

National Aeronautics and Space Administration

Supercomputing A Ground Based Instrument for Exploration

Bob Ciotti

Supercomputing Systems Lead/System Architect

OFA14 Monterey

10010
10010
0001
010
010
10
10
00
0

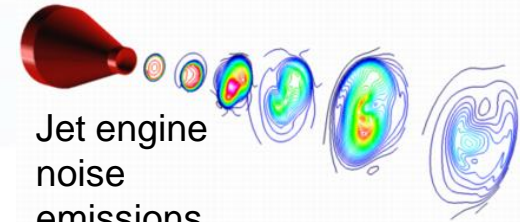
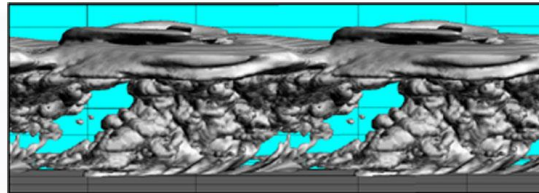
Supercomputing Support for NASA Missions



- Agency wide resource
- Production Supercomputing
 - Focus on availability
- Machines mostly run large ensembles
- Some very large calculations (50k)
 - Typically o500 jobs running

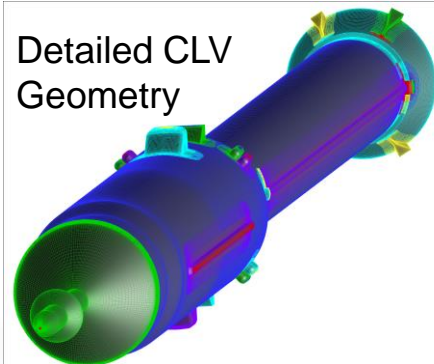
- Example applications
- ARMD
 - LaRC: Jet wake vortex simulations, to increase airport capacity and safety
 - GRC: Understanding jet noise simulations, to decrease airport noise
- ESMD
 - ARC: Launch pad flame trench simulations for Ares vehicle safety analysis
 - MSFC: Correlating wind tunnel tests and simulations of Ares I-X test vehicle
 - ARC/LaRC: High-fidelity CLV flight simulation with detailed protuberances
- SMD
 - Michigan State: Ultra-high-resolution solar surface convection simulation
 - GSFC: Gravity waves from the merger of orbiting, spinning black holes
- SOMD
 - JSC/ARC: Ultra-high-resolution Shuttle ascent analysis
- NESC
 - KSC/ARC: Initial analysis of SRB burn risk in Vehicle Assembly Building

Jet aircraft wake vortices

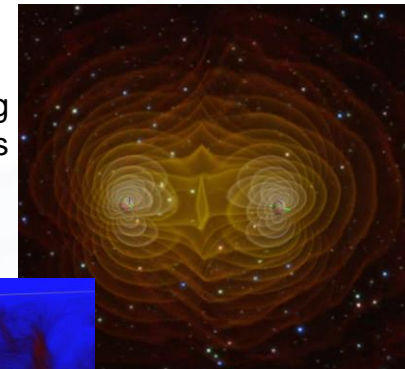


Jet engine noise emissions

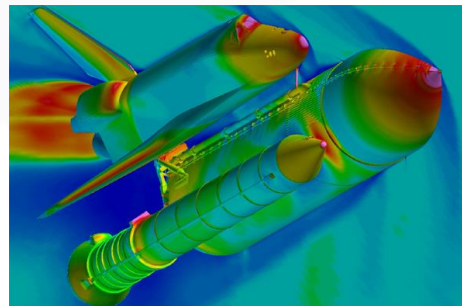
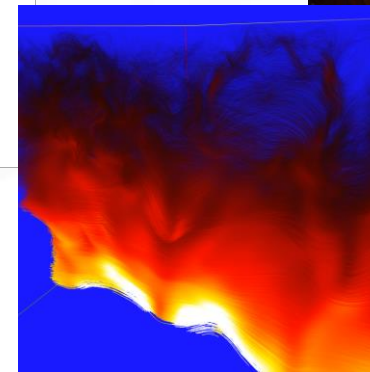
Detailed CLV Geometry



Orbiting, Spinning Black Holes

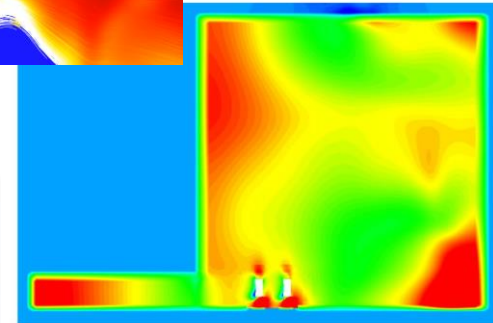


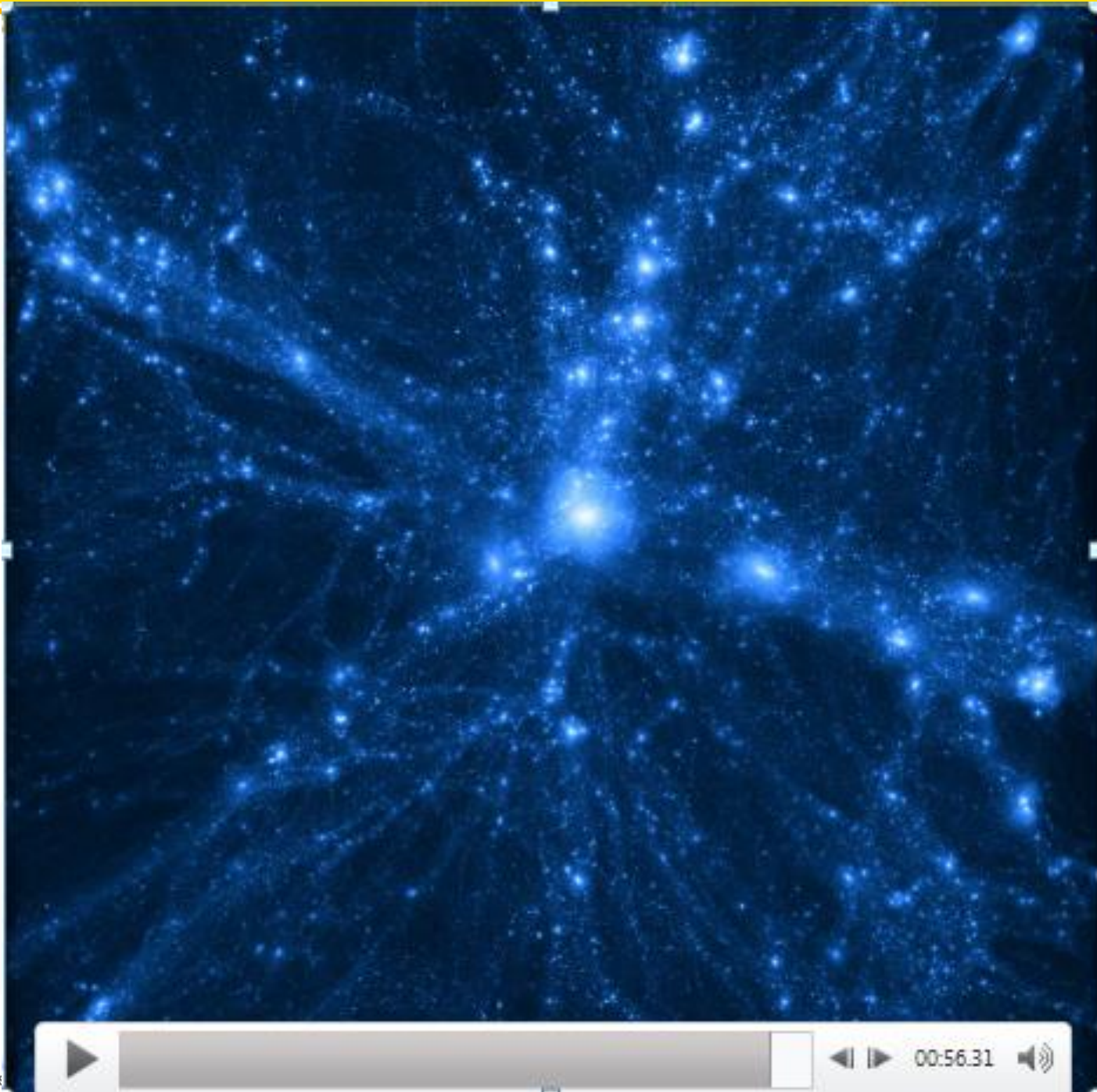
Solar surface convection

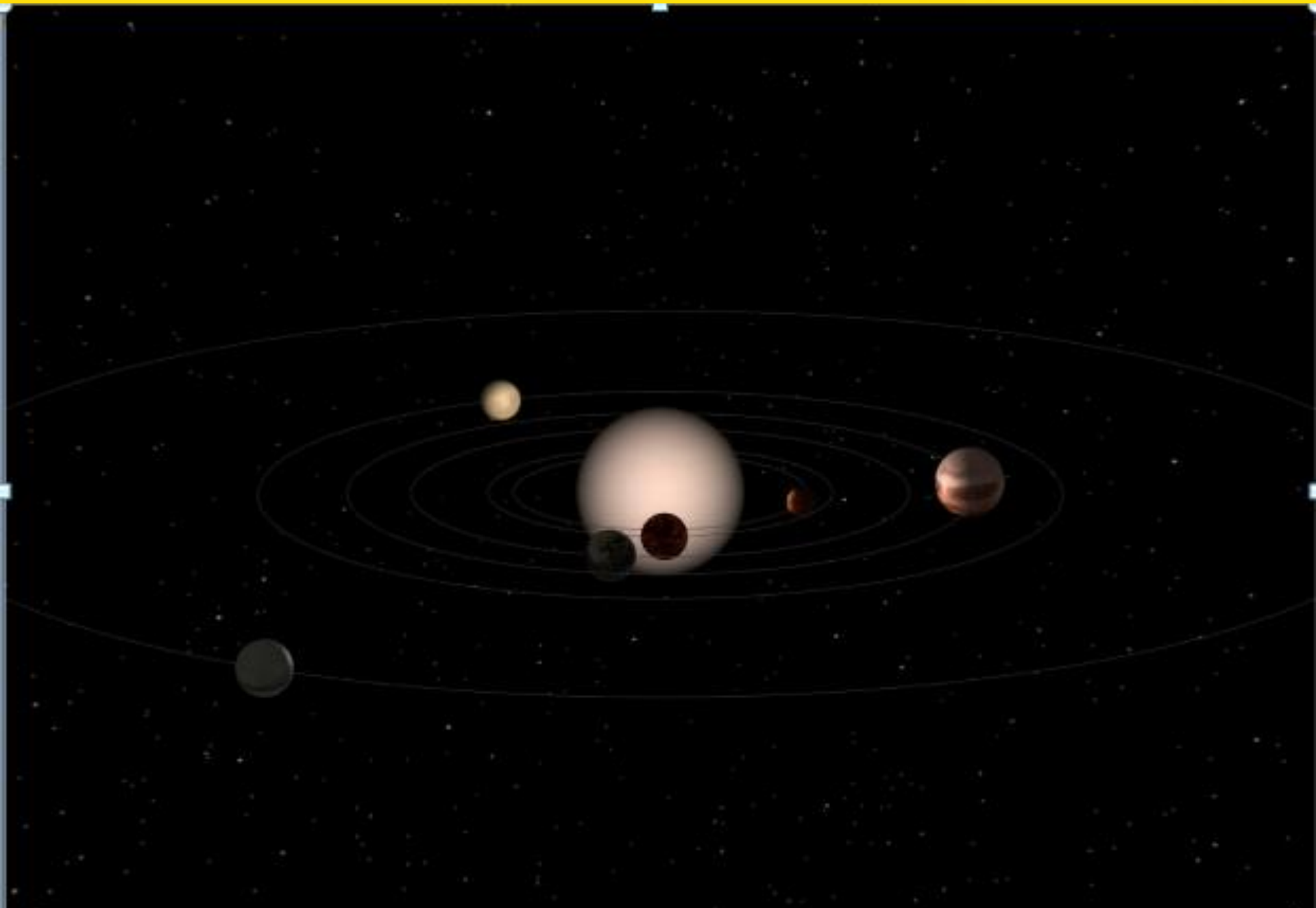


Shuttle Ascent Configuration

2-SRB Burn in VAB

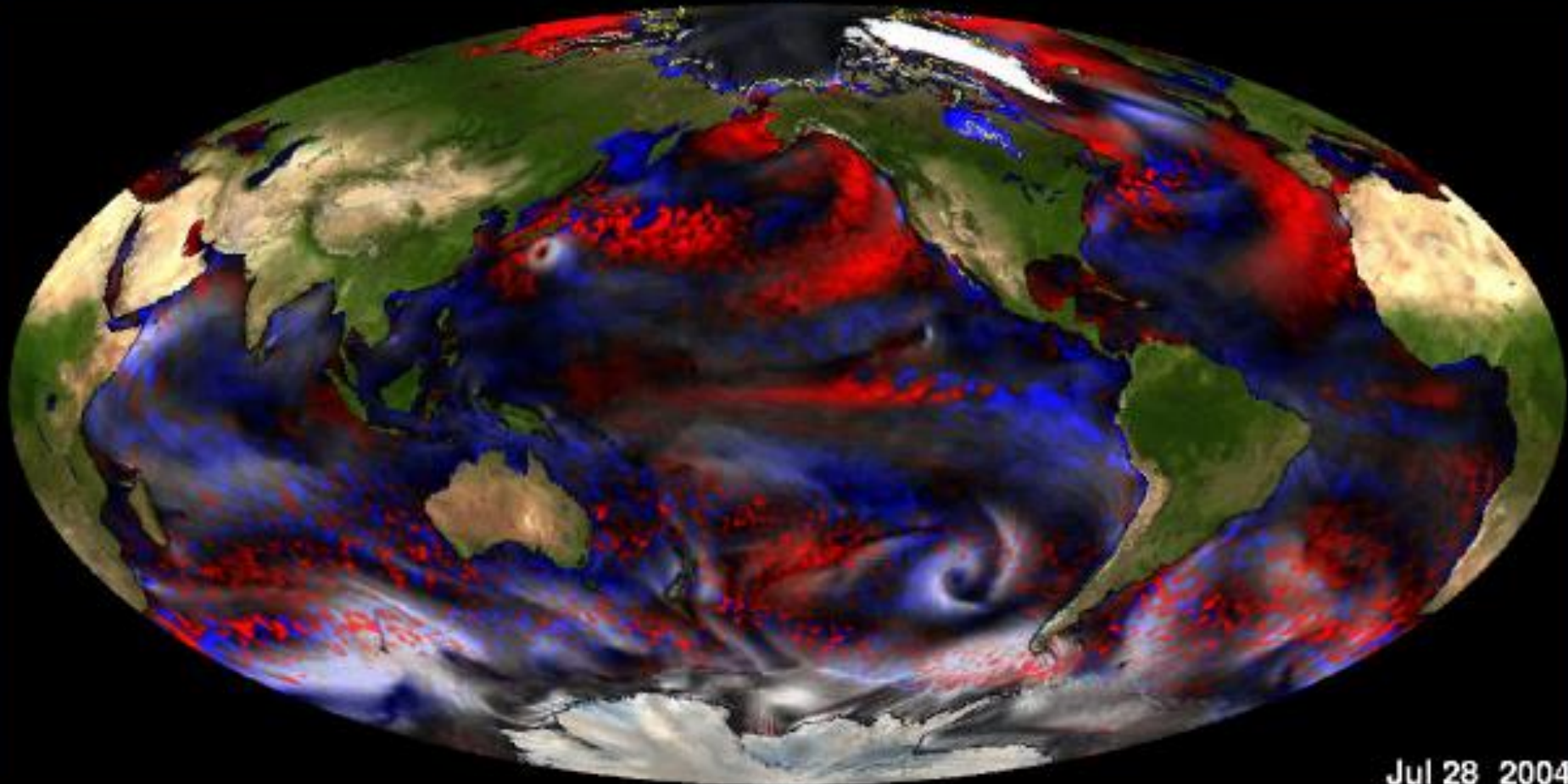






00:33.74

A video player interface. On the left is a play button. Next to it is a progress bar with a slider. On the right is a volume icon and a timestamp of 00:33.74.

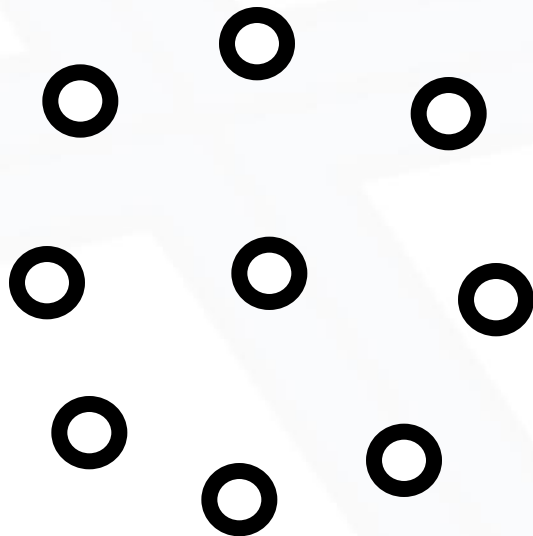


Jul 28 2004

NASA's Computational Landscape

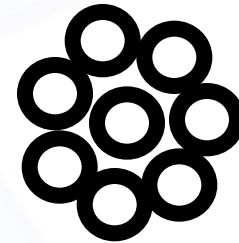
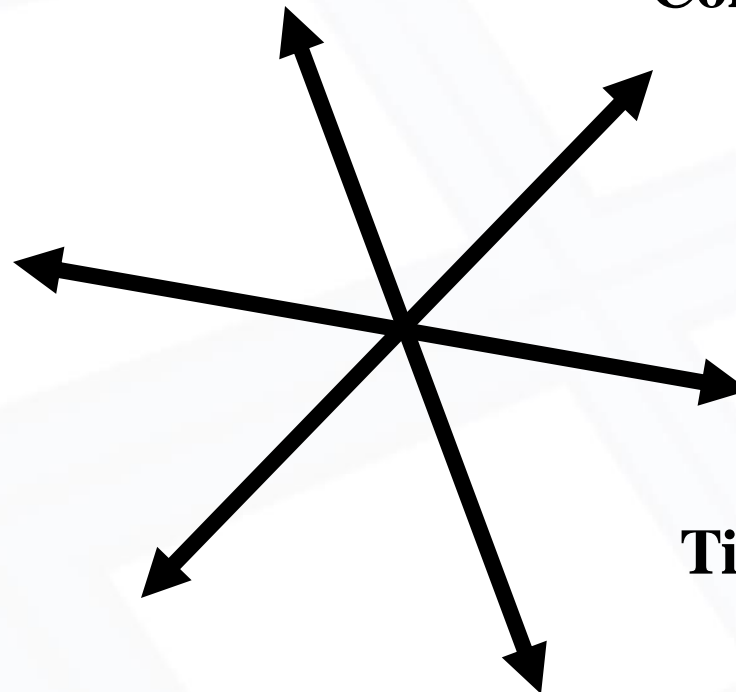


Embarrassingly Parallel



Compute Bound

Simple Well Understood Computations



Tightly Coupled

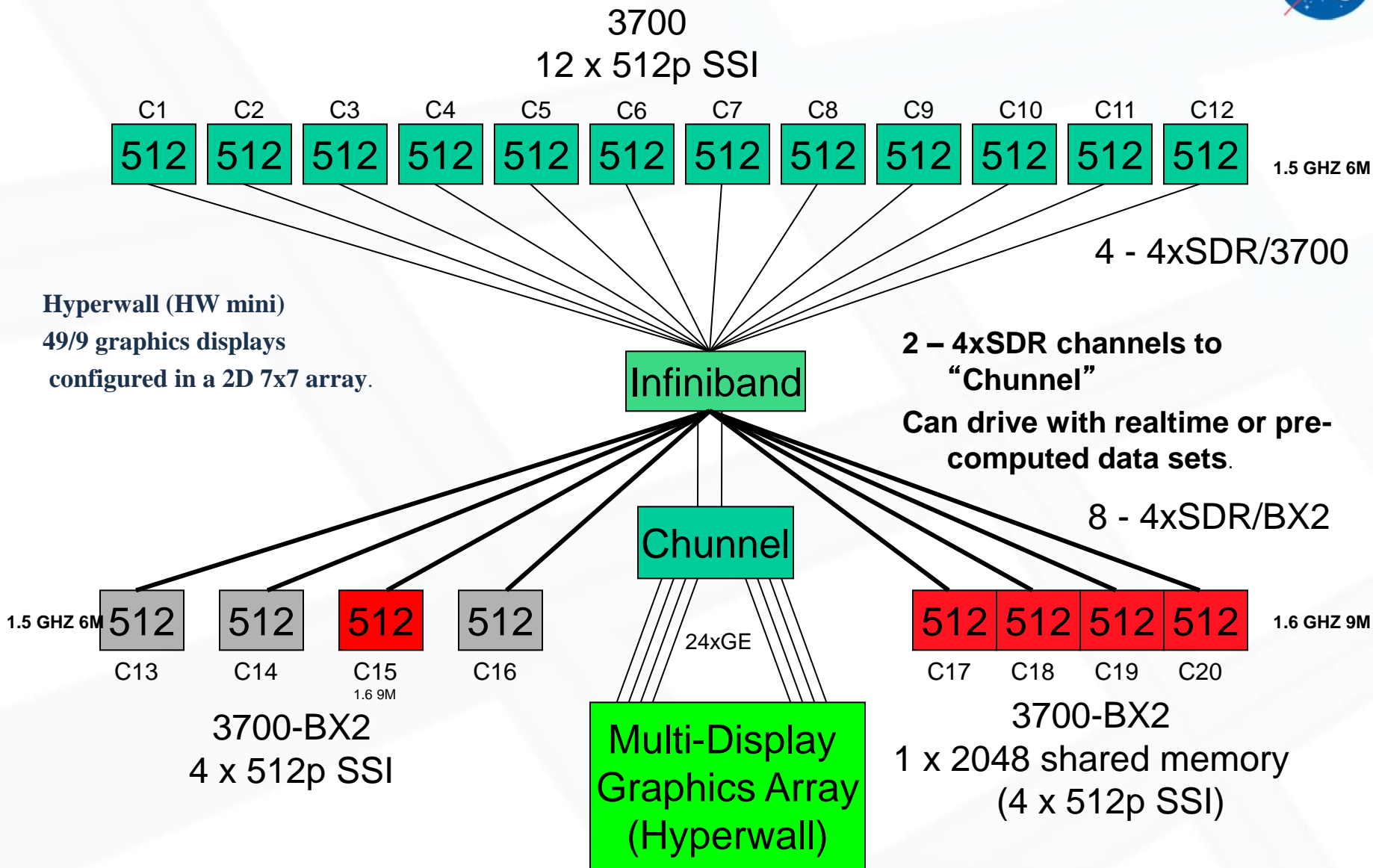
Highly Complex and Evolving Computations

Data/Storage Intensive

Columbia System - October 2004



Columbia

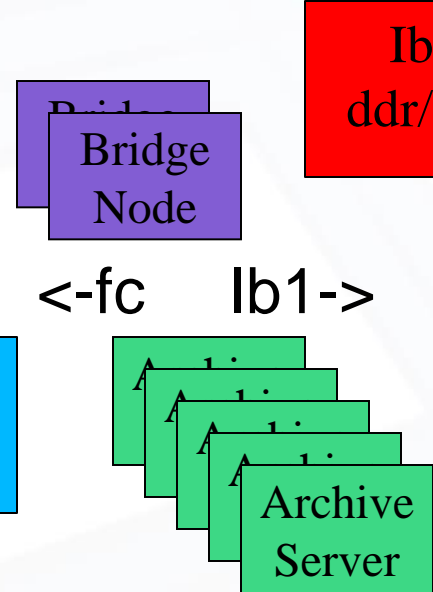
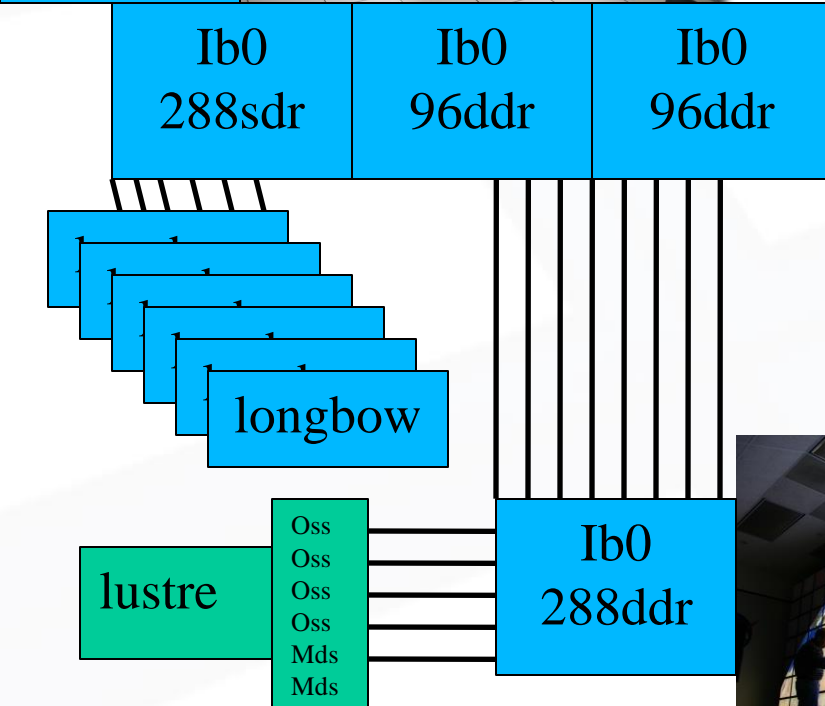




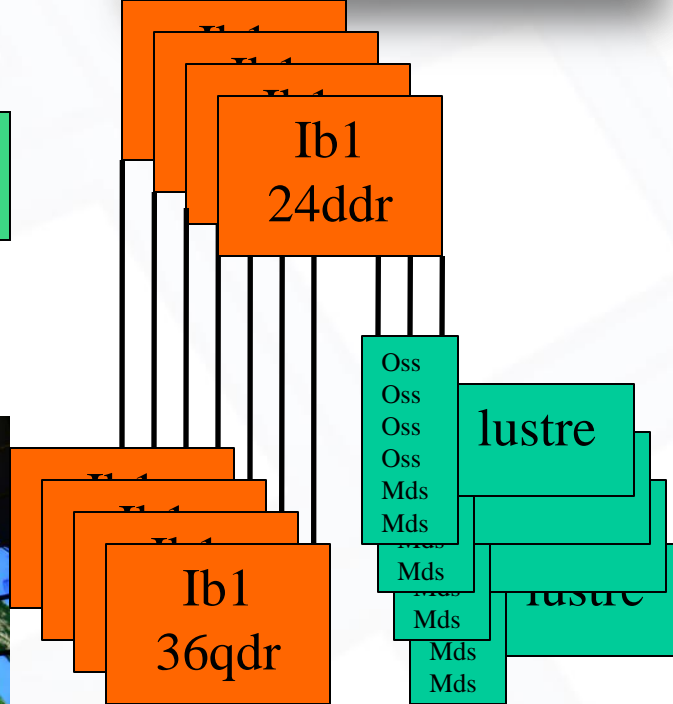
OFA/Infiniband Connectivity

Columbia

Pleiades

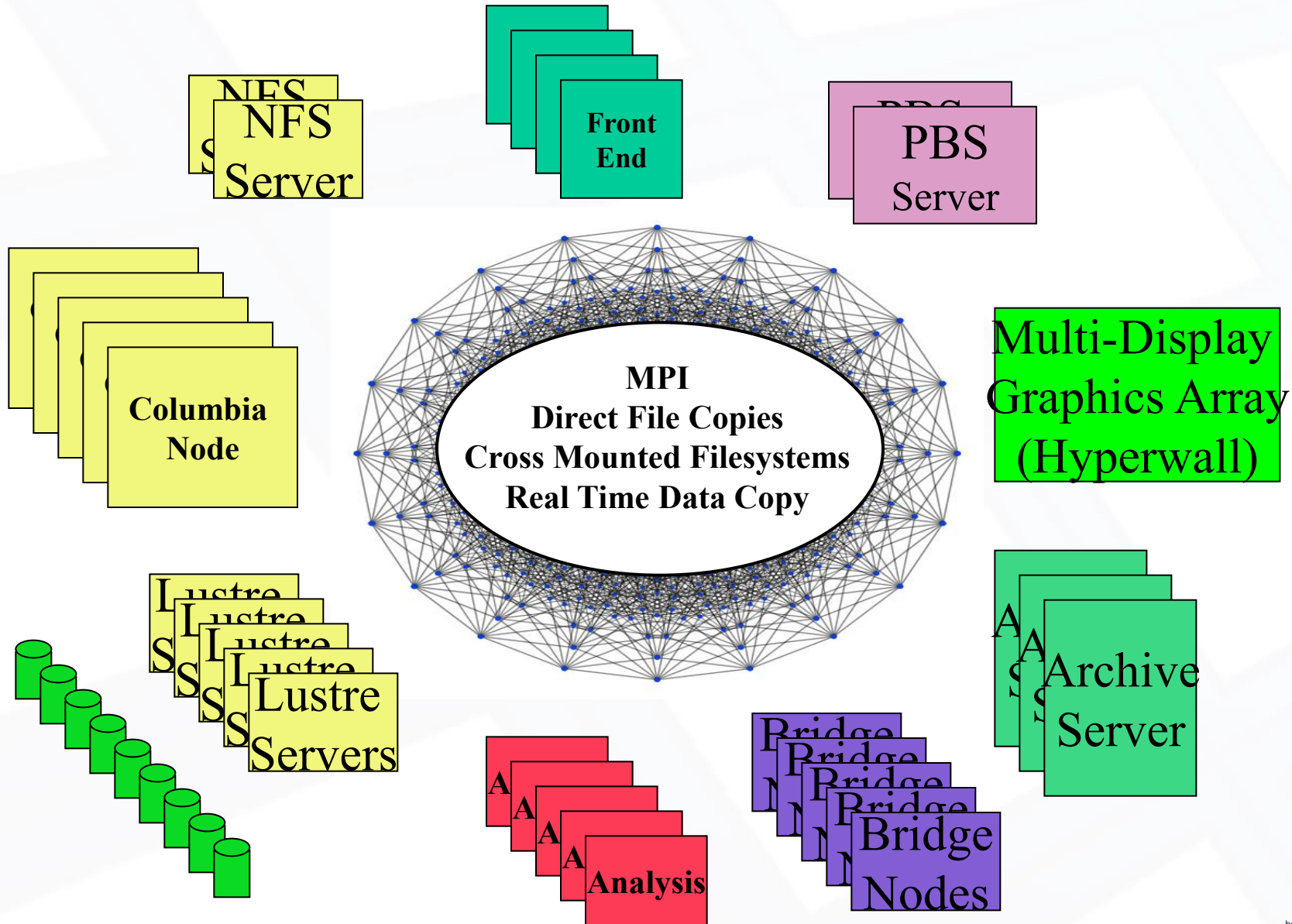


hyperwall



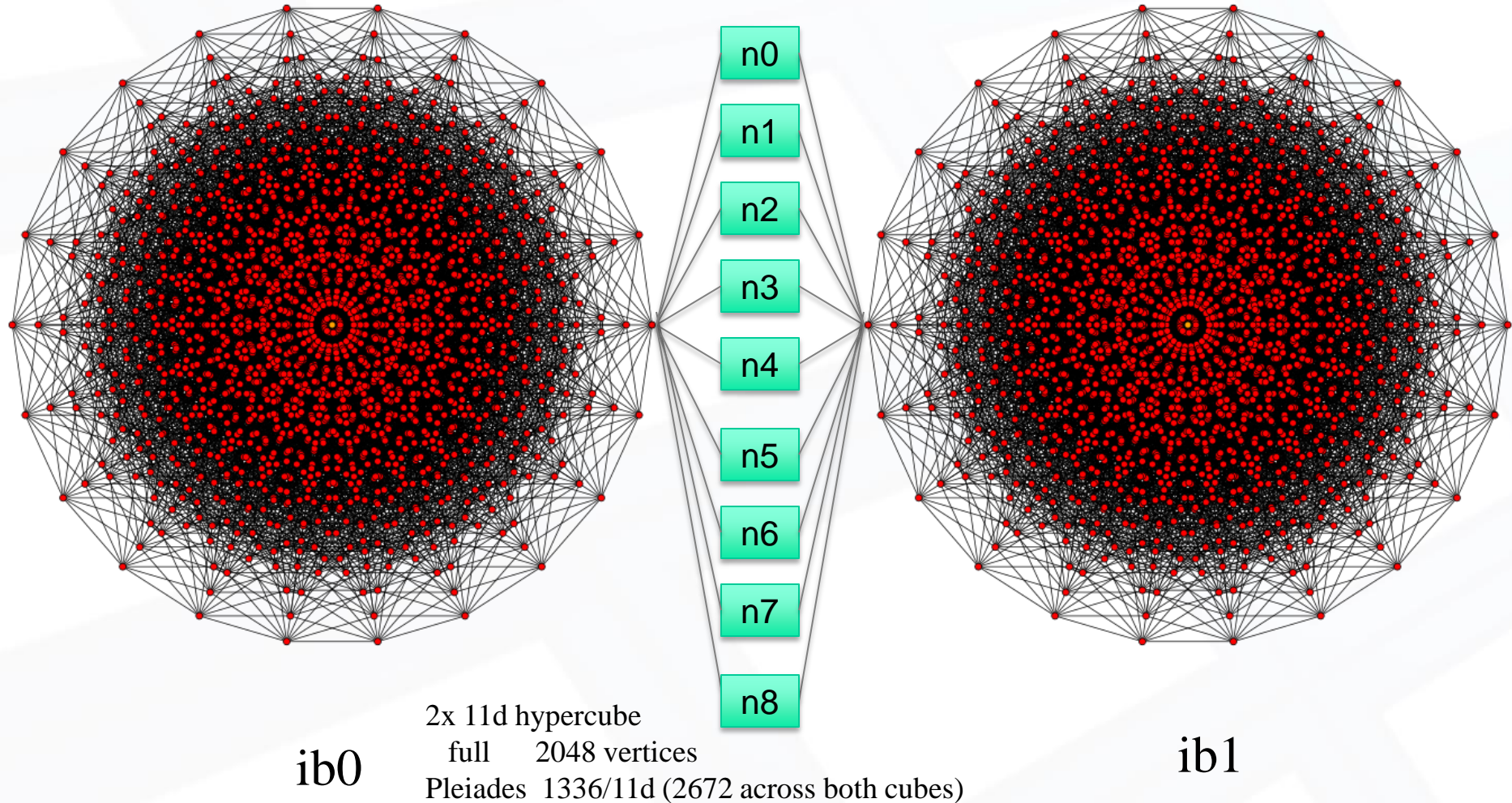


Architecture Target





SGI ICE Dual Plane – Topology

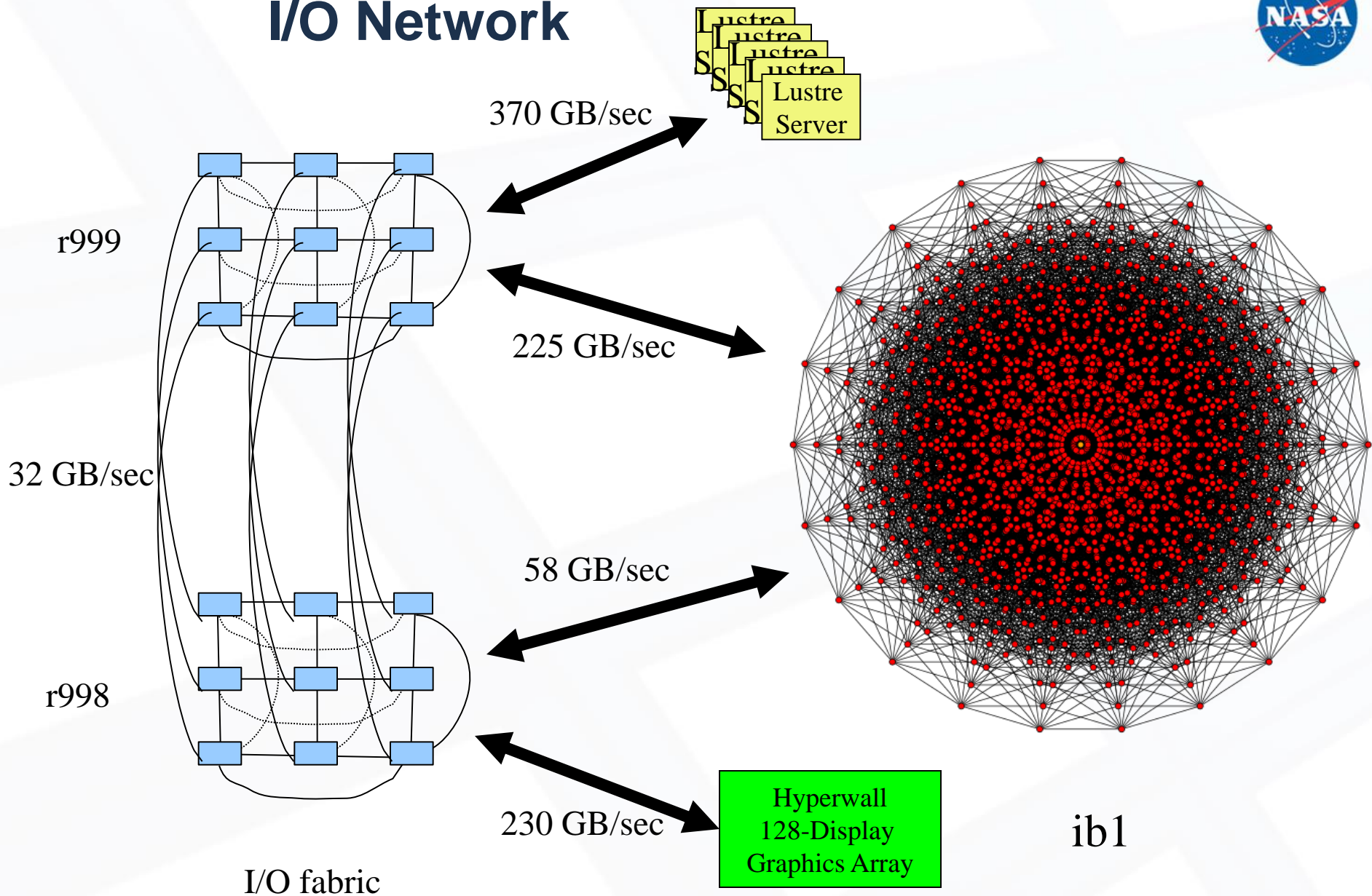


http://en.wikipedia.org/wiki/User:Qef/Orthographic_hypercube_diagram

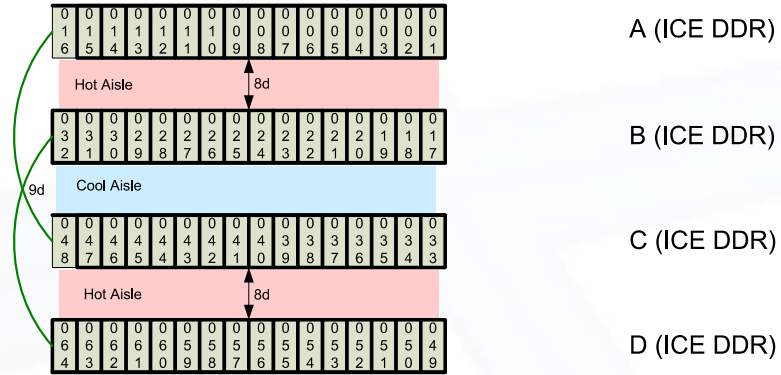
^s OFA14



I/O Network

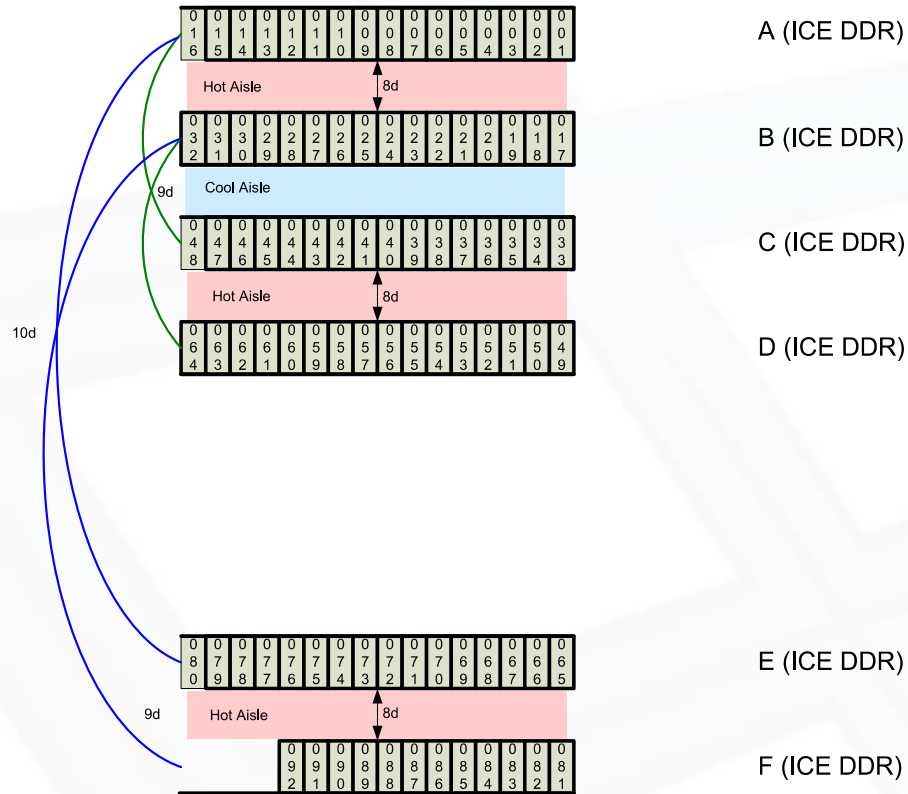


NASA (Pleiades) Rack Layout



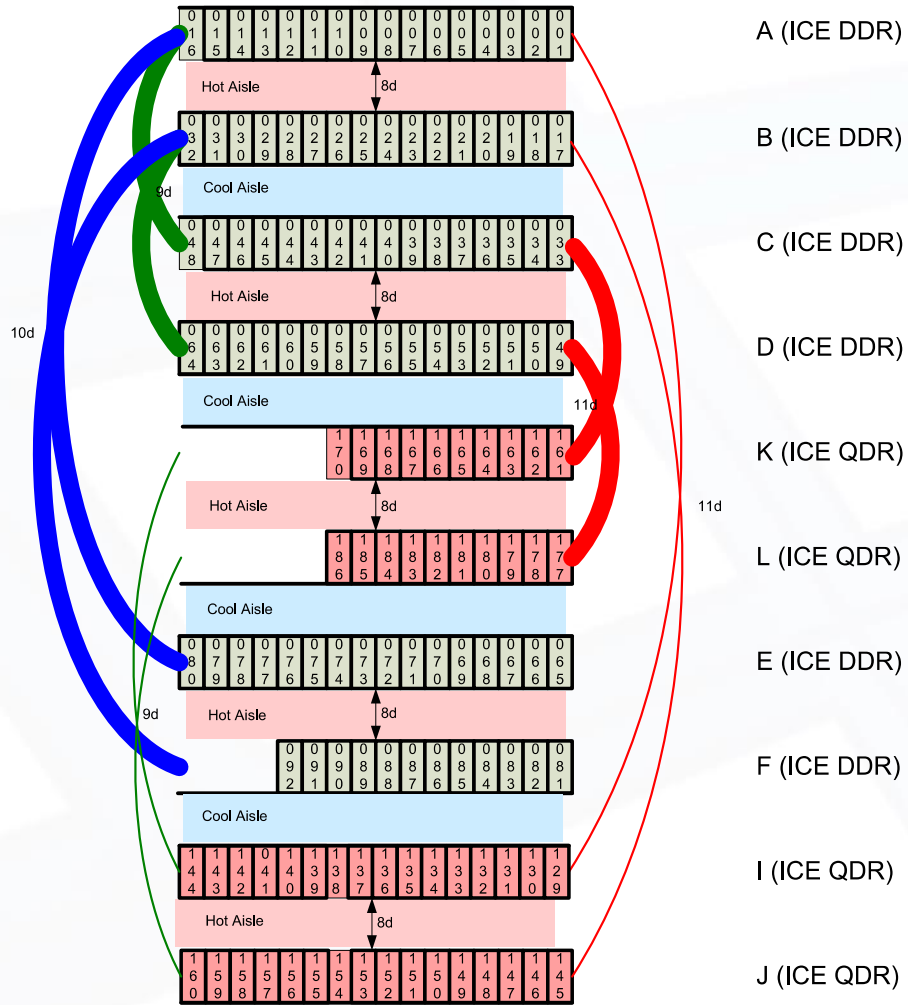
64 racks – 2008
393 teraflops

NASA (Pleiades) Rack Layout



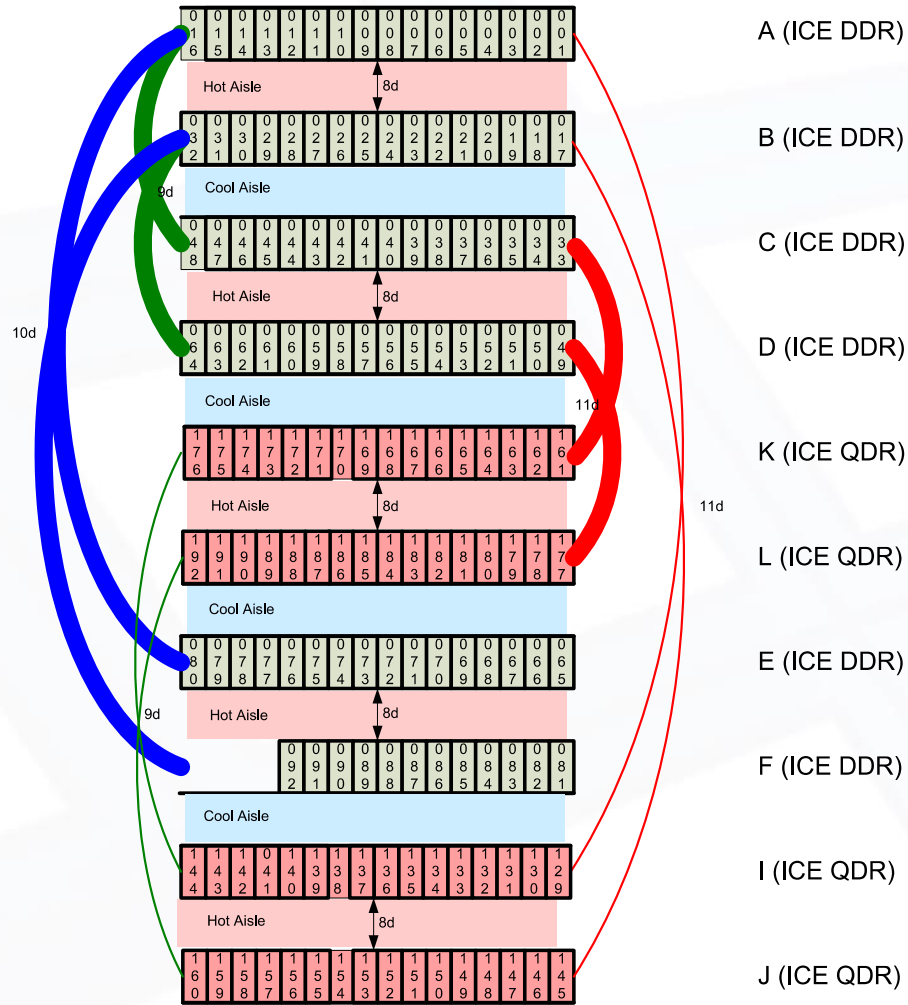
92 racks – 2008
565 teraflops

NASA (Pleiades) Rack Layout



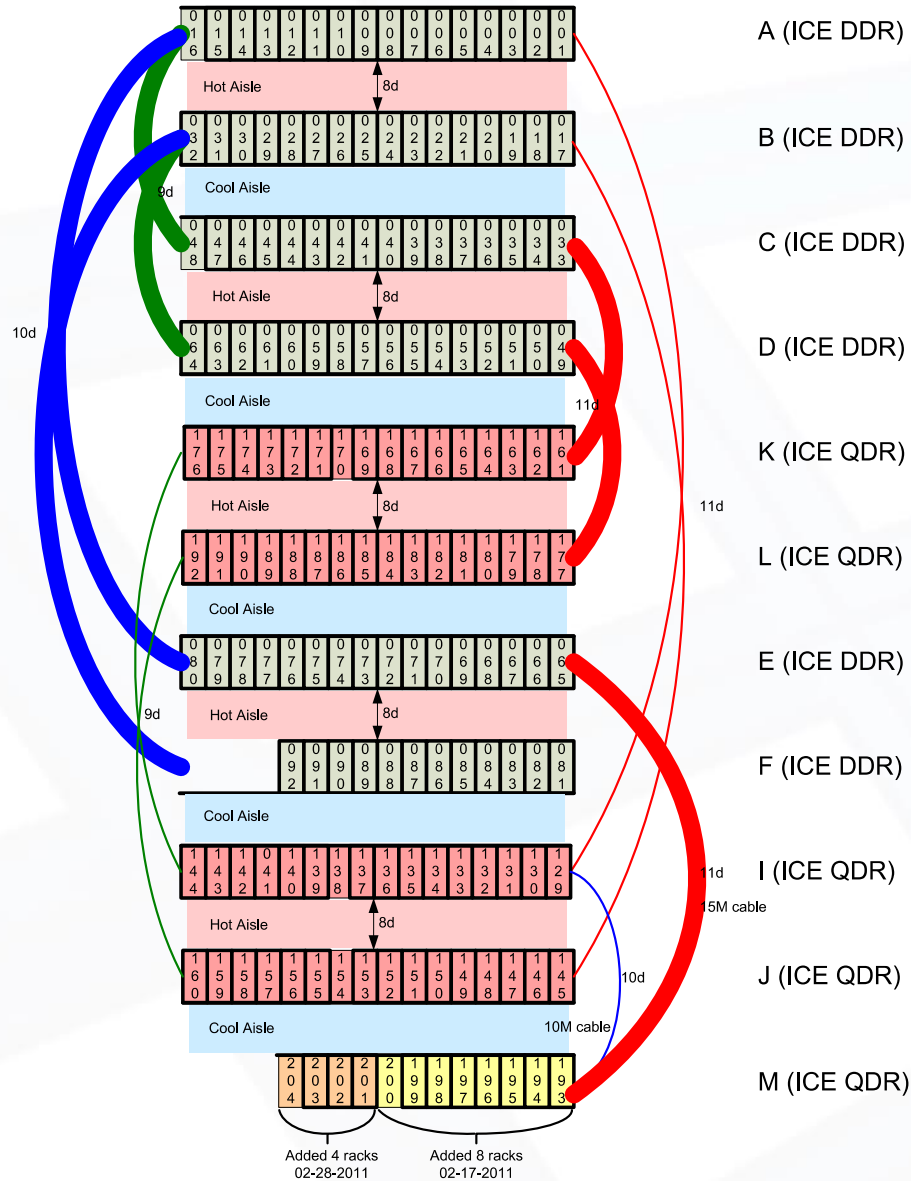
144 racks – 2010
969 teraflops

NASA (Pleiades) Rack Layout



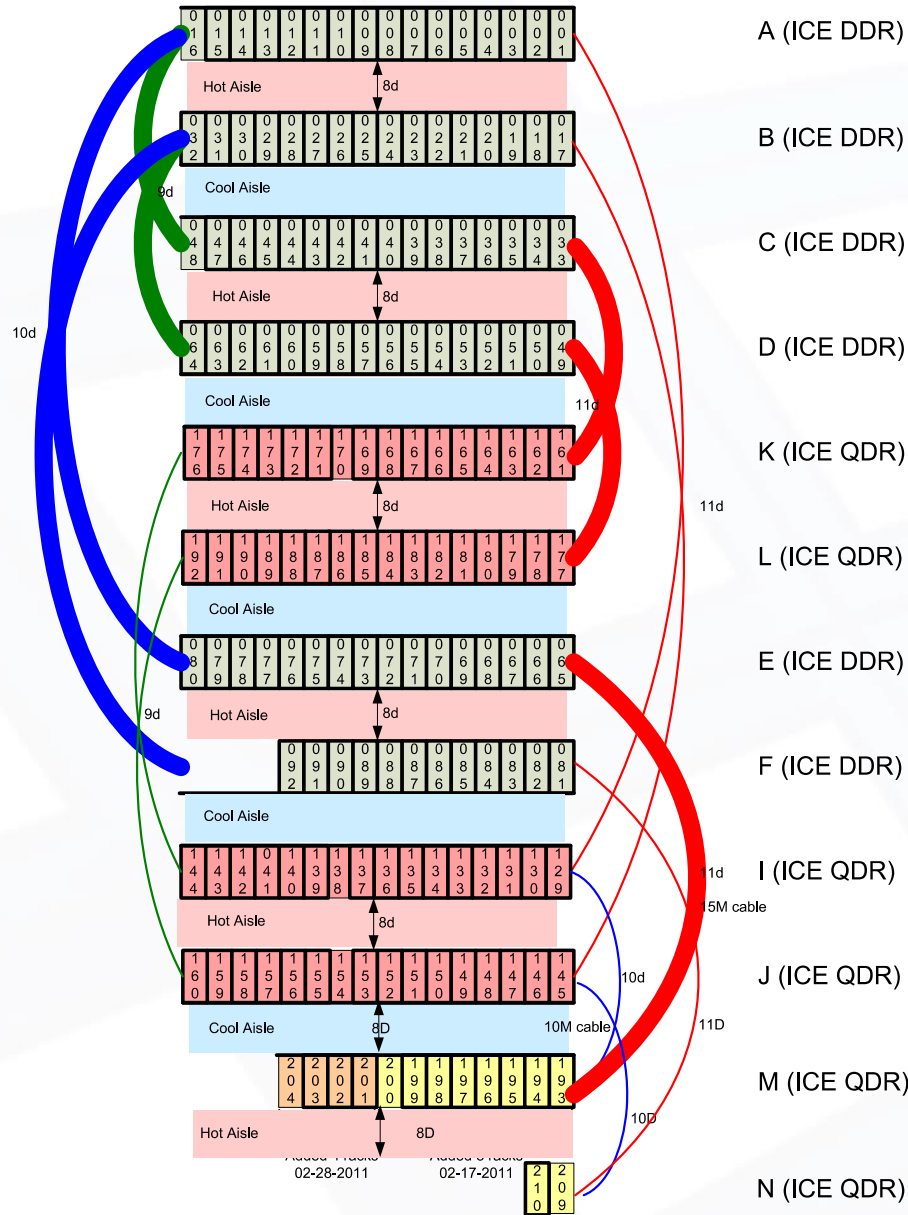
156 racks – 2010
 1.08 petaflops

NASA (Pleiades) Rack Layout



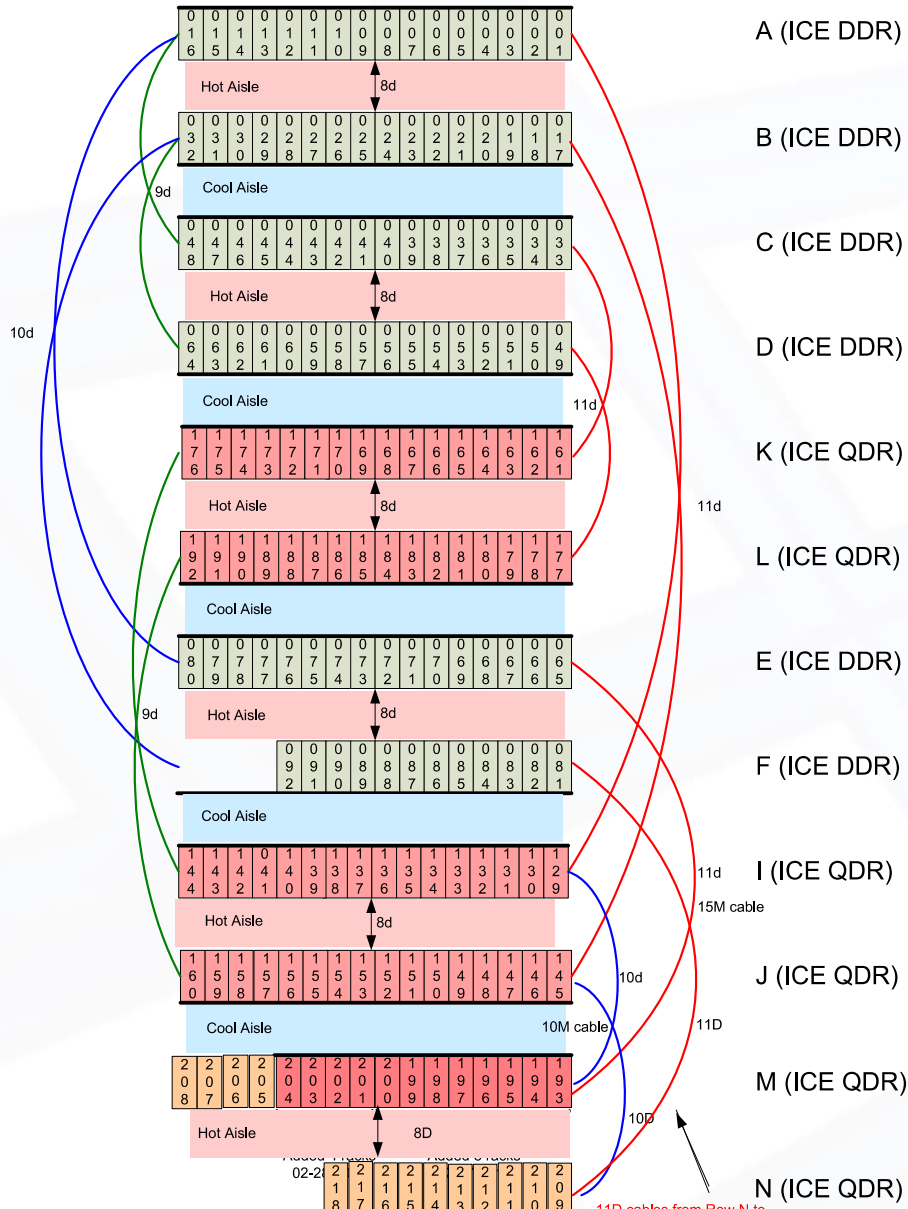
168 racks – 2011
1.18 petaflops

NASA (Pleiades) Rack Layout



170 racks – 2011
 1.20 petaflops

NASA (Pleiades) Rack Layout



182 racks – 2011
1.31 petaflops

Gpgpu racks 223 and 224

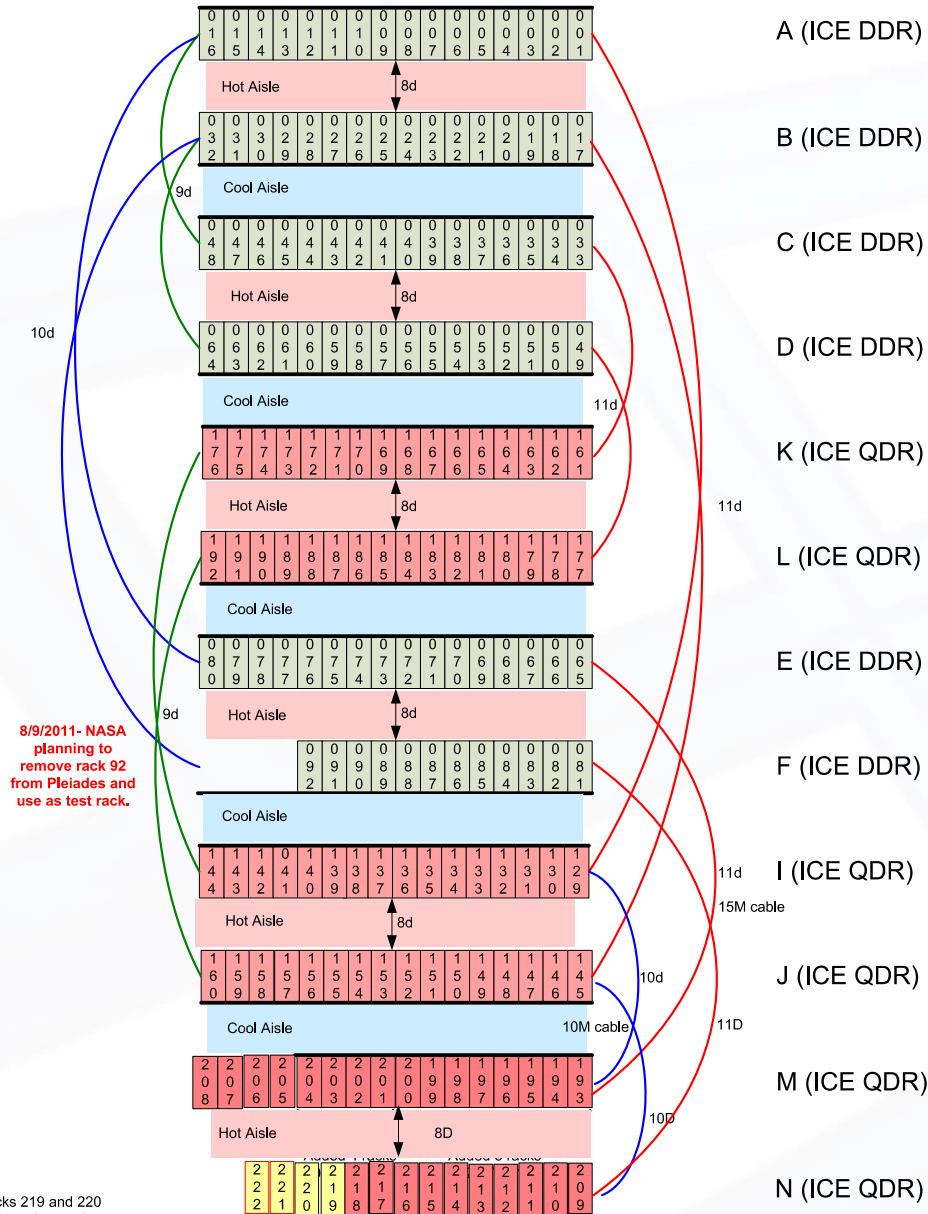
02-2
Rack 205-218 altix ICE
8400/EX westmere – new
addition to pleiades

OFA14

NASA (Pleiades) Rack Layout



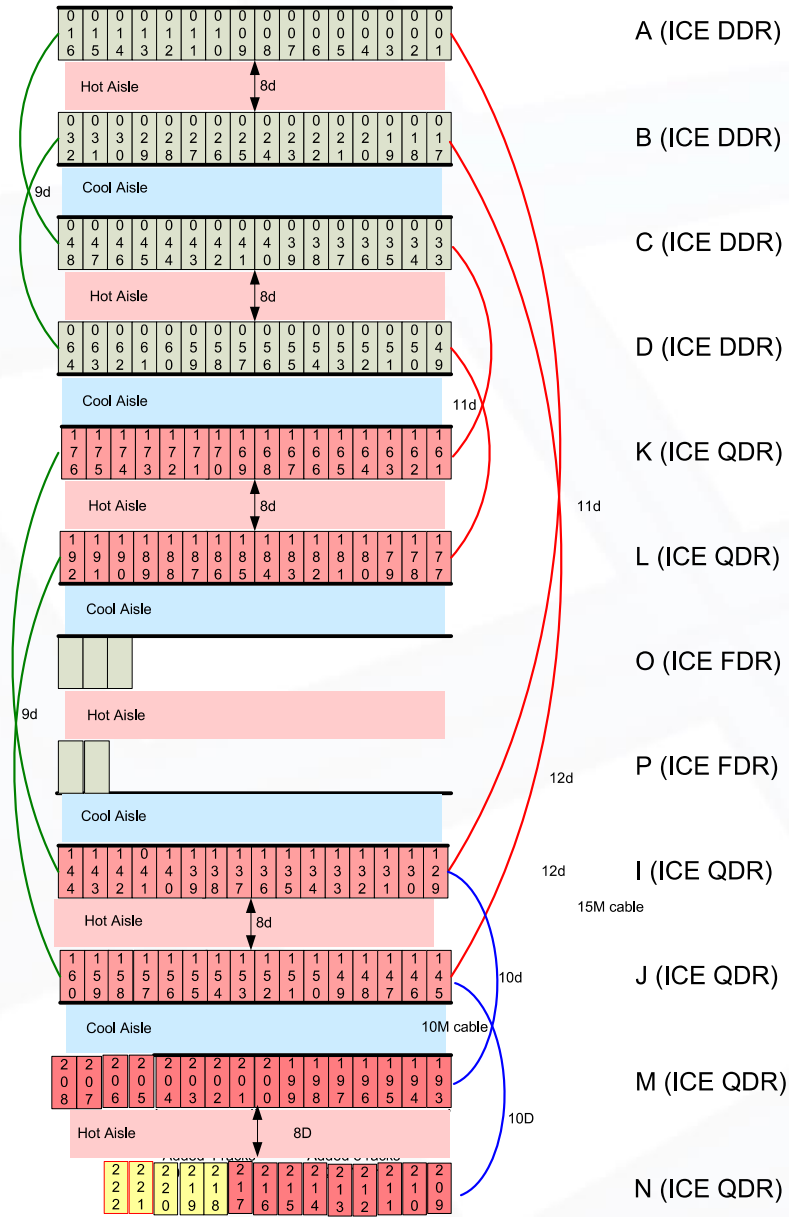
186 racks – 2011
1.33 petaflops



Gpgpu racks 219 and 220 but configured as rack 219. note switches on gpgpu are in rear of rack so cable lengths need to be adjusted to reflect this.

Note: Rack 221 will cable to on 11D to rack 92. There is no 11d for Rack 222. this is a problem. If we remove rack 92 then we have issue with racks 221 & 222.

NASA (Pleiades) Rack Layout

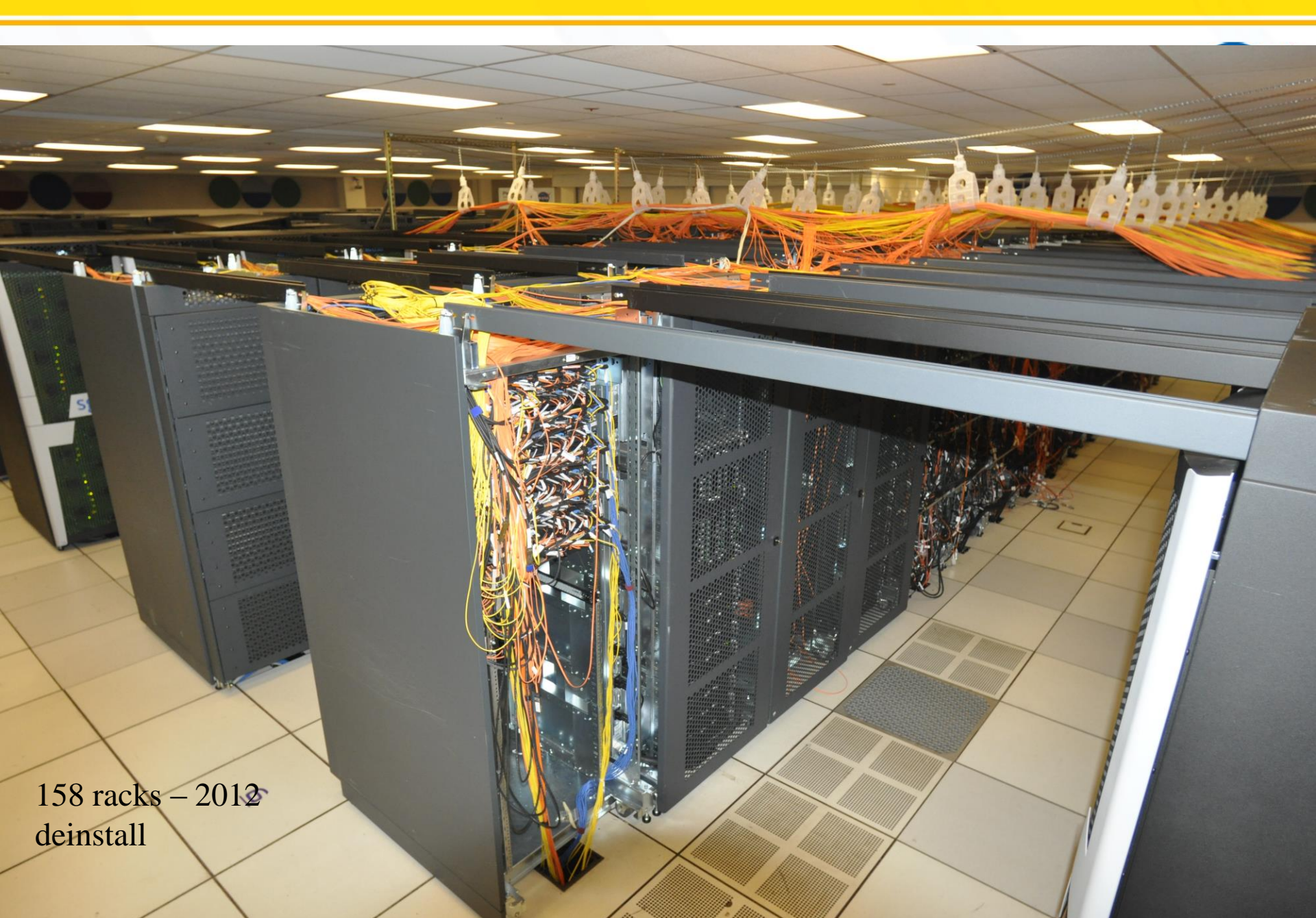


158 racks – 2012
1.15 petaflops
deinstall

*Note: Harpertown Racks
Removed 3/21/2012 in
preparation for SGI ICE X
Racks installation, I/O
Racks remain

Gpgpu racks 219 and 220
but configured as rack
219. note switches on
gpgpu are in rear of rack
so cable lengths needs to
be adjusted to reflect this.

Note: Rack 221 will cable to on 11D to rack 92. There
is no 11d for Rack 222. this is a problem. If we
remove rack 92 then we have issue with racks 221 &
222.



158 racks – 2012
deinstall

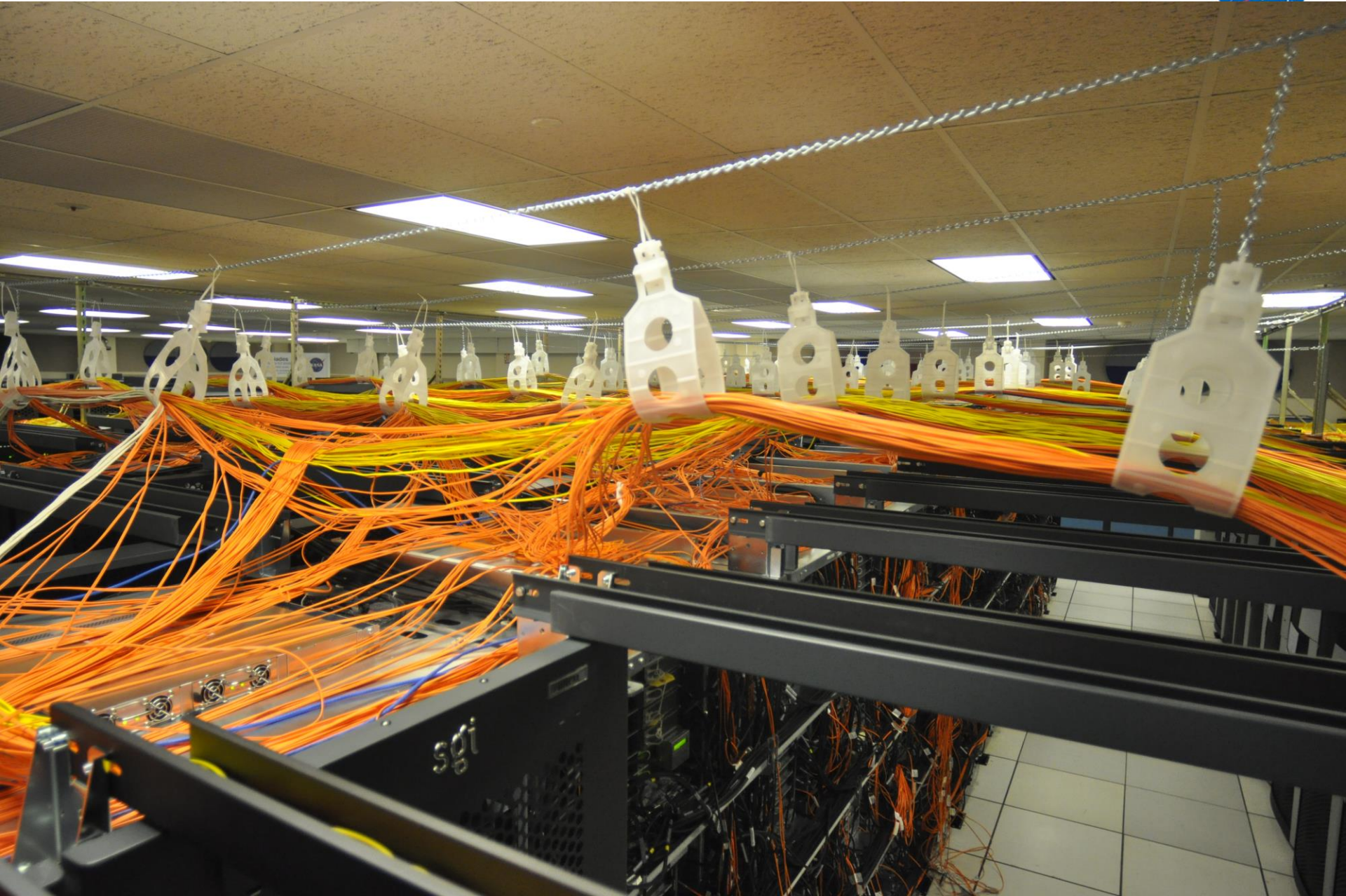
158 racks – 2012
deinstall



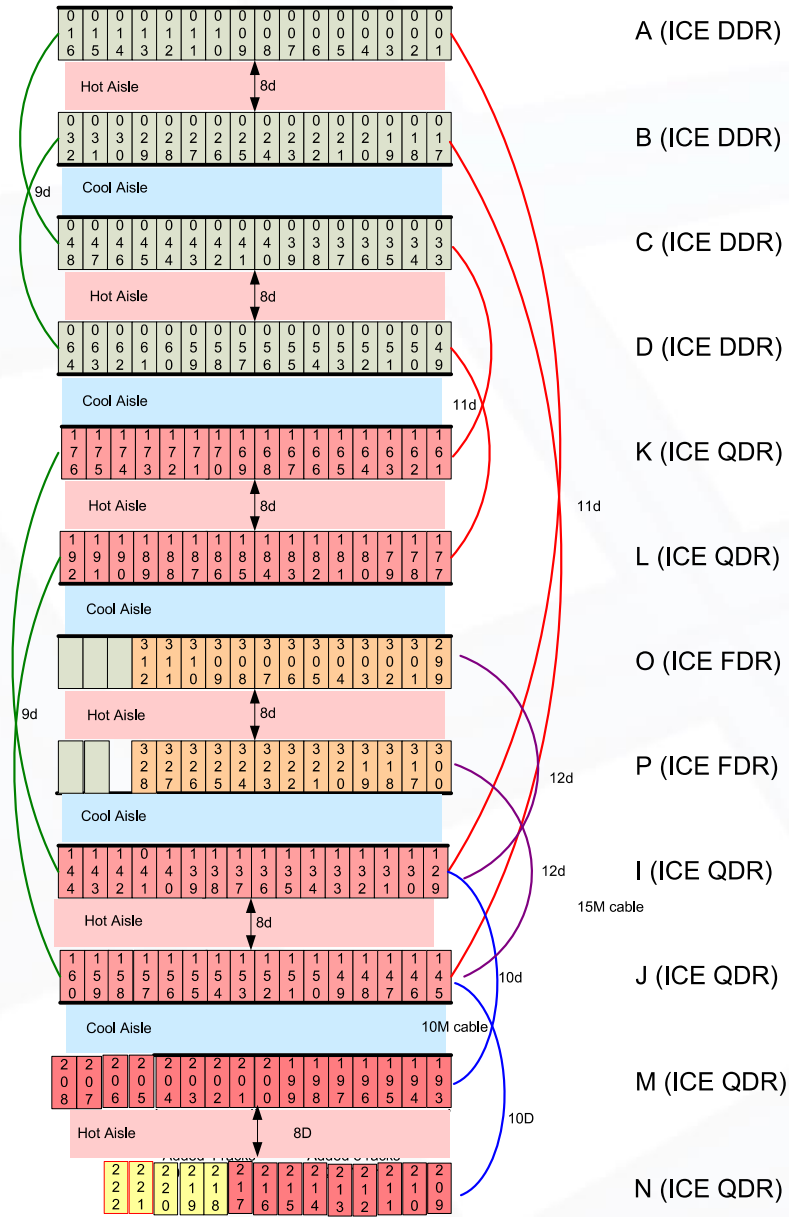




158 racks – 2012
deinstall



NASA (Pleiades) Rack Layout



182 racks – 2012
1.7 petaflops

* Install – 3/30/2012 Note:
RK 299 and RK 300 are
RLC racks. Racks 301-312
and Racks 317-328 are
Intel E5 Processors

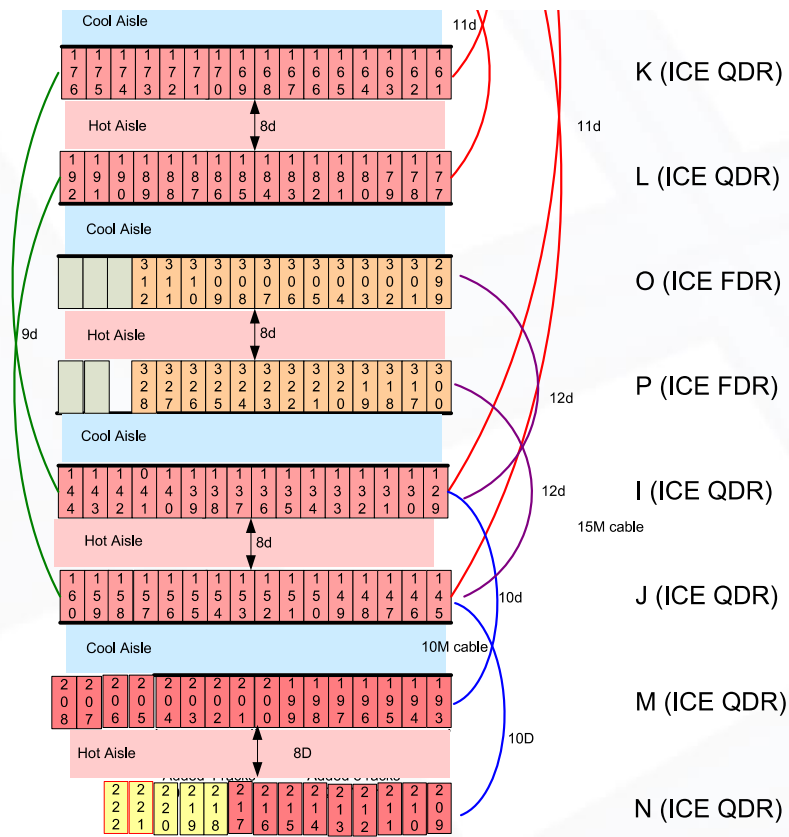
Gpgpu racks 219 and 220
but configured as rack
219. note switches on
gpgpu are in rear of rack
so cable lengths needs to
be adjusted to reflect this.

Note: Rack 221 will cable to on 11D to rack 92. There
is no 11d for Rack 222. this is a problem. If we
remove rack 92 then we have issue with racks 221 &
222.



64 rack deinstall 2013

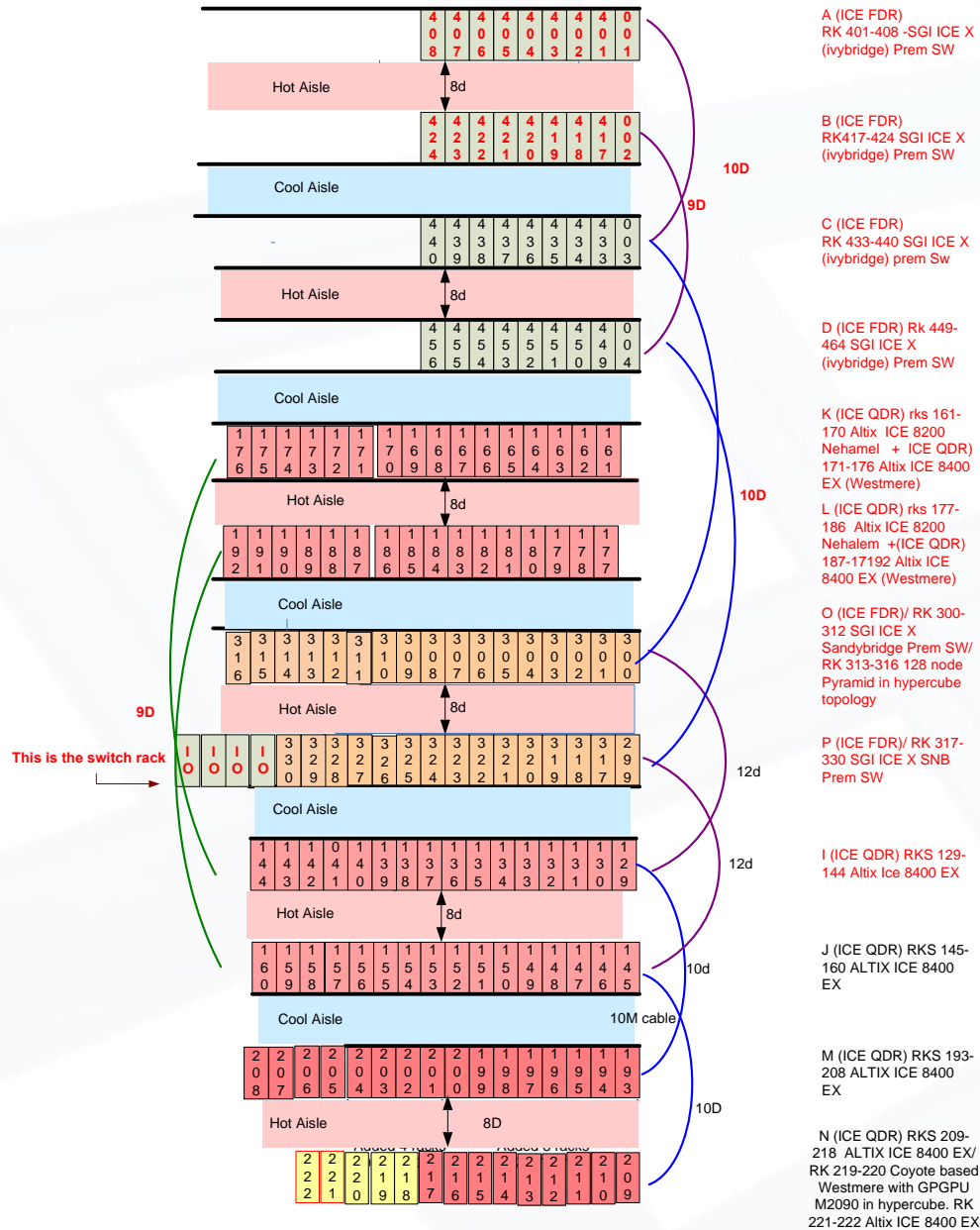
* Install - 3/30/2012 Note:
RK 299 and RK 300 are
RLC racks. Racks 301-312
and Racks 317-328 are
Intel E5 Processors



Gpgpu racks 219 and 220
but configured as rack
219. note switches on
gpgpu are in rear of rack
so cable lengths needs to
be adjusted to reflect this.

Note: Rack 221 will cable to on 11D to rack 92. There
is no 11d for Rack 222. this is a problem. If we
remove rack 92 then we have issue with racks 221 &
222.

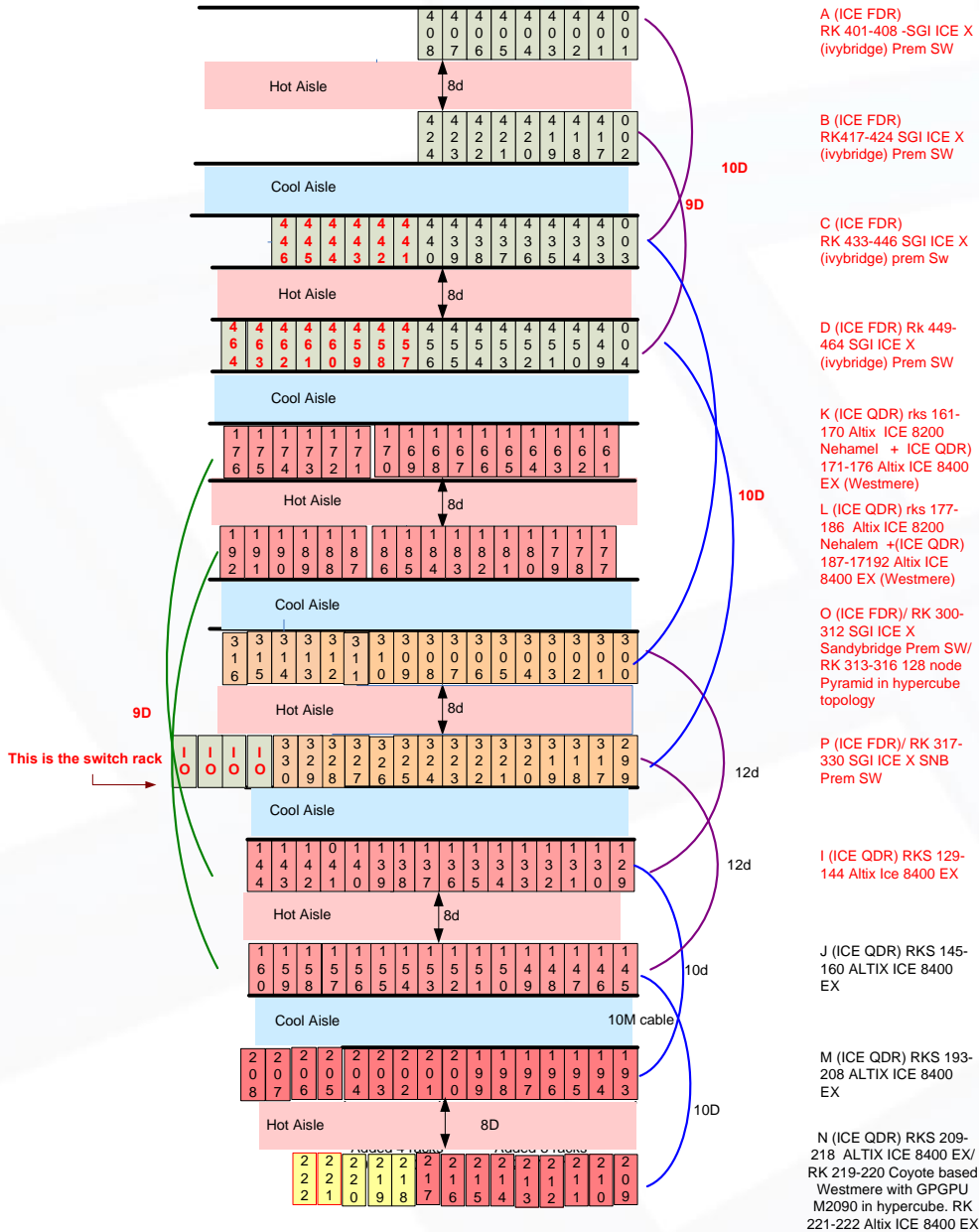
NASA (Pleiades) Rack Layout as of 08/05/2013



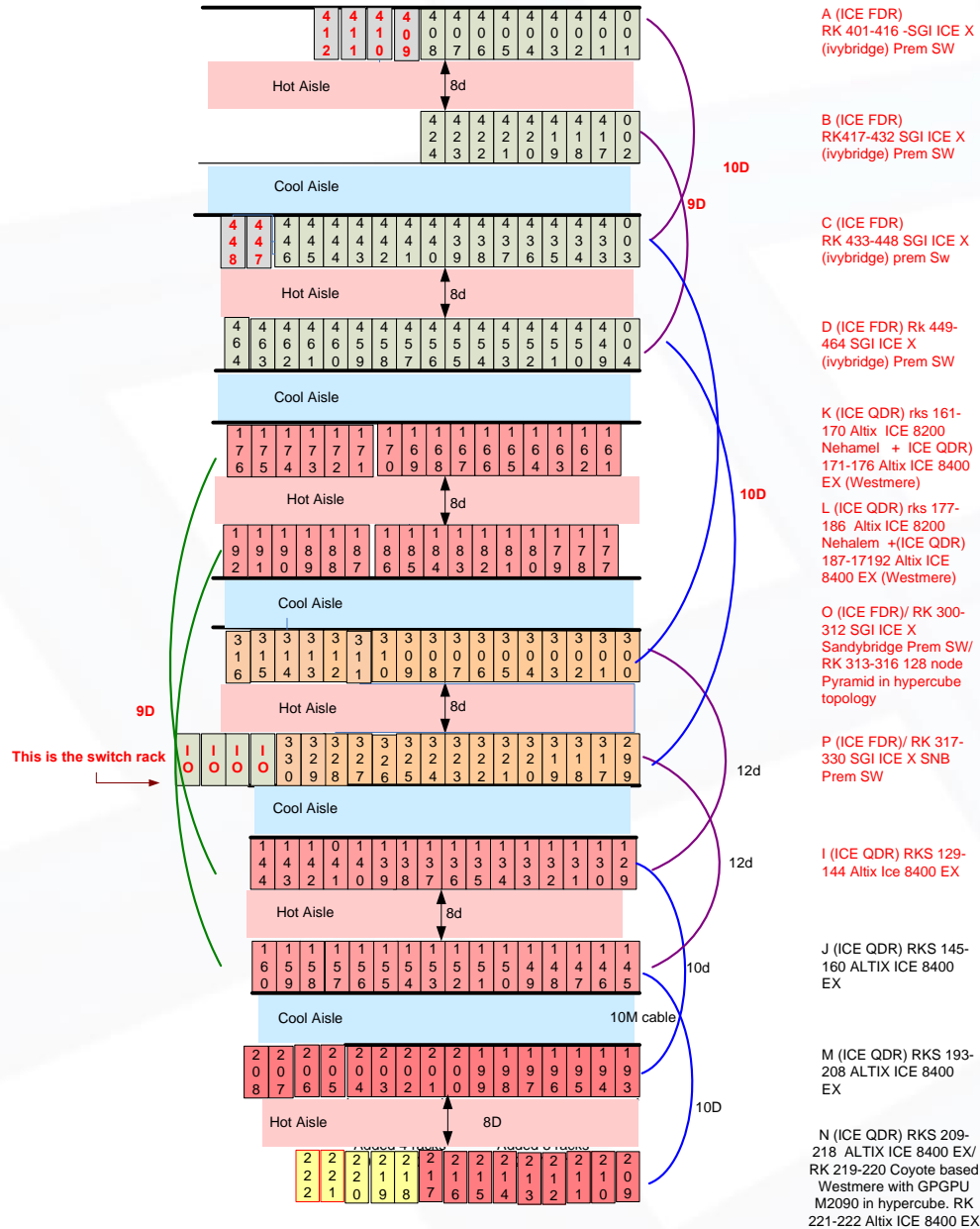
NASA (Pleiades) Rack Layout as of 08/12/2013



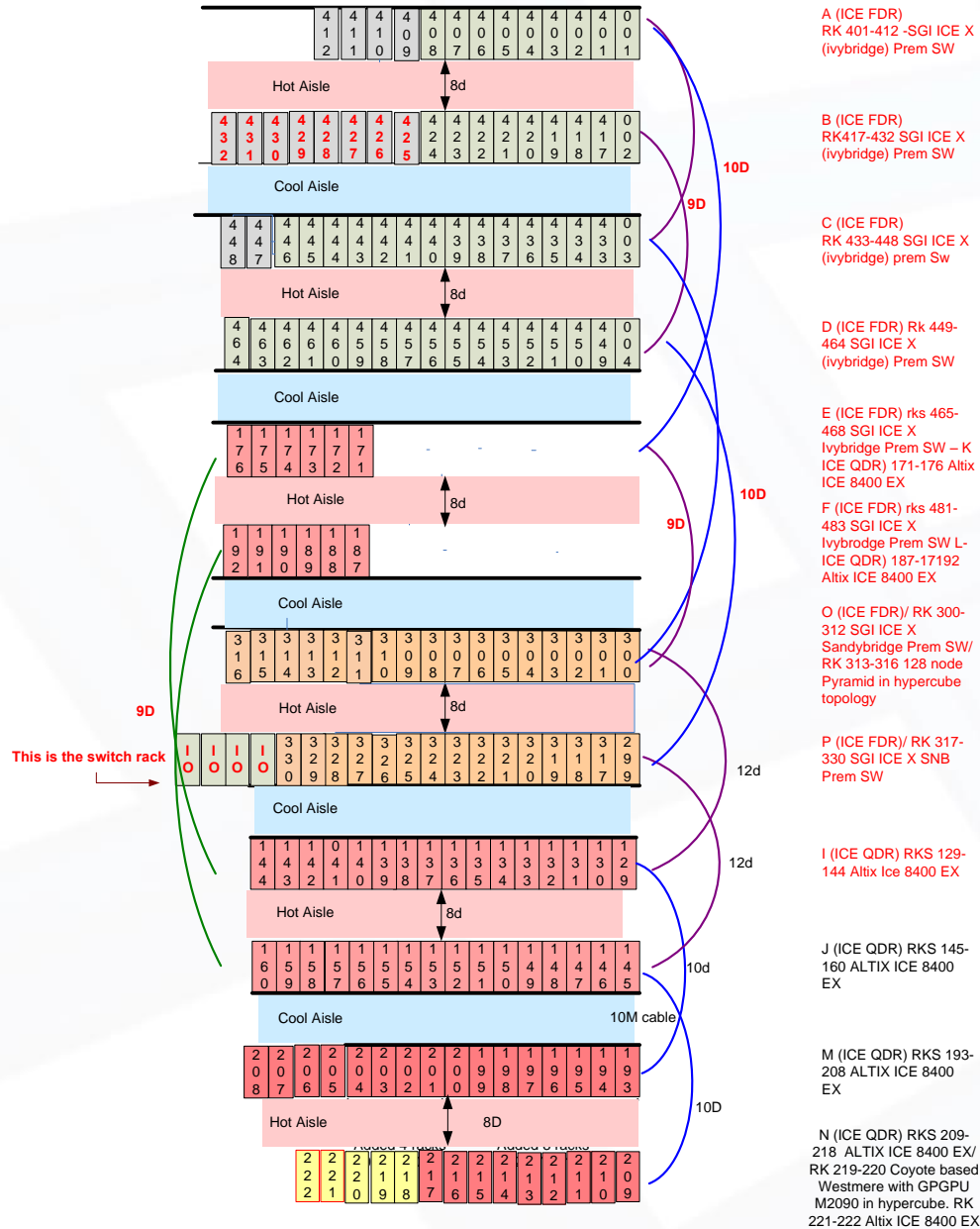
167 racks – 2013
2.9 petaflops



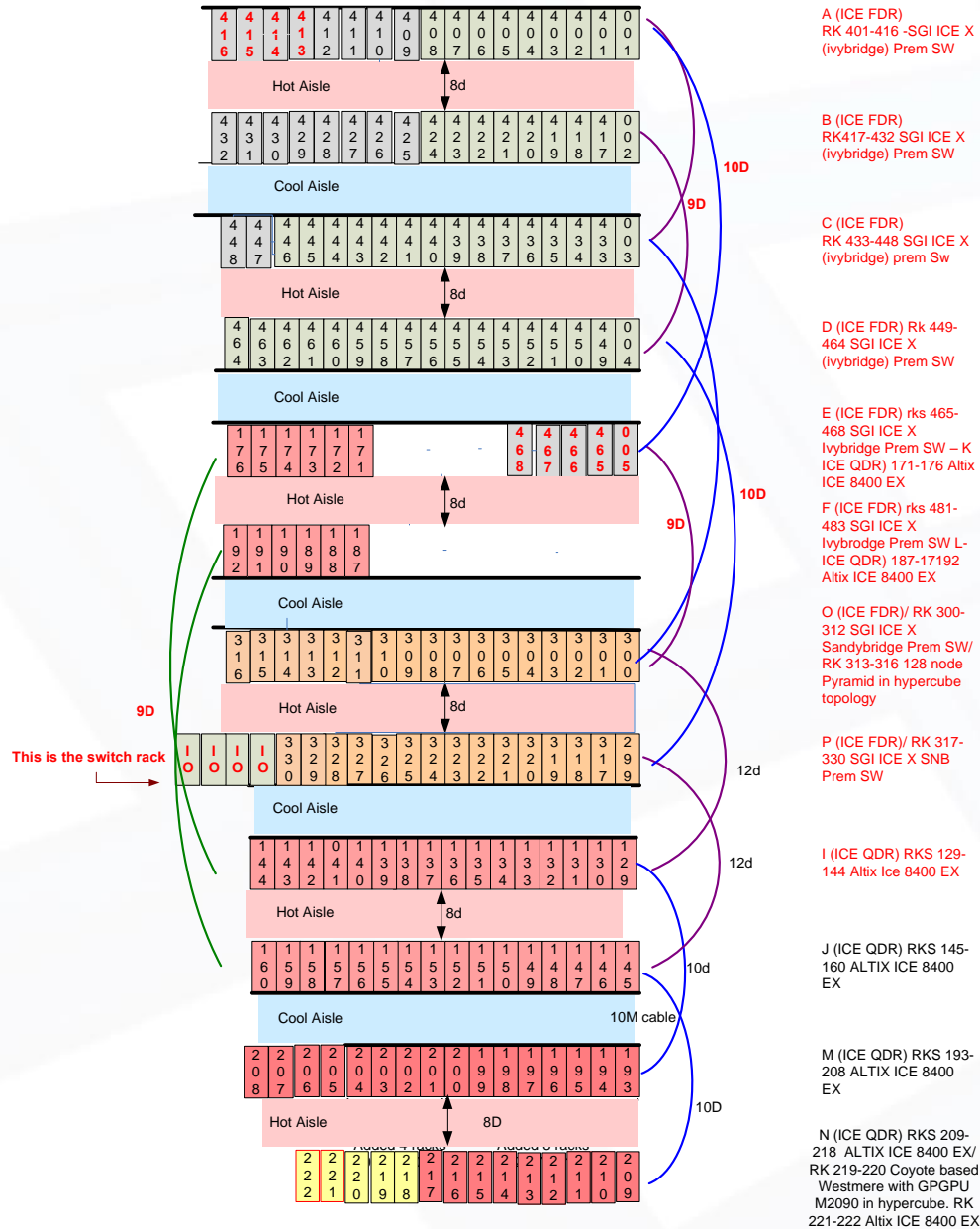
NASA (Pleiades) Rack Layout as of 12/30/2013



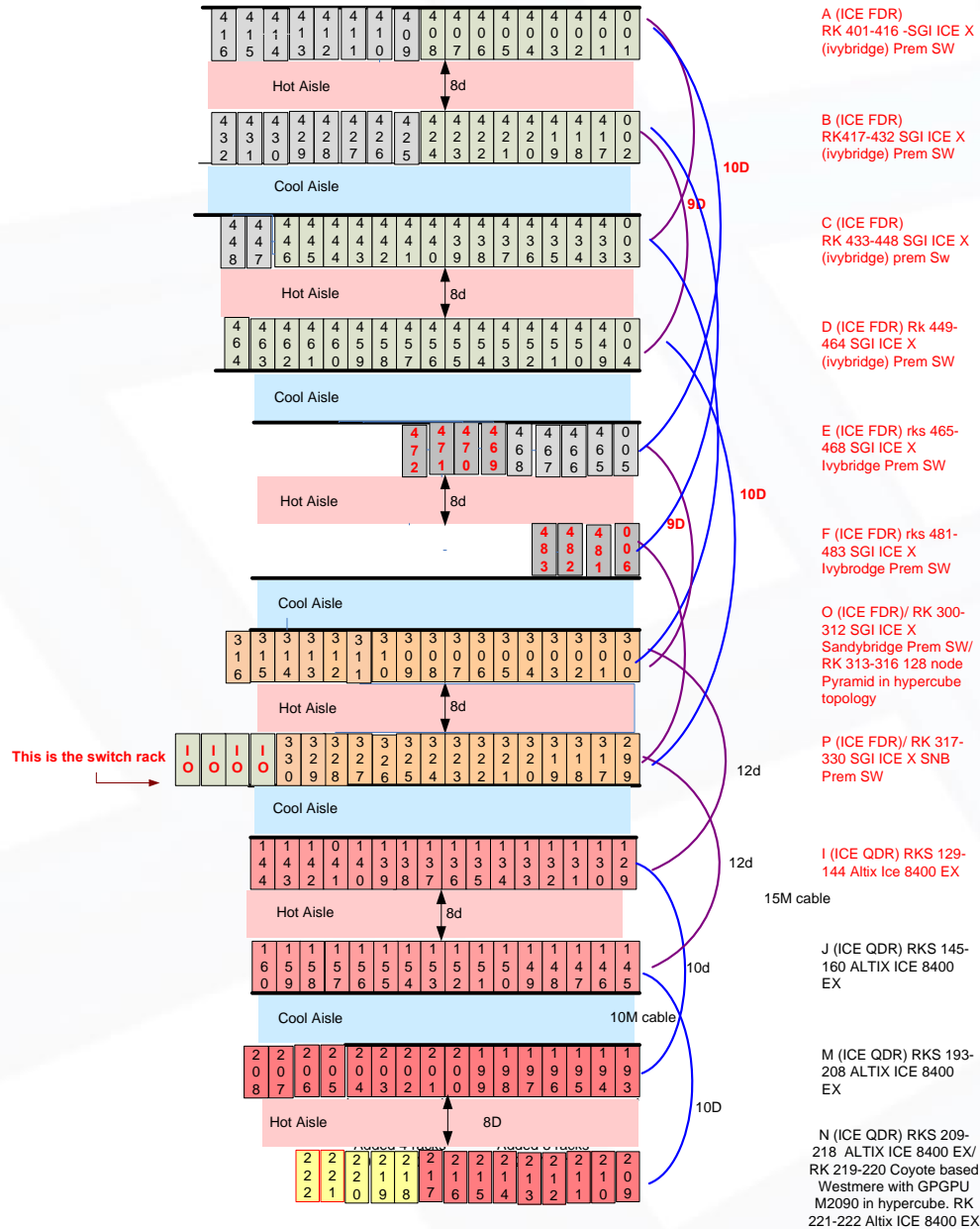
NASA (Pleiades) Rack Layout as of 1/30/2014



NASA (Pleiades) Rack Layout as of 2/18/2014



NASA (Pleiades) Rack Layout as of 2/25/2014



170 racks – 2014
3.5 petaflops

This is the switch rack



Incremental Expansion – Driving Factors

- Annual Funding/Budget Uncertainty
- Synthetic Leases/Sarbanes-Oxley cost
- Risk Mitigation for Fast moving technology
- Supports Short Lead/Opportunistic Strategy
- Timed adoption based on technology readiness
- Decouples technologies on different development cycles
- Dynamic project support

Maintains leading edge components throughout our
“Ground Based Instrument”



Production Filesystems

- 16 different filesystems
 - 8 lustre scratch filesystems spread over about 12PB
 - Speed Optimized
 - IOP Optimized
 - Dedicated
 - 8 nfs filesystems
 - Home, scratch, system images



Production Software Environment

- 4 different production selectable operating systems
 - AOE: 3 sles, centos
 - Additional test images
- 251 different loadable modules
 - 58 different compilers (32 intel, 8 PGI, 4 gcc, 3 cuda, 3 matlab...)
 - 26 different MPIs (10 SGI MPT, 12 Intel MPI, 8 MVAPICH)
 - 23 libraries (13 hdf, 6 netcdf, 4 mkl)
- Various debuggers, performance analyzers, plotting/graphing, editors
- Driven by user requests/requirements

This is an HPC Cloud



What is Today's General Purpose Supercomputer

- 1980s/1990s – a monolithic system with limited access
 - Typically served smaller communities
 - Local dedicated disk with limited network connectivity
- Today – its a collection of heterogeneous elements both SW & HW
 - Supports a wide variety and types of computation
 - Tuned for user productivity
- General Purpose - a compromise in some ways
 - MAY not be the #1 top 500 machine
 - But should be the most productive for highly varied requirements in multiple science and engineering domains.



How does this relate to Exascale

- NASA will get there – Incrementally
 - As long as doing so enables engineering and scientific discovery
- Exascale systems are likely to be more heterogeneous - diverse
 - More specialized components
 - Increased relevance around GPU/ARM/Phi
 - Cost saving realized from leveraging existing infrastructure
- Exascale systems will have significantly more parallelism
 - More things to break
- **Grand'ist challenge – dealing with component failure of HW/SW**

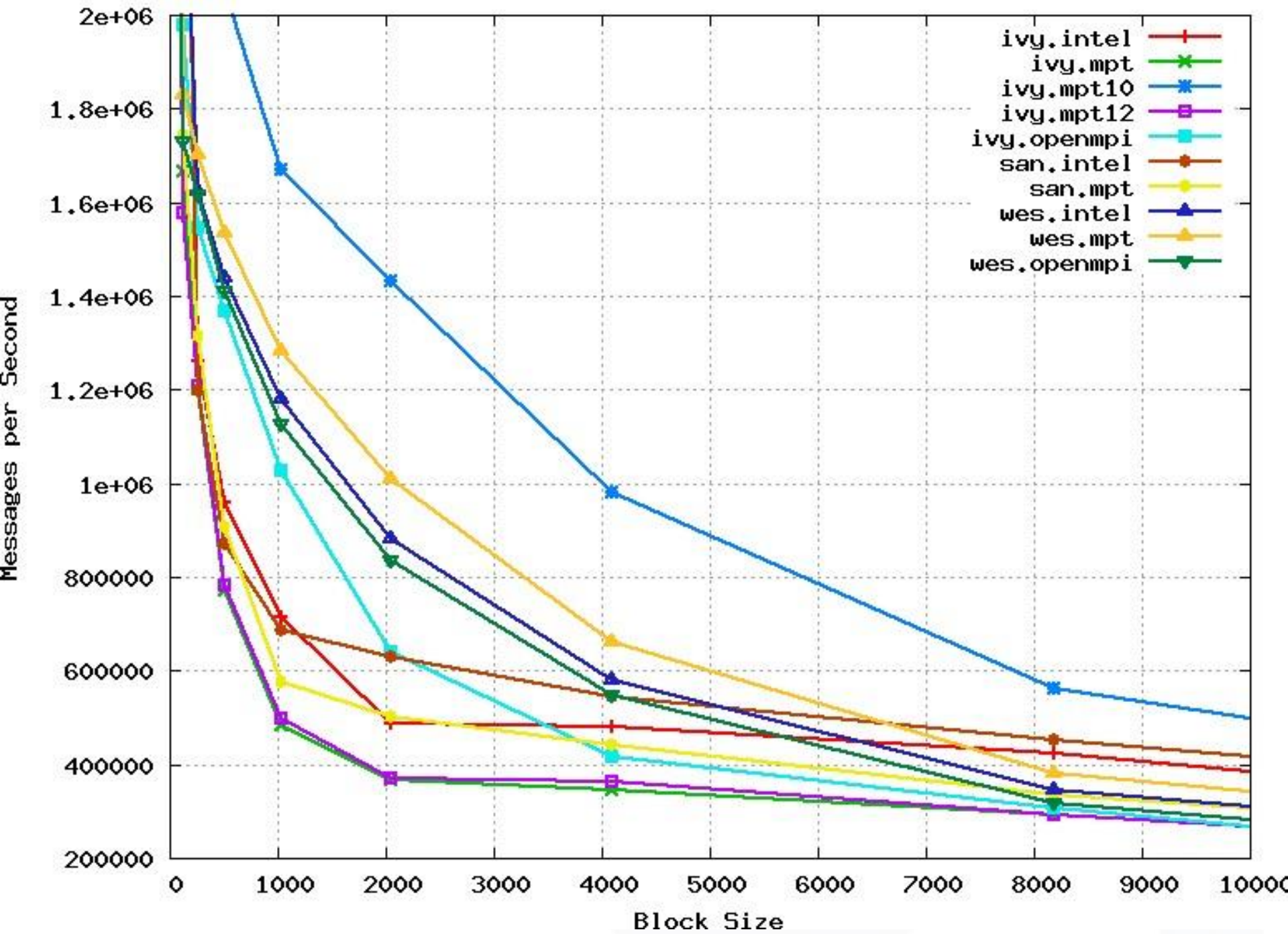


Infiniband Limitations (10x)

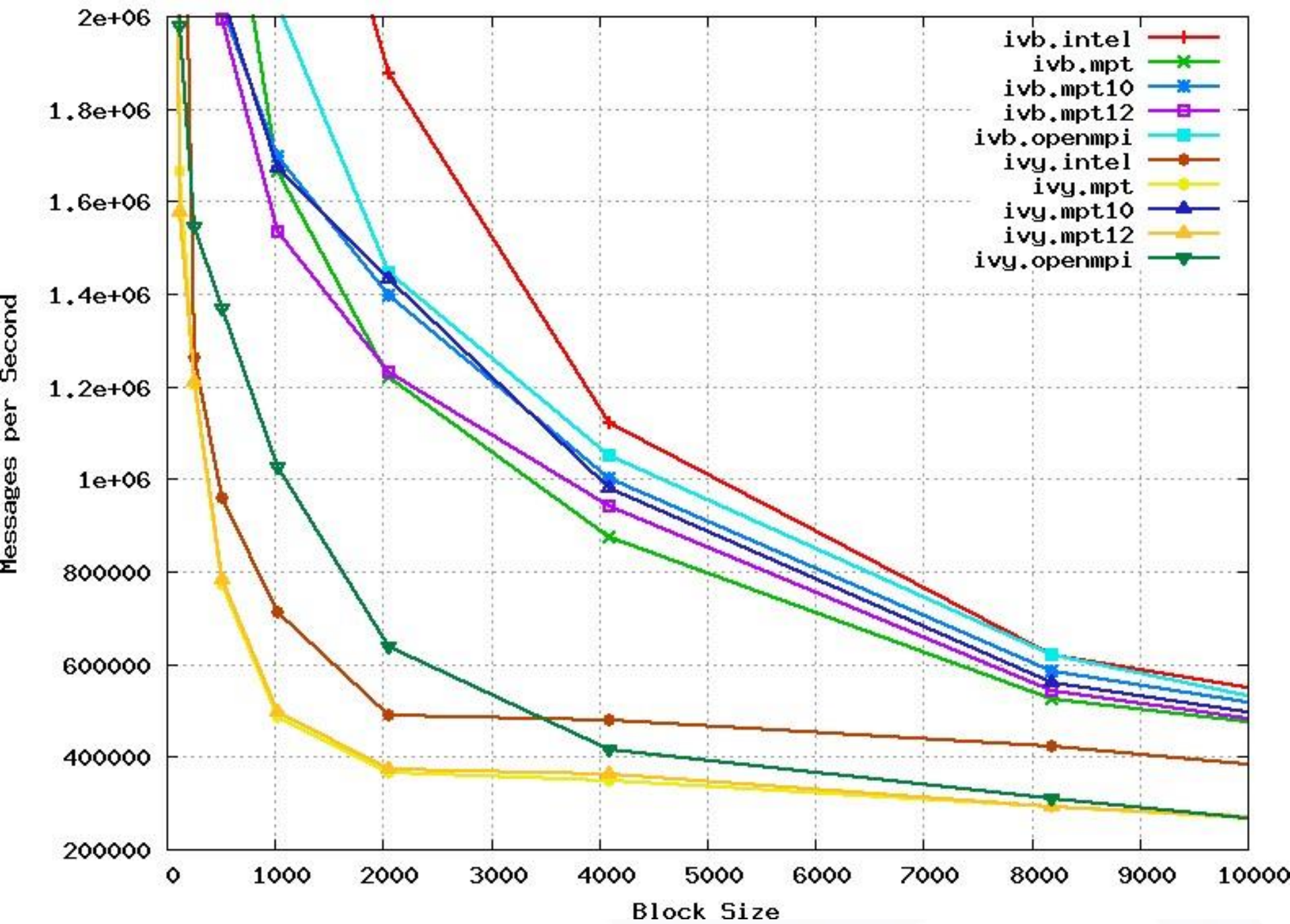
- Infiniband Cables
 - Physical layer issues – diagnostics – light queues – SM support
- Subnet Management
 - Does not scale - Scans, SA, multi-cast, restarts, continuous operation
- Packet Loss/RC not reliable
- Bonding Performance/Failover



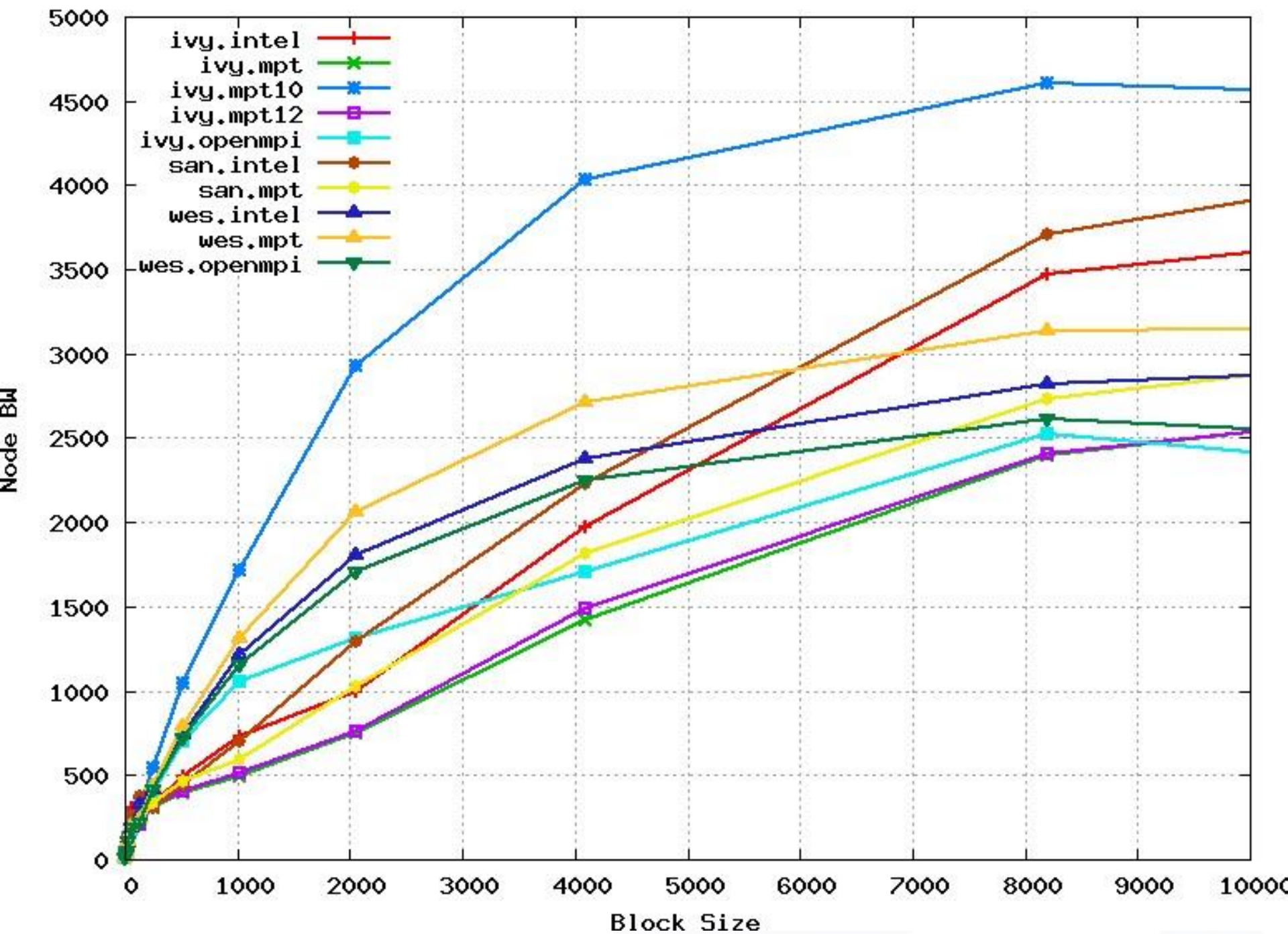
Node Messages/Second



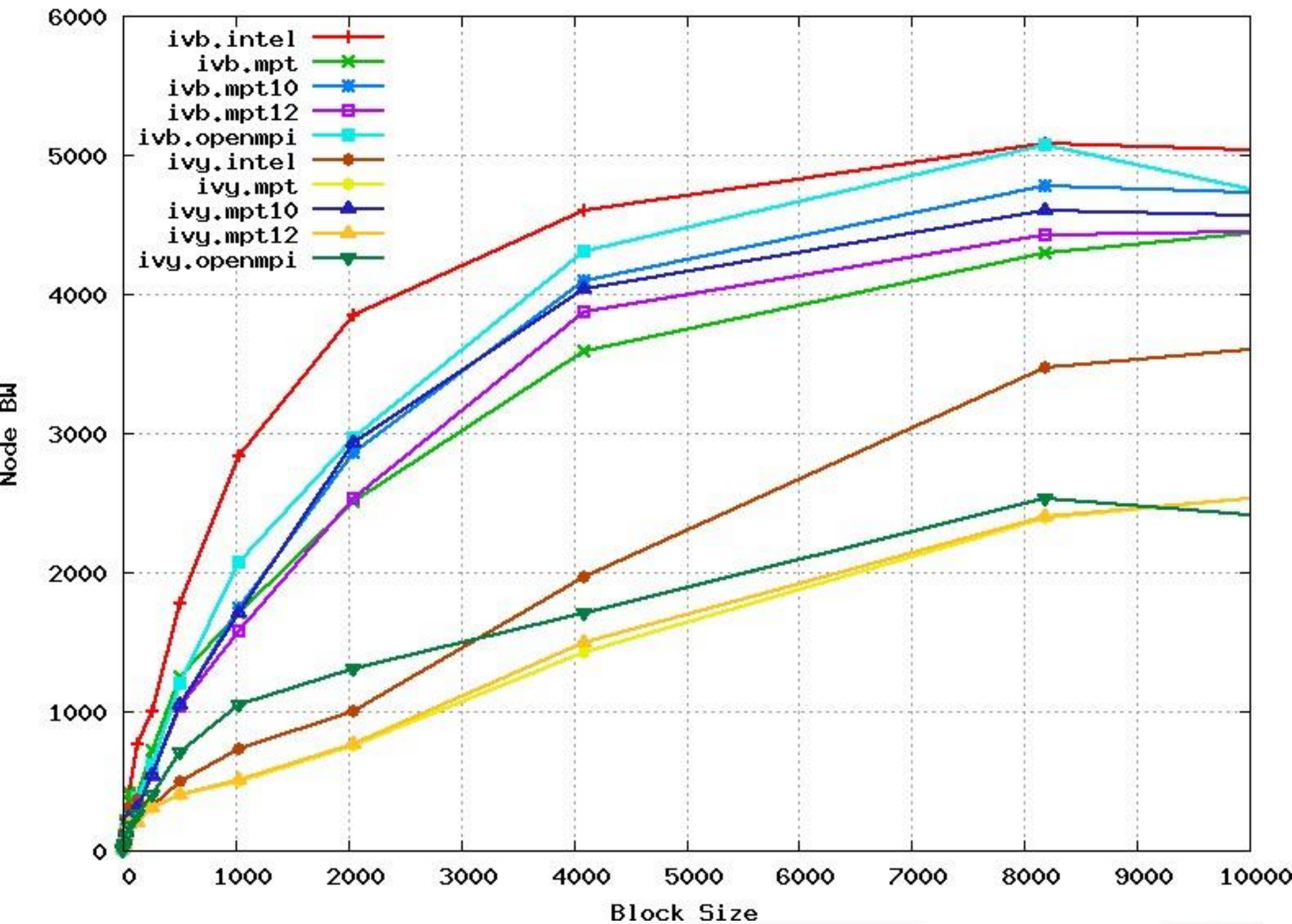
Node Messages/Second BIOS



Total Node BW (mult by 2x to get full-duplex rate)



Total Node BW (mult by 2x to get full-duplex rate) BI05





Infiniband Limitations (10x)

- Congestion Control
 - Getting worse for us
 - Nodes more easily overwhelm fabric
 - Discards more common
 - Particularly problematic in one code
 - Retry mechanisms too simply
 - Who is causing the problem
 - Slow node
 - Improper Software
 - Link Recovers
 - Slow Management



SM/SA Requirements

- Does not scale
 - Scans, SA, multi-cast, restarts, continuous operation
- Must handle 1000's of nodes:
 - Crashing
 - Rebooting
 - Being slow or non-responsive or broken
- System Maintenance and Expansion
 - Active SW GUIDS, Active Ports (cable maint *)
- Large Scale systems have regularity to leverage static information
 - SW Ports known to have HCA (*)
 - Static or locally discoverable) information (LIDs, QP services)
 - Decentralized/Distributed services (H-ARP *),
- Separate Switch and HCA management
- Intelligent Switches (am I connected to the right components)

(* working these)



Conclusions

- Backward and forward compatible high speed networks have big advantages in leveraging existing infrastructure.
- Future networks need to address interoperability at a reasonable cost.
 - Protocol translators cannot require significant HW (e.g. server)
- Converged/Heterogeneous Super Scale environments benefit
 - Agile system implementations
 - Cost saving by leveraging existing infrastructure
 - Quick deployments
- Infiniband has been the key to implementing this strategy
- **Most Important Factor – performance/reliability/cost of compute**

