Rmax and Rpeak in TFlops. Power in KW.

Rank	Site	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	DOE/SC/Oak Ridge National Laboratory United States	Titan - Cray XK7 , Opteron 6274 16C 2.200GHz, Cray Gemini Interconnect, NVIDIA K20x Cray Inc.	560640	17590.0	27112.5	8209
2	DOE/NSA/LLNL United States	Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom IBM	1572864	16324.8	20132.7	7890
3	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VIIIfx 2.00GHz, Tofu interconnect Fujitsu	705024	10510.0	11280.4	12660
4	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	786432	8162.4	10066.3	3945
5	Forschungszentrum Juelich (FZJ) Germany	JUQUEEN - BlueGene/Q, Power BQC 16C 1.600GHz, Custom Interconnect IBM	393216	4141.2	5033.2	1970
6	Leibniz Rechenzentrum Germany	SuperMUC - iDataPlex DX360M4, Xeon E5-2680 8C 2.70GHz, Infiniband FDR IBM	147456	2897.0	3185.1	3423
7	Texas Advanced Computing Center/Univ. of Texas United States	Stampede - PowerEdge C8220, Xeon E5-2680 8C 2.700GHz, Infiniband FDR, Intel Xeon Phi Dell	204900	2660.3	3959.0	
8	National Supercomputing Center in Tianjin China	Tianhe-1A - NUDT YH MPP, Xeon X5670 8C 2.93 GHz, NVIDIA 2050 NUDT	186368	2566.0	4701.0	4040
9	CINECA Italy	Fermi - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	163840	1725.5	2097.2	822
10	IBM Development Engineering United States	DARPA Trial Subset - Power 775, POWER7 8C 3.836GHz, Custom Interconnect IBM	63360	1515.0	1944.4	3576

# State of the Art

Top500, June 2012. Rmax and Rpeak in TFlops. Power in KW.

Rank	Site	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	DOE/NSA/LNL United States	Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom IBM	1572864	16324.8	20132.7	7890
2	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 Vlllfx 2.0GHz, Tofu interconnect Fujitsu	705024	10510.0	11280.4	12660
3	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	786432	8162.4	10066.3	3945
4	Leibniz Rechenzentrum Germany	SuperMUC - iDataPlex DX360M4, Xeon E5-2680 8C 2.70GHz, Infiniband FDR IBM	147456	2897.0	3185.1	3423
5	National Supercomputing Center in Tianjin China	Tianhe-1A - NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 NUDT	186388	2566.0	4701.0	4040
6	DOE/SC/Oak Ridge National Laboratory United States	Jaguar - Cray XK6, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA 2090 Cray Inc.	298592	1941.0	2627.6	5142
7	CINECA Italy	Fermi - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	163840	1725.5	2097.2	822
8	Forschungszentrum Juelich (FZJ) Germany	JUQUEEN - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	131072	1380.4	1677.7	658
9	CEA/TGCC-GENCI France	Curie thin nodes - Bullx B510, Xeon E5-2680 8C 2.700GHz, Infiniband QDR Bull SA	77184	1359.0	1667.2	2251
10	National Supercomputing Centre in Shenzhen (NSCS) China	Nebulae - Dawning TC3600 Blade System, Xeon X5650 6C 2.66GHz, Infiniband QDR, NVIDIA 2050 Dawning	120640	1271.0	2984.3	2580



# State of the Art

Top500, November 2011. Rmax and Rpeak in TFlops. Power in KW.

Rank	Site	Computer/Year Vendor	Cores	R <sub>max</sub>	R <sub>peak</sub>	Power
1	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect / 2011 Fujitsu	705024	10510.00	11280.38	12659.9
2	National Supercomputing Center in Tianjin China	NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010 NUDT	186368	2566.00	4701.00	4040.0
3	DOE/SC/Oak Ridge National Laboratory United States	Cray XT5-HE Opteron 6-core 2.6 GHz / 2009 Cray Inc.	224162	1759.00	2331.00	6950.0
4	National Supercomputing Centre in Shenzhen (NSCS) China	Dawning TC3600 Blade System, Xeon X5650 6C 2.66GHz, Infiniband QDR, NVIDIA 2050 / 2010 Dawning	120640	1271.00	2984.30	2580.0
5	GSIC Center, Tokyo Institute of Technology Japan	HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU, Linux/Windows / 2010 NEC/HP	73278	1192.00	2287.63	1398.6
6	DOE/NNSA/LANL/SNL United States	Cray XE6, Opteron 6136 8C 2.40GHz, Custom / 2011 Cray Inc.	142272	1110.00	1365.81	3980.0
7	NASA/Ames Research Center/NAS United States	SGI Altix ICE 8200EX/8400EX, Xeon HT QC 3.0/Xeon 5570/5670 2.93 GHz, Infiniband / 2011 SGI	111104	1088.00	1315.33	4102.0
8	DOE/SC/LBNL/NERSC United States	Cray XE6, Opteron 6172 12C 2.10GHz, Custom / 2010 Cray Inc.	153408	1054.00	1288.63	2910.0
9	Commissariat a l'Energie Atomique (CEA) France	Bull bulx super-node S6010/S6030 / 2010 Bull	138368	1050.00	1254.55	4590.0
10	DOE/NNSA/LANL United States	BladeCenter QS22/LS21 Cluster, PowerXCell 8i 3.2.Ghz / Opteron DC 1.8 GHz, Voltaire Infiniband / 2009 IBM	122400	1042.00	1375.78	2345.0

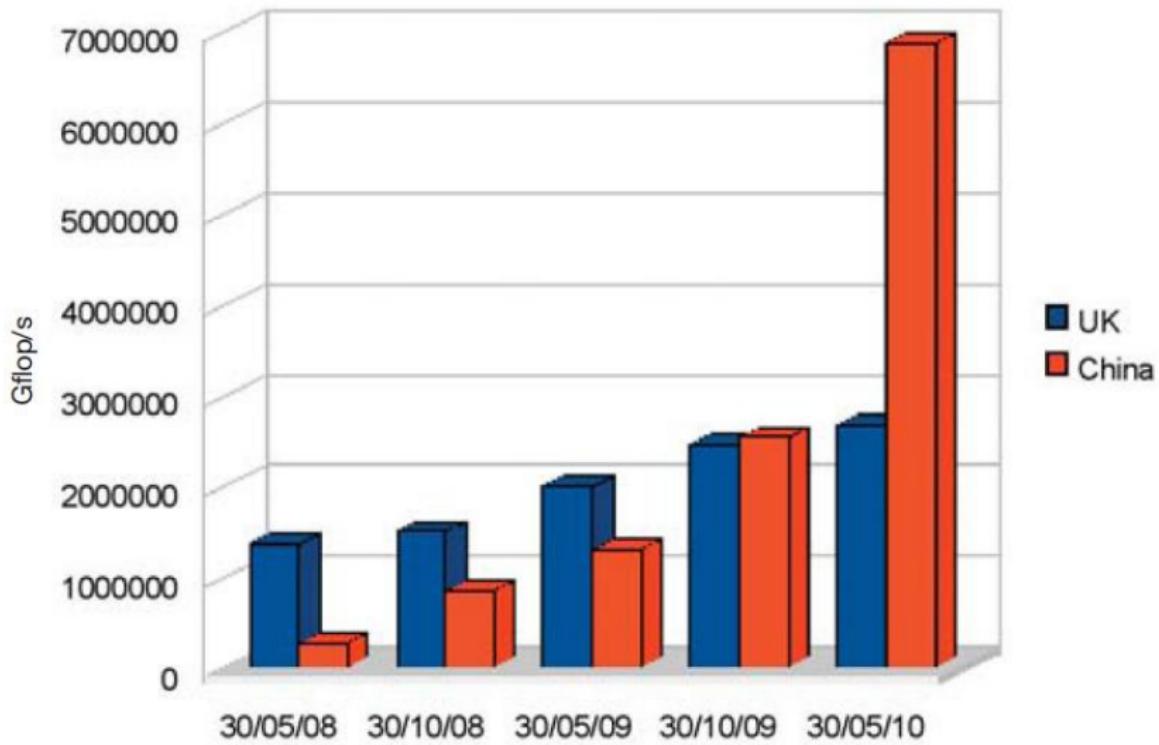


# State of the Art

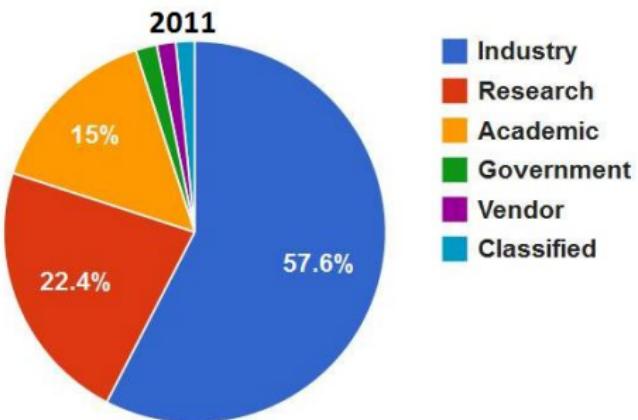
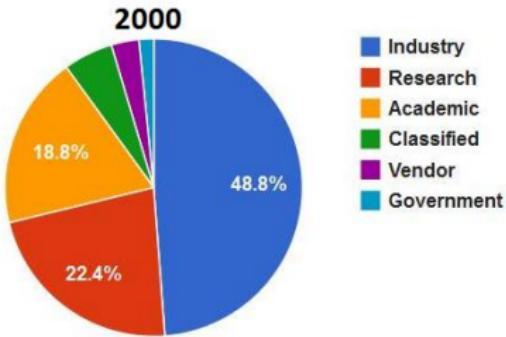
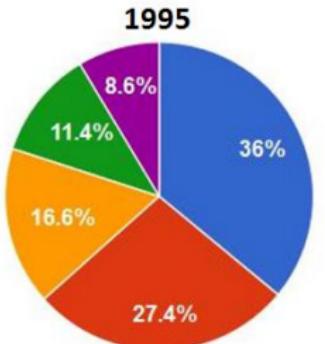
Top500, November 2010. Rmax and Rpeak in TFlops. Power in KW.

Rank	Site	System	Cores	R <sub>max</sub>	R <sub>peak</sub>
1	National Supercomputing Center in Tianjin China	NUDT TH MPP, X5670 2.93Ghz 6C, NVIDIA GPU, FT-1000 8C NUDT	186368	2566	4701
2	DOE/SC/Oak Ridge National Laboratory United States	Cray XT5-HE Opteron 6-core 2.6 GHz Cray Inc.	224162	1759	2331
3	National Supercomputing Centre in Shenzhen (NSCS) China	Dawning TC3600 Blade, Intel X5650, NVidia Tesla C2050 GPU Dawning	120640	1271	2984.3
4	GSIC Center, Tokyo Institute of Technology Japan	HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU, Linux/Windows NEC/HP	73278	1192	2287.63
5	DOE/SC/LBNL/NERSC United States	Cray XE6 12-core 2.1 GHz Cray Inc.	153408	1054	1288.63
6	Commissariat a l'Energie Atomique (CEA) France	Bull bulx super-node S6010/S6030 Bull SA	138368	1050	1254.55
7	DOE/NNSA/LANL United States	BladeCenter QS22/LS21 Cluster, PowerXCell 8i 3.2 Ghz / Opteron DC 1.8 GHz, Voltaire Infiniband IBM	122400	1042	1375.78
8	National Institute for Computational Sciences/University of Tennessee United States	Cray XT5-HE Opteron 6-core 2.6 GHz Cray Inc.	98928	831.7	1028.85
9	Forschungszentrum Juelich (FZJ) Germany	Blue Gene/P Solution IBM	294912	825.5	1002.7
10	DOE/NNSA/LANL/SNL United States	Cray XE6 8-core 2.4 GHz Cray Inc.	107152	816.6	1028.66

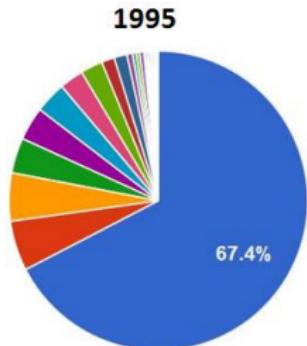
# State of the Art



# State of the Art: Segment

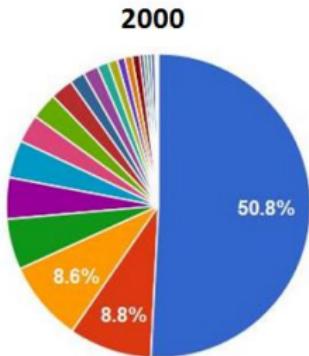


# State of the Art: Application Areas



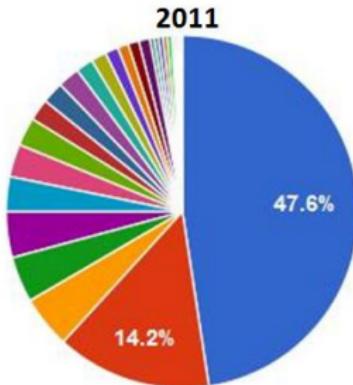
- Not Specified
- Weather and ...
- Geophysics
- Automotive
- Aerospace
- Chemistry
- Benchmarking
- Software

▲ 1/3 ▼



- Not Specified
- Telecomm
- Finance
- Weather and C...
- Database
- Automotive
- WWW
- Geophysics

▲ 1/3 ▼

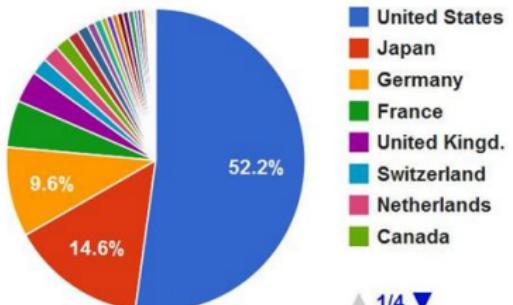


- Not Specified
- Research
- Finance
- Service
- Logistic Services
- Defense
- WWW
- Geophysics

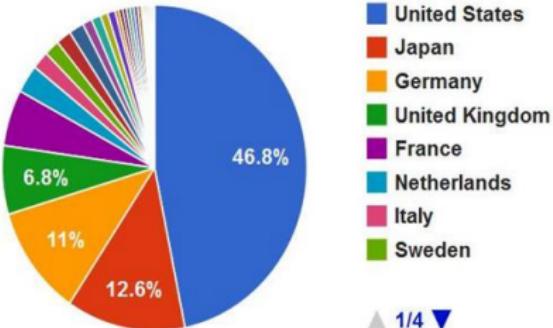
▲ 1/4 ▼

# State of the Art: Countries

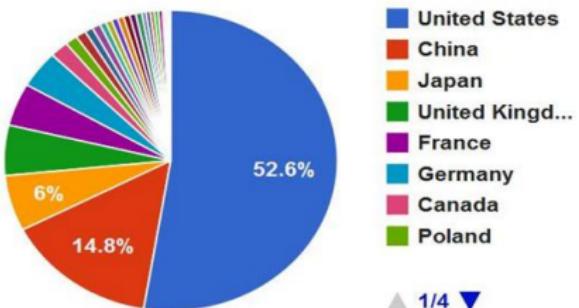
1995



2000

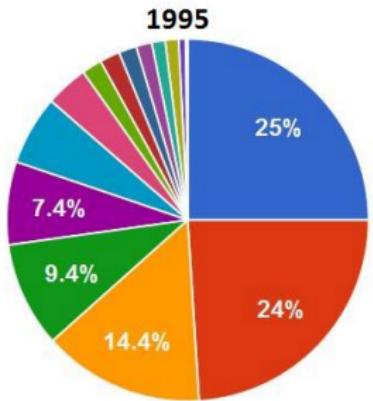


2011

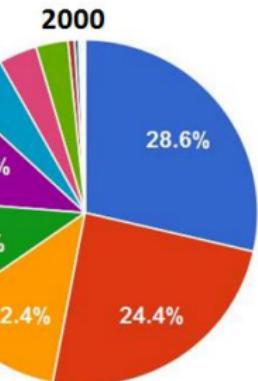


▲ 1/4 ▼

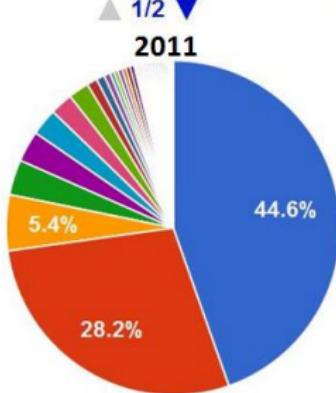
# State of the Art: Vendors



Cray Inc.  
SGI  
IBM  
Intel  
TMC  
Fujitsu  
NEC  
KSR



IBM  
Oracle  
SGI  
Cray Inc.  
HP  
NEC  
Fujitsu  
Hitachi

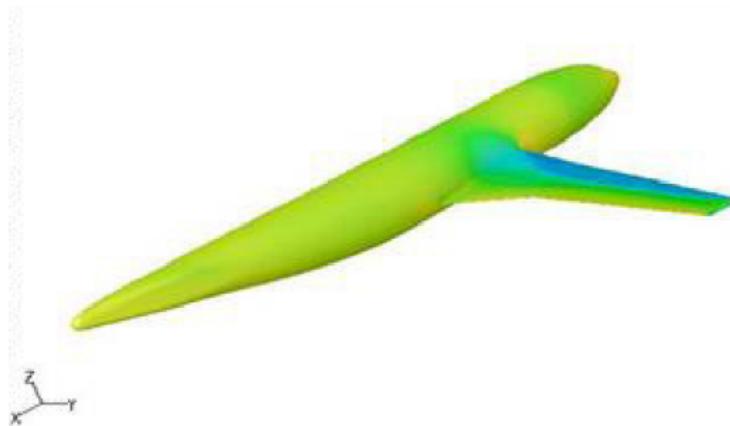


IBM  
HP  
Cray Inc.  
SGI  
Bull  
Appro  
Dell  
Oracle

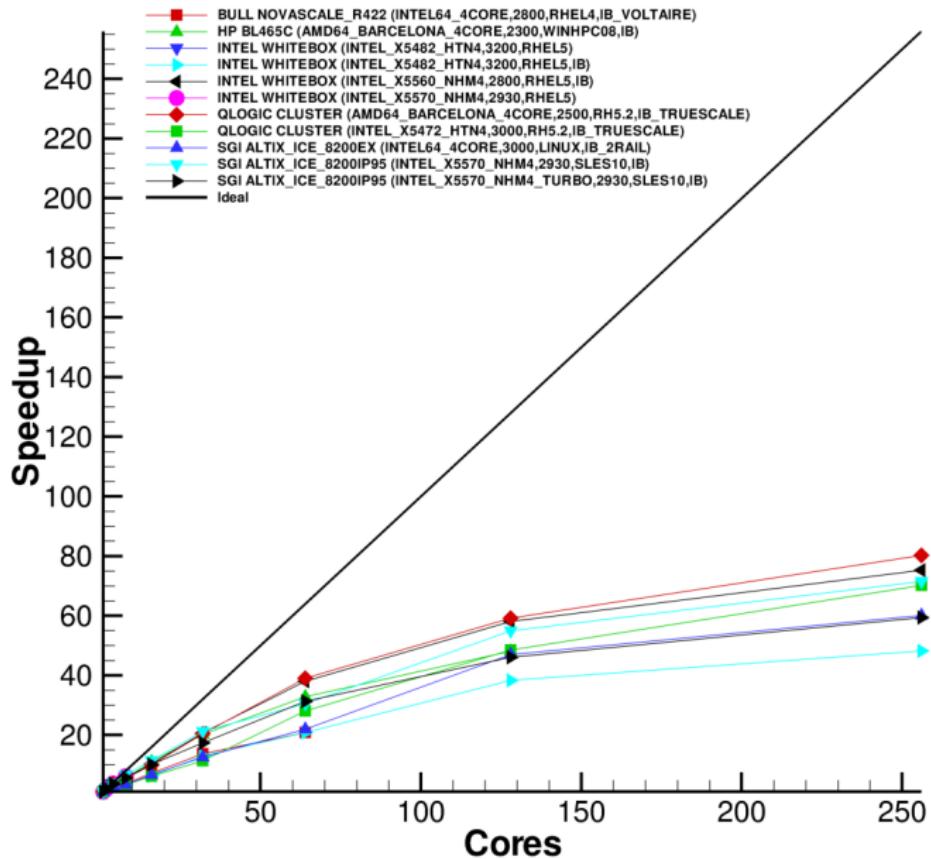
# ANSYS FLUENT 12: Small grid

May 2009. Test case: Aircraft wing

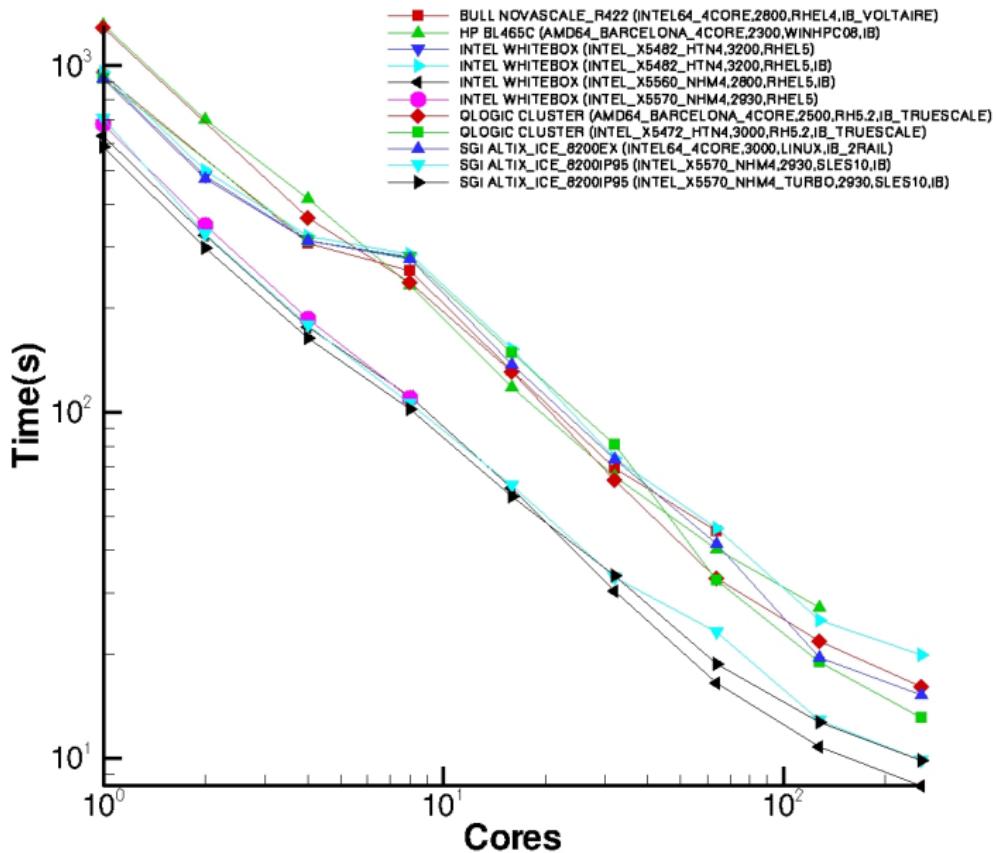
- Grid: 1,800,000 cells, hexahedral
- Model: realizable k- $\epsilon$ /coupled implicit solver



## FLUENT 12, Wing flow,k-e 1,8m cells



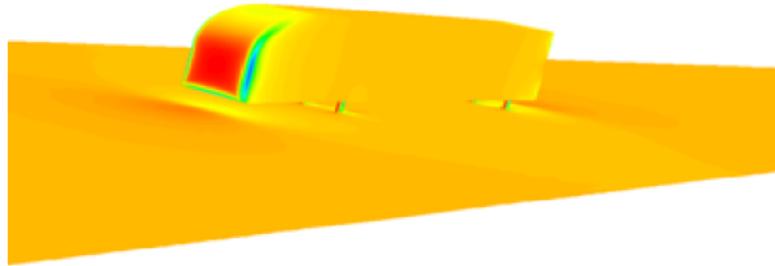
## FLUENT 12, Wing flow,k-e 1,8m cells



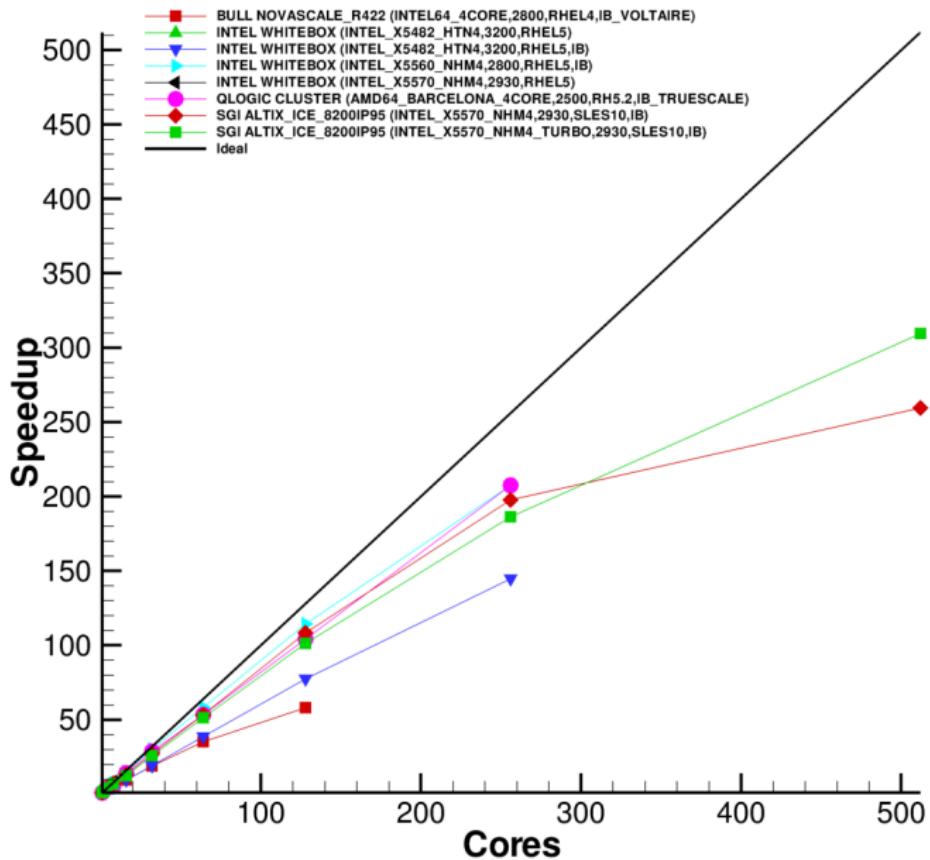
# ANSYS FLUENT 12: Large Grid

May 2009. External flow over a truck body

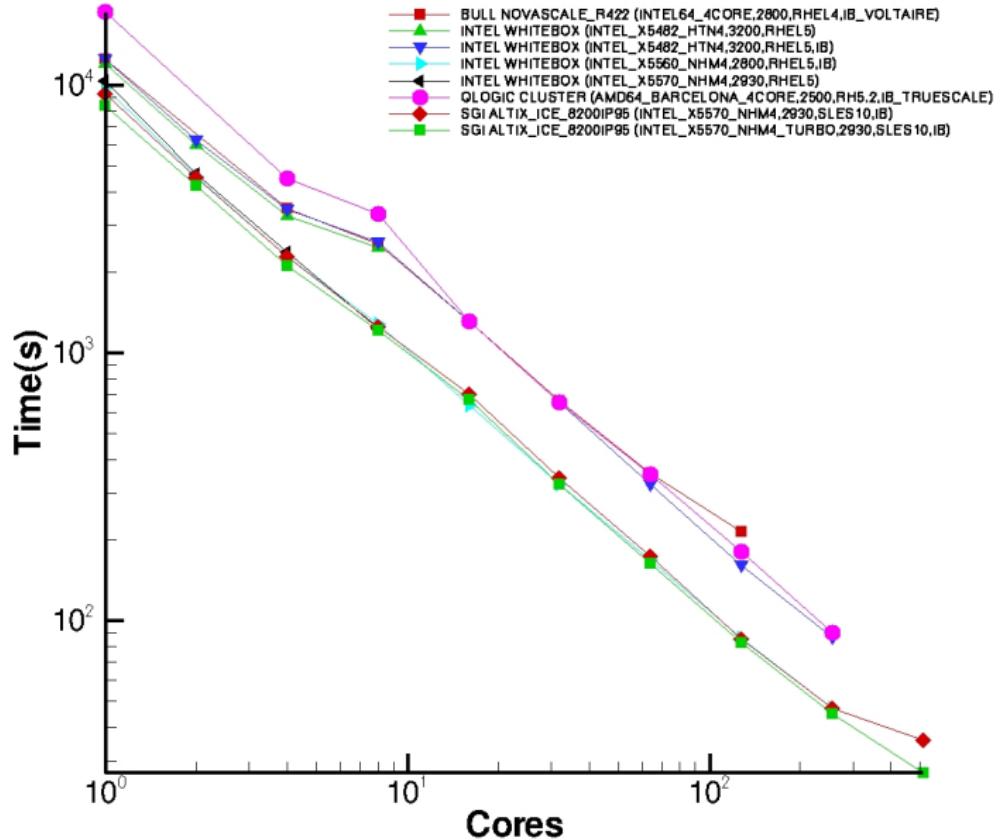
- Grid: 14,000,000 cells, mixed
- Model: DES/segregated implicit solver



## FLUENT 12, Truck, DES, 14m cells



# FLUENT 12, Truck, DES, 14m cells



# Evolution of Performance (May 2009)

Test system (Source: HPC Advisory council report,  
[www.hpcadvisorycouncil.com](http://www.hpcadvisorycouncil.com))

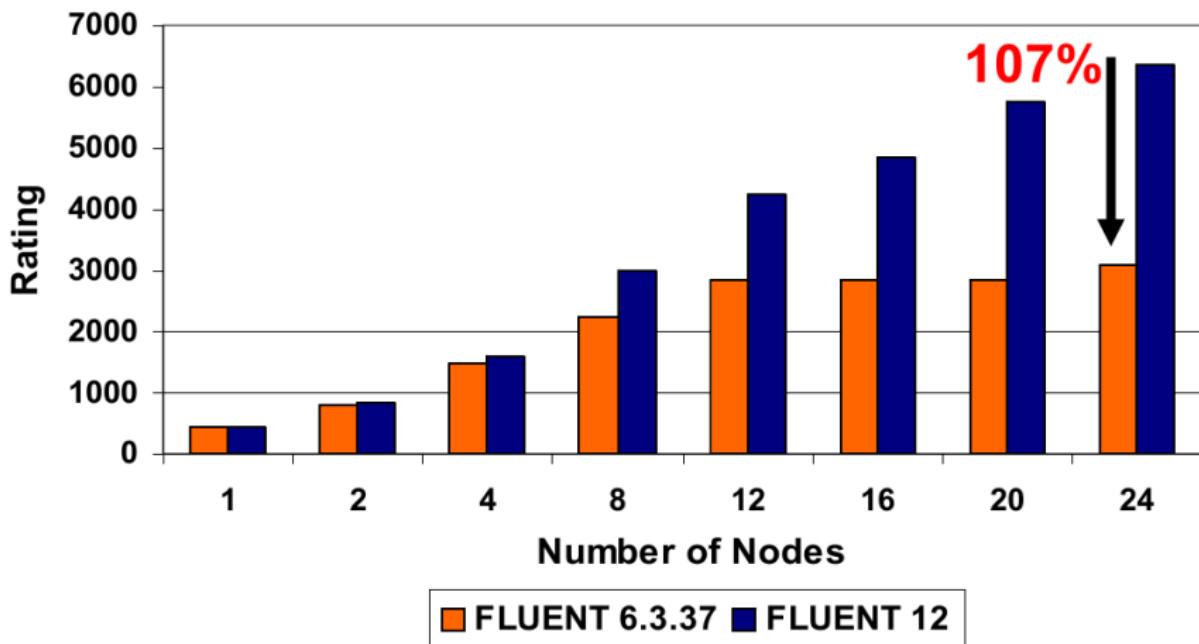
- Dell™ PowerEdge™ SC 1435 24-node cluster
- Quad-Core AMD Opteron™ 2382 (“Shanghai”) CPUs
- Mellanox® InfiniBand ConnectX® 20Gb/s (DDR) HCAs
- Mellanox® InfiniBand DDR Switch
- Memory: 16GB memory, DDR2 800MHz per node
- OS: RHEL5U2
- MPI: HP-MPI 2.3
- Application: FLUENT 6.3.37, FLUENT 12.0



# Evolution of Performance

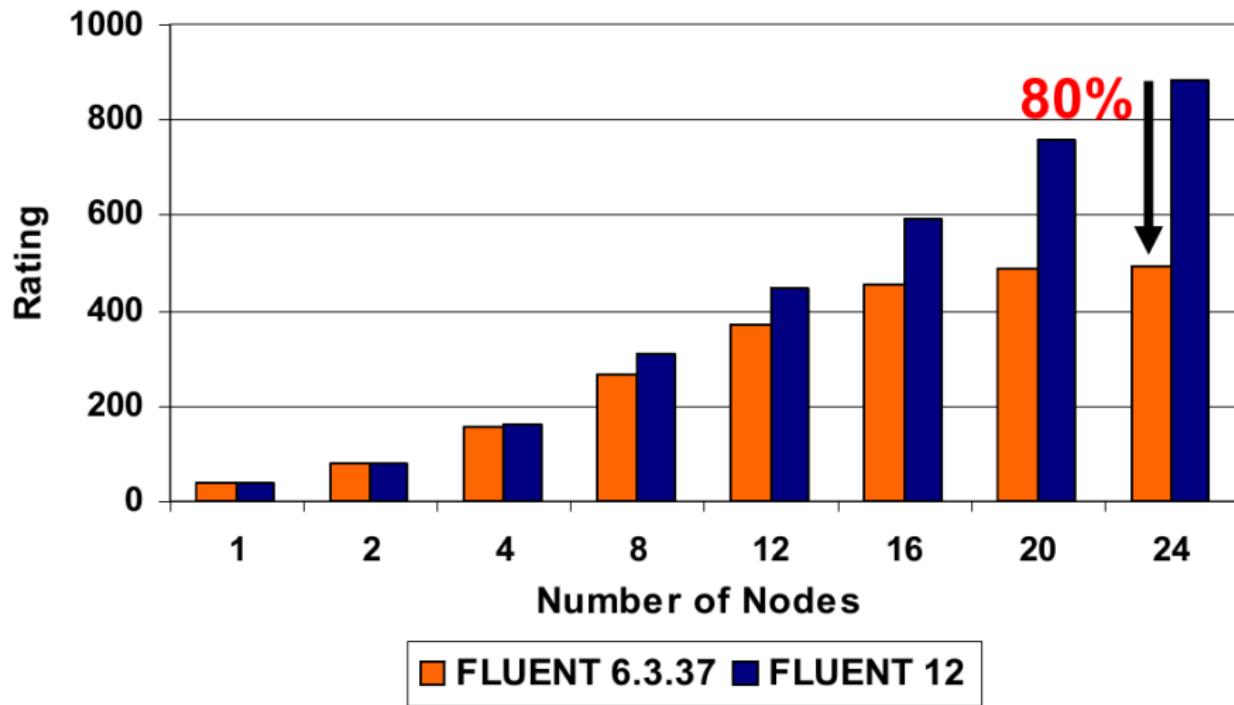
*Rating* - The number of benchmarks that can be run on a given machine (in sequence) in a 24 hour period.

**FLUENT Benchmark Result  
(Aircraft\_2M)**



# Evolution of Performance

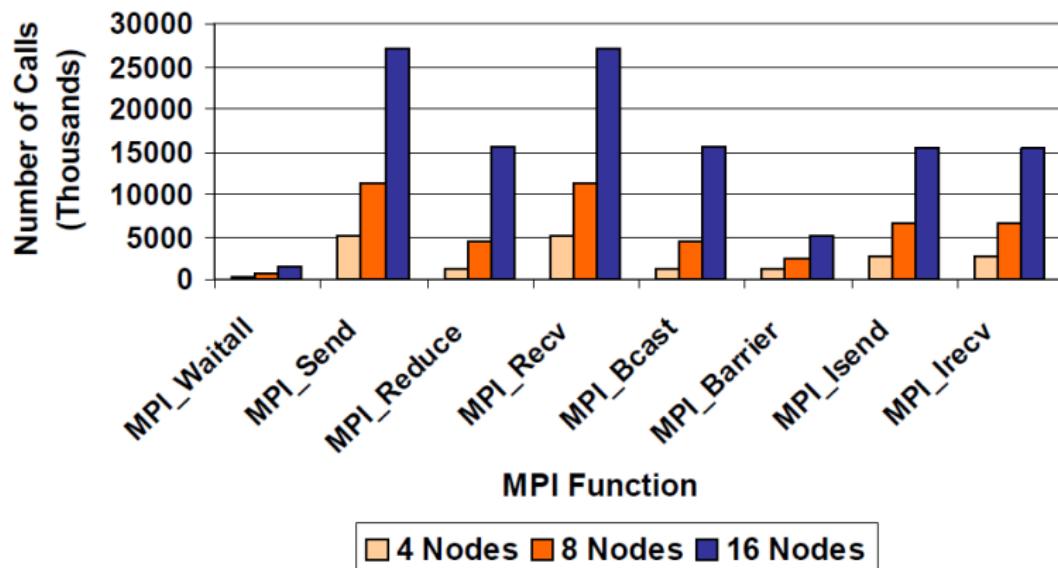
## FLUENT Benchmark Result (Truck\_14M)



# Evolution of Performance

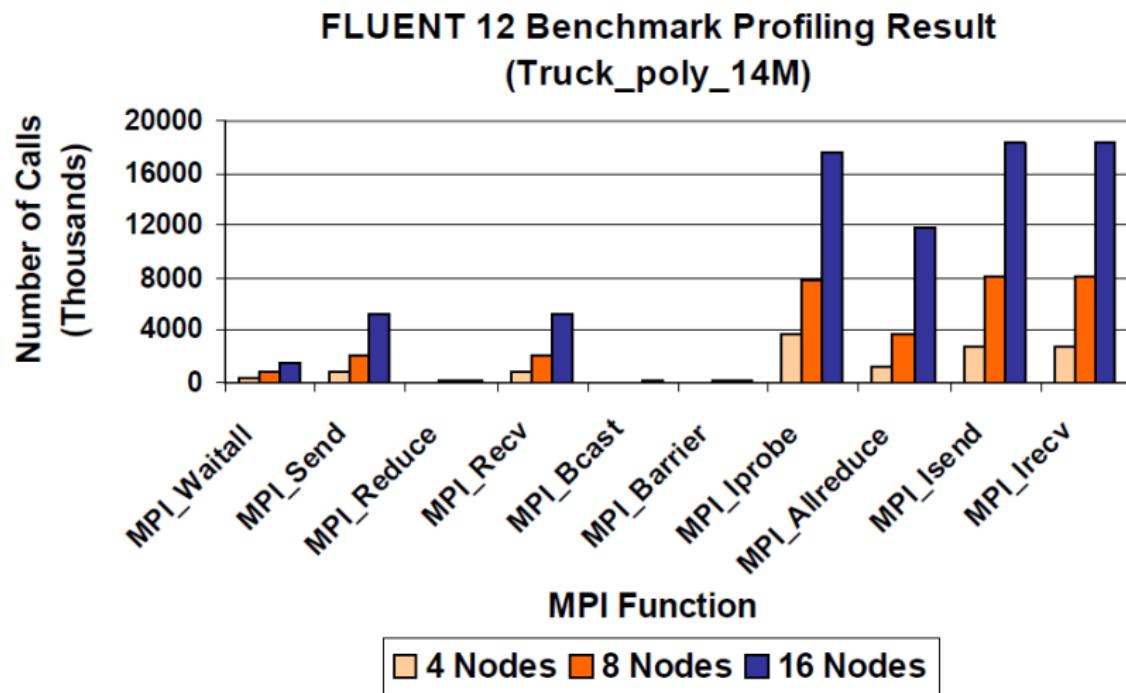
FLUENT 6.3.37 (Truck case) - Most used functions: MPI\_Send, MPI\_Recv, MPI\_Reduce, and MPI\_Bcast

FLUENT 6.3.37 Benchmark Profiling Result  
(Truck\_poly\_14M)



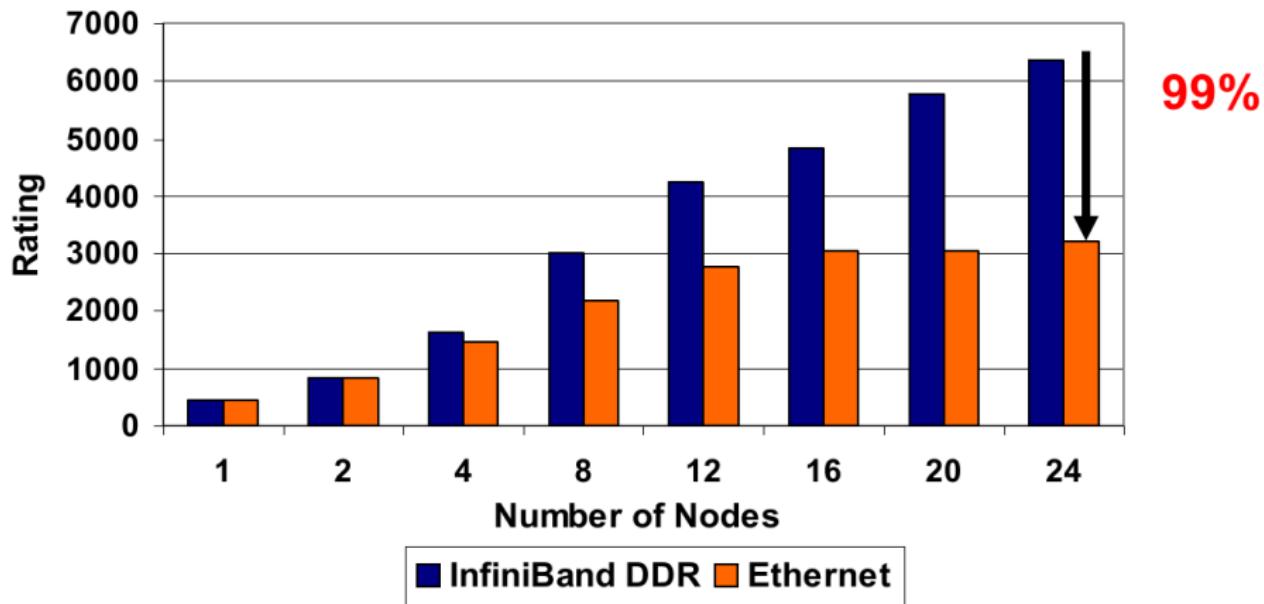
# Evolution of Performance

FLUENT 12 (Truck case) - Most used functions: MPI\_Iprobe, MPI\_Allreduce, MPI\_Isend, and MPI\_Irecv



# Effect of the Interconnect

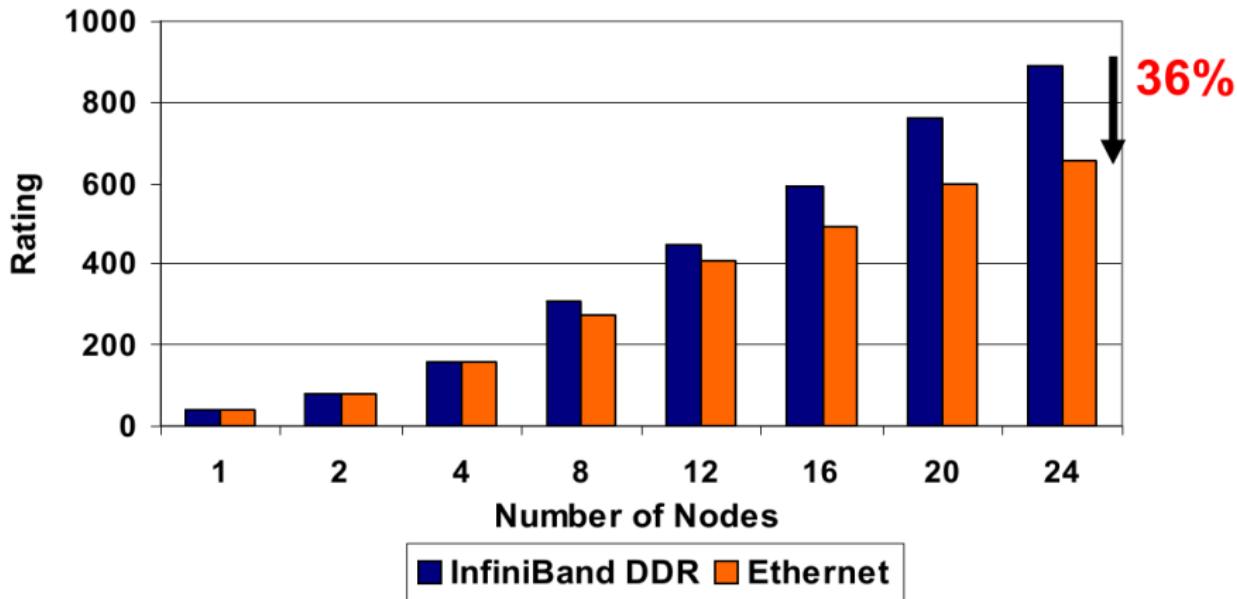
**FLUENT 12.0 Benchmark Result  
(Aircraft\_2M)**



99%

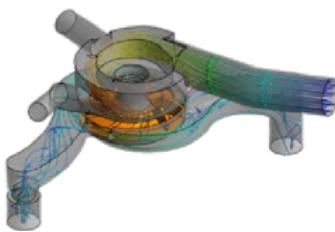
# Effect of the Interconnect

**FLUENT 12.0 Benchmark Result  
(Truck\_14M)**



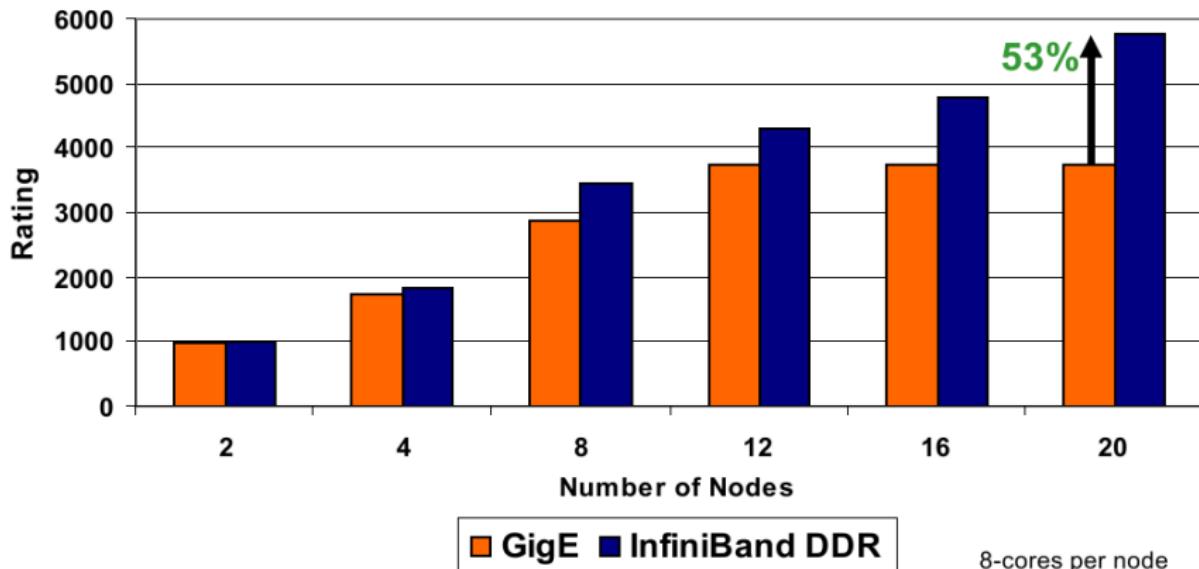
# ANSYS CFX: Pump Benchmark, July 2009

- Same hardware as above.
- Multiple frames of reference/Rotating and stationary components
- Unstructured mesh: ~600000 nodes



# ANSYS CFX: Pump Benchmark

**ANSYS CFX Benchmark Result  
(Pump)**



# STAR CCM+

- Vendor: CD-Adapco
- Integrated environment for engineering simulations
- CAD Embedding (SolidWorks, Pro/E, CATIA V5 or Unigraphics NX)

# Hardware

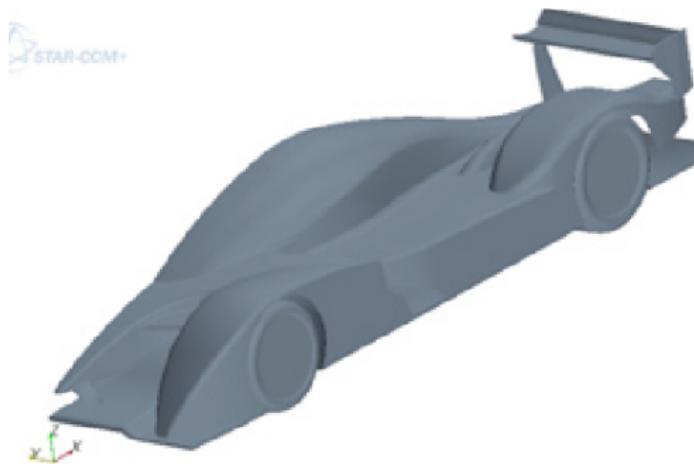
Cluster hardware (white paper: CD-adapco STAR-CCM+ with TrueScale):

- 32 compute nodes + Network File System (NFS) server node
- Dual Quad-Core AMD Opteron™ 2360 SE Barcelona CPU
- 16 GB of DDR2-667Mhz
- TrueScale InfiniBand/20-Gbps
- Gigabit Ethernet interconnect

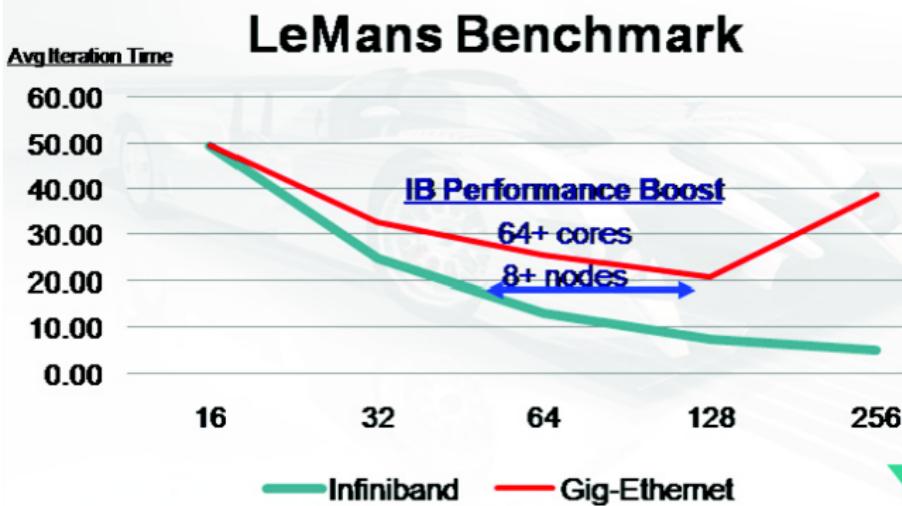
# Low Speed Benchmark

## LeMans Race Car Benchmark

- Number of cells: 17,000,000/polyhedral
- Solver: Segregated/Low speed external aerodynamics



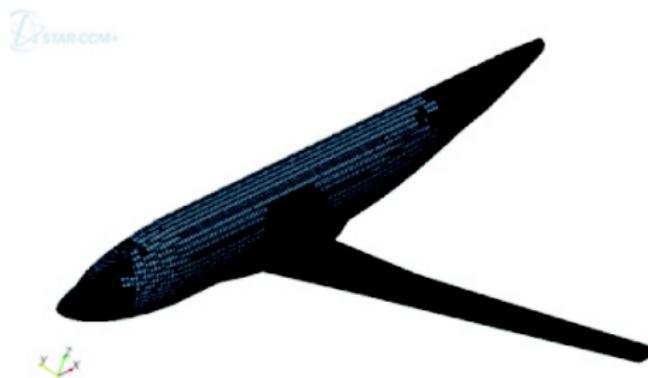
# Scalability



# High Speed Benchmark

## Civilian Airline Model

- Number of cells: 20,000,000/trimmed
- Solver: Segregated/High speed external aerodynamics (Mach 0.8)



# Scalability

