# NASA SuperComputing
# A Ground Based Instrument for Exploration and Discovery

## LUG 2015

Bob Ciotti

Chief Architect/Supercomputing Systems Lead

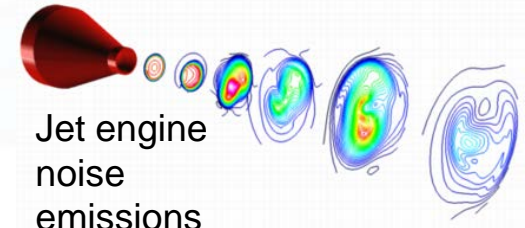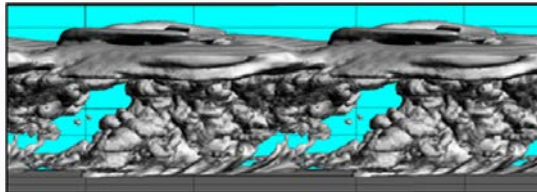LUG 2015 - Denver

# Discussion

- What is Pleiades
- The NASA Workload
- System Build Strategy
- Operational Strategy
- Tools and Analysis Software
- Issues Do We See
- Whats Lustre Does
- What We Want

# Supercomputing Support for NASA Missions

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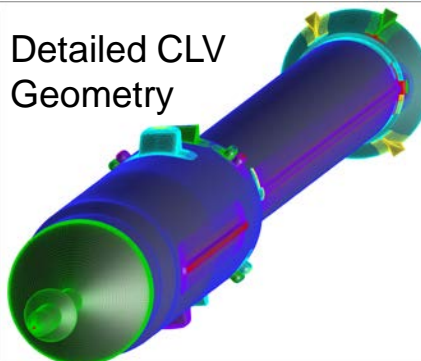  - Typically o500 jobs running

- 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
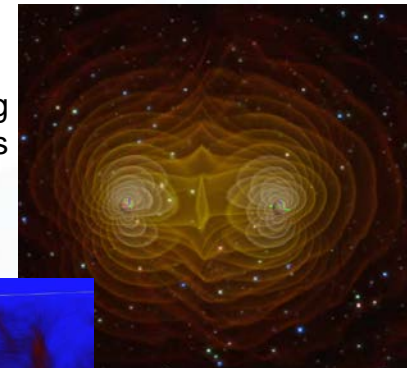
Jet aircraft wake vortices

Jet engine noise emissions

Detailed CLV Geometry
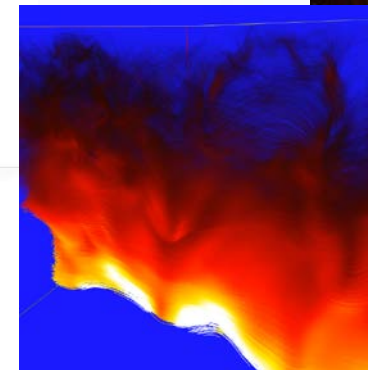
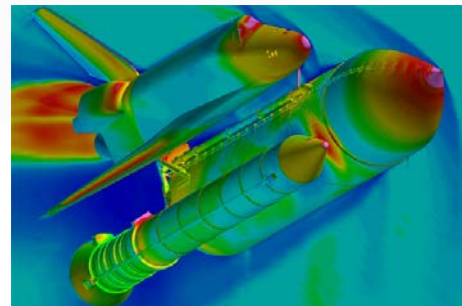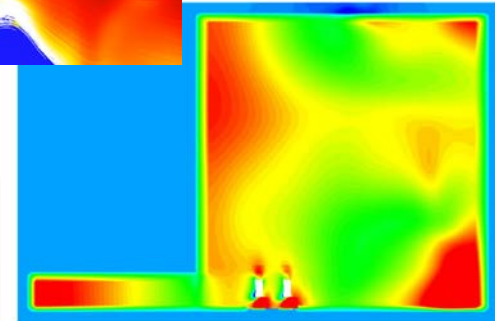Orbiting, Spinning Black Holes

Solar surface convection

Shuttle Ascent Configuration

2-SRB Burn in VAB

KOI-157 :: Teff = 5685   logg = 4.38   Rs = 1.06

# ECCO – Ocean Modeling



Jan 1    2003

# Planetary Defense

DON DAVIS
3-27-91

time = 0.01 secs [0000200]

temp

National Aeronautics and Space Administration          LUG - Denver                    Apr 2015

# NASA's Computational Landscape

**Embarrassingly Parallel**

**Compute Bound**

**Simple Well Understood Computations**

**Tightly Coupled**

**Highly Complex and Evolving Computations**

**Data/Storage Intensive**

# NASA Advanced Supercomputing | Systems Status

## ———————— Pleiades Mission Directorate Usage ————————

| Summary | ARMD | HEOMD NESC | SMD | Other |
|---------|------|------------|-----|-------|
| 183198 / 208116 | 51504 / 57232 | 45404 / 55151 | 84614 / 89490 | 1676 / 6243 |
| 88   671+1084 | 90   114+780 | 82   40+118 | 95   348+591 | 27   1+18 |

### Pleiades Special Queues

| Devel | GPU |
|-------|-----|
| 18936 / 26960 | 0 / 1488 |
| 70   14 | 0   0 |

### Pleiades Node Utilization   *(bars scaled by computing capability)*
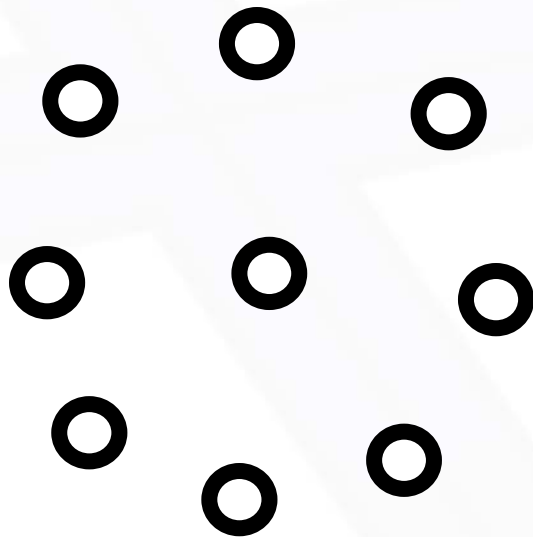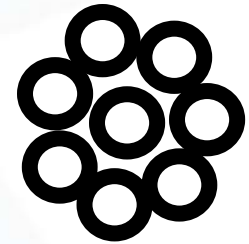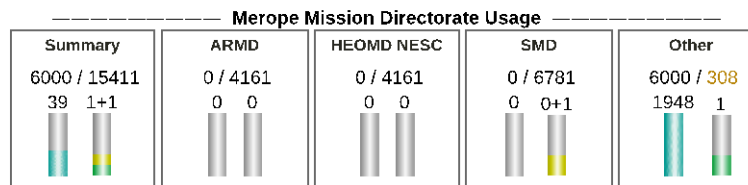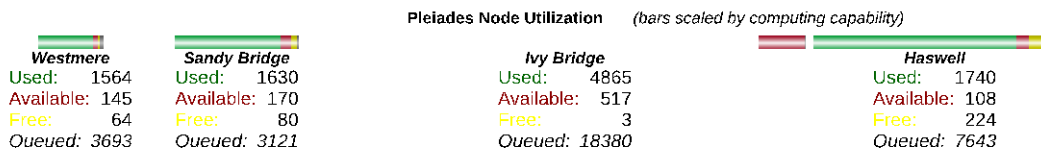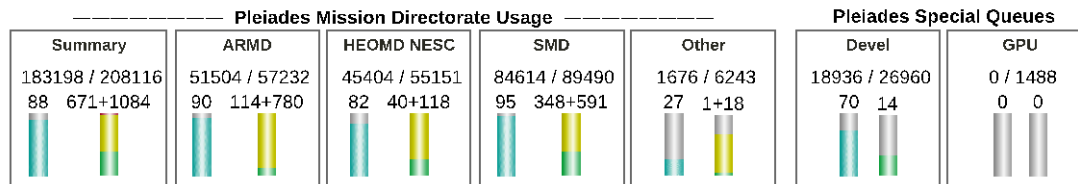
**Westmere**
Used:        1564
Available:   145
Free:         64
Queued: 3693

**Sandy Bridge**
Used:        1630
Available:   170
Free:         80
Queued: 3121

**Ivy Bridge**
Used:        4865
Available:   517
Free:          3
Queued: 18380

**Haswell**
Used:        1740
Available:   108
Free:         224
Queued: 7643

## ———————— Merope Mission Directorate Usage ————————

| Summary | ARMD | HEOMD NESC | SMD | Other |
|---------|------|------------|-----|-------|
| 6000 / 15411 | 0 / 4161 | 0 / 4161 | 0 / 6781 | 6000 / 308 |
| 39   1+1 | 0   0 | 0   0 | 0   0+1 | 1948   1 |

Merope node utilization data is currently unavailable.

## ———————— Endeavour Mission Directorate Usage ————————

| Summary | ARMD | HEOMD NESC | SMD | Other |
|---------|------|------------|-----|-------|
| 1152 / 1512 | 48 / 408 | 0 / 408 | 1104 / 665 | 0 / 30 |
| 76   33+263 | 12   2+2 | 0   0 | 166   31+261 | 0   0 |

Pleiades / Merope / Endeavour Mission "Used/Total CPUs":        Used / Currently Available        Used / *Allocated*

### ———————— Host Groups ————————

Lou Cluster    NAS Front Ends    Pleiades Front Ends    Bridge Nodes    NFS Systems

### — Filesystems —

Pleiades

# Pleiades

# Pleiades

# Pleiades

Snapshot Size Distribution

# SGI ICE Dual Plane – Topology



n0
n1
n2
n3
n4
n5
n6
n7
n8

ib0

2x 11d hypercube
full     2048 vertices
Pleiades   1336/11d (2672 across both cubes)

ib1

http://en.wikipedia.org/wiki/User:Qef/Orthographic_hypercube_diagram
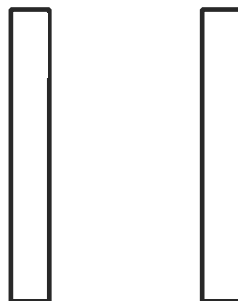
# Infiniband Subnet LAN

LAN Implemented with out board IB switches

Archive Servers

NFS File Servers

Hyperwall Graphics System

Data Transfer Nodes

Bridge Nodes

Front End Nodes

Data Analysis Nodes

Lustre Filesystems



Orthographic demidekeract
by Claudio Rocchini, wikipedia
Copyright GNU http://en.wikipedia.org/wiki/GNU_Free_Documentation_License
Creative Commons 3.0  http://creativecommons.org/licenses/by/3.0

# I/O Network

105 OSS+MDS

480 GB/sec — Lustre Server

r999

382 GB/sec

107 GB/sec

428 GB/sec ib0+ib1

r998

857 GB/sec — Hyperwall 128-Display Graphics Array

I/O fabric

ib1

64 racks – 2008
393 teraflops

NASA (Pleiades) Rack Layout

112 racks – 2009
683 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

NASA (Pleiades) Rack Layout

186 racks – 2011
1.33 petaflops

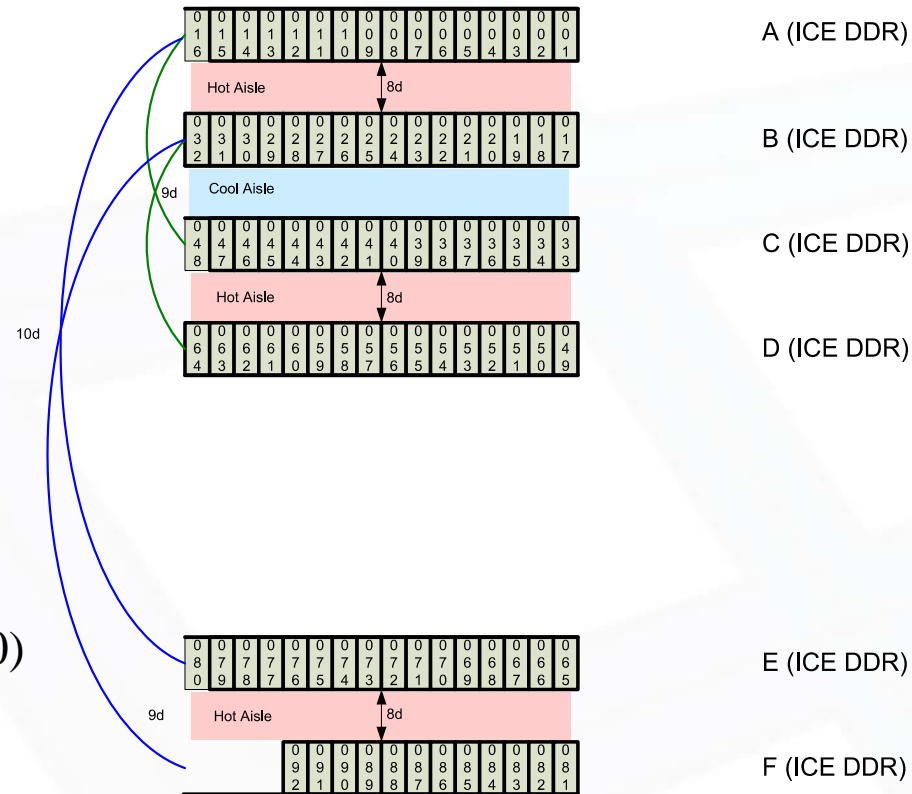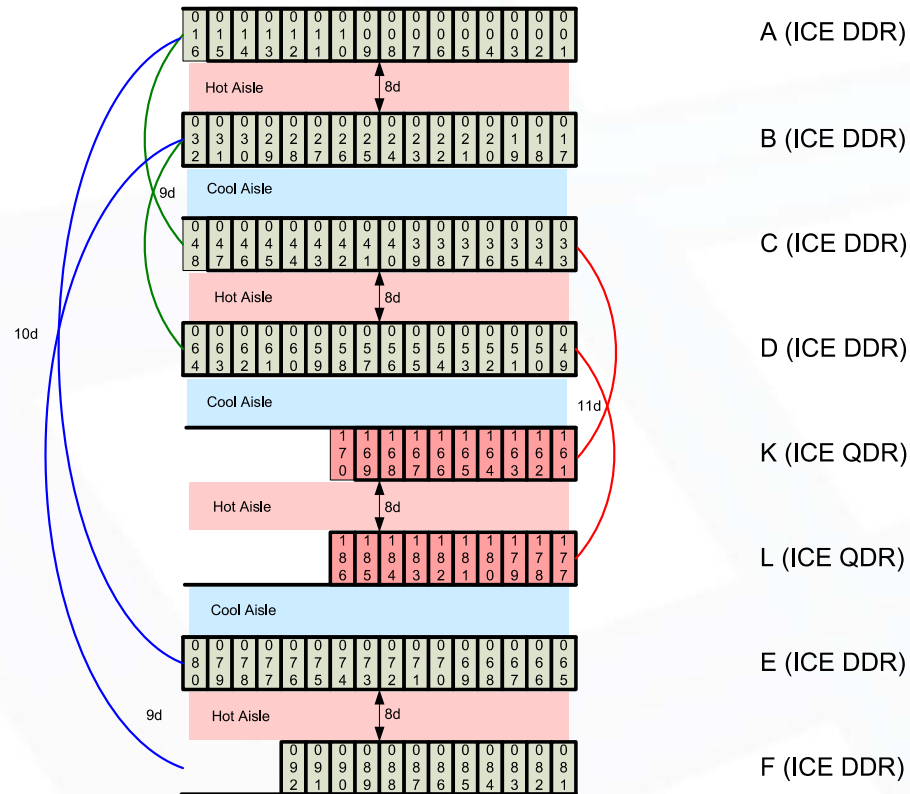# Pleiades - Sustained SpecFP rate base

- **SpecFP rate base <u>estimates</u>** (eliminates cell/GPU/blue-gene/SX vec)

| Spec | Top500 | Machine | CPU | #Sockets | FPR/Socket | TSpec |
|---|---|---|---|---|---|---|
| 1 | 2 | Jaguar | AMD-2435 | 37,360 | 65.2 | 2,436,246 |
| 2 | **6** | Tera-100 | Intel-7560 | 17,296 | 133.4 | 2,307,805 |
| 3 | **5** | Hopper | AMD-6176 | 12,784 | 149.8 | 1,800,115 |
| 4 | **1** | Tianhe-1a | Intel-x5670 | 14,336 | 119.5 | 1,713,868 |
| **5** | **11** | **Pleiades** | **Intel-x** | **21,632** | **72.2** | **1,562,510** |
| 6 | **10** | Cielo | AMD-6136 | 13,394 | 115.5 | 1,547,408 |
| 7 | 8 | Kraken | AMD-2435 | 16,448 | 65.2 | 1,075,182 |
| 8 | 14 | RedSky | Intel-x5570 | 10,610 | 90.3 | 958,401 |
| 9 | 17 | Lomonosov | Intel-x5570 | 8,840 | 90.3 | 798,517 |
| 10 | 15 | Ranger | AMD-2356 | 15,744 | 37.3 | 588,196 |

- Tspec == number of 2-core 296mhz UltraSPARC II

NASA (Pleiades) Rack Layout

158 racks – 2012
1.15 petaflops
deinstall

A (ICE DDR)
B (ICE DDR)
C (ICE DDR)
D (ICE DDR)
K (ICE QDR)
L (ICE QDR)
O (ICE FDR)
P (ICE FDR)
I (ICE QDR)
J (ICE QDR)
M (ICE QDR)
N (ICE QDR)

*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.

National Aeronautics and Space Administration

Apr 2015

NASA (Pleiades) Rack Layout

158 racks – 2012 deinstall

158 racks – 2012 deinstall

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

Cool Aisle

K (ICE QDR)

Hot Aisle

L (ICE QDR)

Cool Aisle

O (ICE FDR)

Hot Aisle

P (ICE FDR)

Cool Aisle

I (ICE QDR)

Hot Aisle

J (ICE QDR)

Cool Aisle

M (ICE QDR)

Hot Aisle

N (ICE QDR)

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

167 racks – 2013
2.9 petaflops

National Aeronautics and Space Administration

LUG - Denver

Apr 2015

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

160 racks – 2013
3.1 petaflops

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

168 racks – 2013
3.2 petaflops

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

168 racks – 2014
3.3 petaflops

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

170 racks – 2014
3.5 petaflops

NASA (Pleiades) Rack Layout

168 racks – 2014
4.5 petaflops

NASA (Pleiades) Rack Layout

168 racks – 2015
5.4 petaflops

# Pleiades 2015 – Based on MemoryBW (ignore GPU/PHI)

| Machine | Type | 11/14 T500 | Sockets | Type | Mem BW Socket | Spec Socket | Mem BW (PB/Sec) | Mega Spec | Rmax | Rpeak | PctPeak |
|---|---|---|---|---|---|---|---|---|---|---|---|
| K computer | Sparc64 | 4 | 88,128 | VIII fx | 64.0 | 373.2 | 5,640 | 32.9 | 10,510 | 11,280 | 93.2% |
| Sequoia | BGQ/Power | 3 | 98,304 | BGQ-A2 | 42.7 | 144.3 | 4,198 | 14.2 | 17,173 | 20,132 | 85.3% |
| BlueWater | XK6/XK7 | | 49,200 | 6276 | 51.2 | 176.0 | 2,519 | 8.7 | | 71,378 | |
| Mira | BGQ /Power | 5 | 49,152 | BGQ-A2 | 42.7 | 144.3 | 2,099 | 7.1 | 8,586 | 10,066 | 85.3% |
| Tianhe-2 | Xeon/Xeon Phi | 1 | 32,000 | E5-2692v2 | 59.7 | 321.5 | 1,910 | 10.3 | 33,862 | 54,902 | 61.7% |
| **Pleiades** | **SGI/Xeon Mix** | **11** | **22,896** | **XeonMix** | **54.8** | **283.7** | **1,255** | **6.5** | **3,375** | **3,987** | **84.7%** |
| Juqueen | BGQ/Power | 8 | 28,672 | BGQ-A2 | 42.7 | 144.3 | 1,224 | 4.1 | 5,008 | 5,872 | 85.3% |
| Secret2 | XC30/Xeon | 13 | 18,832 | E5-2697v2 | 59.7 | 341.0 | 1,124 | 6.4 | 3,143 | 4,881 | 64.4% |
| Vulcan | BGQ/Power | 9 | 24,576 | BGQ-A2 | 42.7 | 144.3 | 1,049 | 3.5 | 4,293 | 5,033 | 85.3% |
| Titan | XK7/Opteron/K20x | 2 | 18,688 | 6274 | 51.2 | 173.0 | 957 | 3.2 | 17,590 | 27,112 | 64.9% |
| SuperMUC | iData/Xeon | 14 | 18,432 | E5-2680 | 51.2 | 244.5 | 944 | 4.5 | 2,897 | 3,185 | 91.0% |
| Pangea | SGI/Xeon | 20 | 13,800 | E5-2670 | 51.2 | 240.5 | 707 | 3.3 | 2,098 | 2,296 | 91.4% |
| Stampede | Dell/Xeon/Phi | 7 | 12,800 | E5-2680 | 51.2 | 244.5 | 655 | 3.1 | 5,168 | 8,520 | 60.7% |
| Hornet | XC40/Xeon | 16 | 7,884 | E5-2680v3 | 68.0 | 396.5 | 536 | 3.1 | 2,763 | 3,784 | 73.0% |
| Tianhe-1A | Xeon/Nvidia2050 | 17 | 14,336 | X5670 | 32.0 | 132.0 | 459 | 1.9 | 2,566 | 4,701 | 54.6% |
| Secret1 | CS/Xeon/K40 | 10 | 7,280 | E5-2660v2 | 59.7 | 287.5 | 435 | 2.1 | 3,577 | 6,131 | 58.3% |
| HPC2 | iData/Xeon/K20x | 12 | 7,200 | E5-2680v2 | 59.7 | 313.0 | 430 | 2.3 | 3,188 | 4,605 | 69.2% |
| Excalibur | XC40/Xeon | 19 | 6,254 | E5-2698v3 | 68.0 | 434.0 | 425 | 2.7 | 2,485 | 3,682 | 67.5% |
| Piz Daint | XC30/Xeon/K20x | 6 | 5,272 | E5-2670 snb | 51.2 | 240.5 | 270 | 1.3 | 6,271 | 7,788 | 80.5% |
| Cascade | Xeon/Xeon Phi | 18 | 1,880 | E5-2670 | 51.2 | 240.5 | 96 | 0.5 | 2,539 | 3,388 | 74.9% |
| Tsubame | Nec/Xeon/K20x | 15 | 2,816 | X5670 | 32.0 | 132.0 | 90 | 0.4 | 2,785 | 5,735 | 48.6% |

Numbers in Red are sWAG

# Pleiades Environment

- 11,280 compute nodes – 22,560 sockets - 211,360 x86 cores
- 128 visualization nodes
- 192 GPU Nodes
- 192 Xeon Phi Nodes
- 10 Front End Nodes
- 4 "Bridge Nodes"
- 4 Archive Front Ends
- 8 Data Analysis Nodes
- 8 Archive Nodes
- 2 large memory nodes 2 TB + 4 TB

- Everything cross mounted. NFS Home, Modules, Nobackup (NFS, lustre)

- + a couple hundred administration/management nodes of various types.

# Pleiades /nobackup Filesystems
## (production)

Lustre

| Filesystem | #OST | #OSS | Size (PB) | Write BW | Read BW | controller |
|---|---|---|---|---|---|---|
| p5 | 180 | 12 | 2.2 | 19 | 15 | DDN SFA10k |
| p6 | 180 | 12 | 1.3 | 17 | 14 | DDN SFA10k |
| p7 | 84 | 18 | 1.9 | 18 | 18 | NetApp 5400 |
| p8 | 312 | 26 | 7.1 | 65 | 52 | NetApp 5400 |
| p9 | 240 | 18 | 3.5 | 22 | 21 | DDN SFA10k |
| **996** | **86** | **16.0** | **141** | **120** | |

NFS

# 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 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 Todays
# 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.

# Continuous Availability
# 24/7 Operations

– Goal – never take the whole system down

  • Outages are very disruptive

  • Dedicated time very costly

  • Not even possible to update entire system in one dedicated session.

  • Things go wrong

• Components

  • Lustre, NFS, CXFS, OFED, OpenSM, Linux Distro patches, cluster management software,

  • Firmware

    • its in everything – including the cables.

# Continuous Availability
# 24/7 Operations

- Rolling updates of various components
  - Lustre/NFS clients/compute node images
    - Easy – simply done at end of user job
  - NFS, Lustre servers
    - Hot swap
      - Nfs hard mounts
      - Lustre recovery
    - Suspend/Resume

  - Schedule filesystems as a resource in addition to nodes
    - Allow us to use all compute nodes and figure out share later
  - Various admin, front ends, bridge nodes are easier or less urgent.

# Continuous Availability
# 24/7 Operations

- Hot Plug - Grow system while in operation
  - Cable up new components powered off
    - Check cabling
  - Signal OpenSM to turn off sweep
  - Power on equipment
  - Run ibnetdiscover to verify cabling
  - Signal OpenSM to sweep
  - Mount file systems and go
- Cable Maintenenace

# PBS Lustre Stats

**Exiting at : Sat Mar 14 14:34:55 2015**

================================================================================

**LUSTRE Filesystem Statistics**

--------------------------------------------------------------------------

nbp8 Metadata Operations

| open | close | stat | statfs | read (GB) | write (GB) |
|------|-------|------|--------|-----------|------------|
| 1056469 | 1056469 | 1058349 | 0 | 274 | 312 |

| Read | 4KB | 8KB | 16KB | 32KB | 64KB | 128KB | 256KB | 512KB | 1024KB |
|------|-----|-----|------|------|------|-------|-------|-------|--------|
|  | 114 | 147 | 1 | 16 | 9 | 29 | 144 | 748 | 48185 |
| Write | 4KB | 8KB | 16KB | 32KB | 64KB | 128KB | 256KB | 512KB | 1024KB |
|  | 5091 | 51 | 51 | 353 | 36 | 48 | 2120 | 49 | 297141 |

---

Job Resource Usage Summary for 3075801.pbspl1.nas.nasa.gov

| | |
|---|---|
| CPU Time Used | : 259:36:54 |
| Real Memory Used | : 37024436kb |
| Walltime Used | : 10:52:49 |
| Exit Status | : 0 |
| Number of CPUs Requested | : 816 |
| Walltime Requested | : 24:00:00 |
| Execution Queue | : sls_aero1 |
| Charged To | : e0847 |
| Job Stopped | : Sat Mar 14 14:35:36 2015 |

---

# File Transfer - Shiftc

- File transfers have become quite complex:
  - Best source/destination
    - Systems have multiple interfaces – want to pick best one
  - Threading
    - Big performance wins by parallelizing within a node
    - Big performance wins by parallelizing across nodes
  - Error checking
    - Checksum
      - Partial resend for hash mismatches
      - Ability to save partial hash to detect location of corruptions
  - Restart/Completion
    - Systems fail or reboot
      - Will restart transfer and notify upon completion

  - Alternative to lustre-hsm, but some potential application…
  - Multi GB/sec transfer within a filesystem
  - Working on similar capability to DMF Archive

- Credit: Paul Kolano

# Log File Analysis

- Lumber - Tool written to go through all the log file data (GBs/day)
    - Lustre logs
        - Server and Clients
    - PBS Logs
    - Console Logs
    - System Logs

    - Absolutely necessary to track system issues

- Can specify a job ID and get all the log information across all systems during that timeframe.

- Can do arbitrary searches – across all logs

- Credit: Dave Barker

# Daily Failure Logs – Past 24 hours

Daily Report for 04/10/2015 on pbspl1

Job Failure Summary from Fri Apr 10 00:00:00 2015 to Fri Apr 10 23:59:59 2015

There were 3197 jobs in the time region, of which 22 indicate as failed.

The total SBUs of these jobs was 500795.64, of which 6.70 (%0.00) belonged to the failed jobs.

Job Failure Summary Sorted by Frequency of Failure Types:

| Count | UID/GID | SBUs | Failure type |
|-------|---------|------|--------------|
| 8 | 6/6 | 0.00 (% 0.00) | head node lost connection with a sister node |
| 6 | 5/5 | 6.38 (% 0.00) | job experienced out of memory (oom) |
| 5 | 3/3 | 0.00 (% 0.00) | the PBS Server discarded the job because it appeared a node was down |
| 1 | 1/1 | 0.05 (% 0.00) | job produced too much spool output (stdout/stderr) |
| 1 | 1/1 | 0.28 (% 0.00) | PBS unable to start job |
| 1 | 1/1 | 0.00 (% 0.00) | PBS server lost connection with head node |

# Weekly Failure Logs – Past 24 hours

Daily Report for last 7 days to 04/10/2015 on pbspl1

Job Failure Summary from Sat Apr  4 00:00:00 2015 to Fri Apr 10 23:59:59 2015

There were 14650 jobs in the time region, of which 148 indicate as failed.

The total SBUs of these jobs was 3598210.40, of which 239289.38 (%6.65) belonged to the failed jobs.

Job Failure Summary Sorted by Frequency of Failure Types:

| Count | UID/GID | SBUs | Failure type |
|---|---|---|---|
| 54 | 19/17 | 480.60 (% 0.01) | job experienced out of memory (oom) |
| 24 | 3/3 | 0.00 (% 0.00) | job start error 15010, node could not JOIN_JOB successfully |
| 8 | 5/5 | 1361.84 (% 0.04) | job produced too much spool output (stdout/stderr) |
| 8 | 6/6 | 0.00 (% 0.00) | the PBS Server discarded the job because it appeared a node was down |
| 8 | 6/6 | 0.00 (% 0.00) | head node lost connection with a sister node |
| 7 | 5/5 | 0.00 (% 0.00) | the PBS Server discarded the job for unknown reasons |
| 6 | 4/2 | 145034.72 (% 4.03) | MPT error - receive completion flushed |
| 6 | 2/2 | 210.25 (% 0.01) | node had RCU sched stalls |
| 5 | 3/2 | 46686.32 (% 1.30) | MPT error - MPI_SGI_ctrl_recv failure |
| 5 | 5/5 | 1553.60 (% 0.04) | node dropped into kdb |
| 4 | 3/3 | 6074.78 (% 0.17) | MPT error - xmpi_net_send failure |
| 3 | 3/3 | 3584.49 (% 0.10) | job experienced uncorrectable ecc memory error |
| 2 | 2/2 | 90.62 (% 0.00) | at least one node associated with the job booted for unknown reasons |
| 2 | 2/2 | 0.00 (% 0.00) | mlx4 internal error |
| 2 | 2/2 | 0.26 (% 0.00) | PBS server lost connection with head node |
| 1 | 1/1 | 34110.72 (% 0.95) | MPT error - continuous IB fabric problems |
| 1 | 1/1 | 47.64 (% 0.00) | MPT error - network error in starting shepherd |
| 1 | 1/1 | 53.27 (% 0.00) | MPT error - shepherd terminated |
| 1 | 1/1 | 0.28 (% 0.00) | PBS unable to start job |

# Real Time I/O Monitor

```
Every 1.0s: abracadabra -i 1
Mar 26 00:31:37 2012

    io_swx     nbp1      .      nbp2      .     nbp3/4     .      nbp5      .      nbp6      .       tot       .
       .       read    write    read    write    read    write    read    write    read    write    read    write

  r999i_mds     0.7      0.4     2.4      1.4     16.7     11.5     0.3      0.3     1.3      0.7     20.7     13.9
  r999i_oss1    2.3      6.5    18.4    208.5      4.1     11.6     2.2      2.2     2.3      2.3     11.0     22.6
  r999i_oss2    3.5    122.1     2.8     51.3      2.5      7.0     2.2      2.3     2.3      2.3     13.4    184.9
  r999i_oss3    2.3      9.7    16.0     39.7      2.5      4.8     2.2      2.2     2.3      3.2     25.3     59.6
  r999i_oss4    2.3      8.1    79.9     34.1      2.4      4.0     2.2      2.2     2.3      2.2     89.2     50.7
  r999i_oss5    2.4      9.0     2.7     42.5      2.2     10.4     2.2      2.2     2.2      2.3     11.7     66.4
  r999i_oss6    2.3     10.6     6.4     38.7      2.2      5.6     2.2      2.2     2.2      2.2     15.5     59.4
  r999i_oss7    2.3     10.6     6.3     23.5      2.2     12.3     2.2      2.2     2.2      2.2     15.3     50.8
  r999i_oss8    2.3     10.2   270.5     35.7      2.2      7.1     2.2      2.2     2.2      3.2    279.3     58.4
      Total    20.4    187.2   405.4    475.4     37.0     74.3    17.9     18.0    19.3     20.6    481.4    566.7
        Max  2809.2  16138.9  5943.9   5003.6   2310.6   4719.3    50.9    171.3 14930.3  15173.6 15127.3  16845.9


   Max  RcvData: 1514.8 8451.6 3319.8 1252.6 6261.4 7874.4 14207.8 3903.5 10441.4 8181.3 6720.7 5473.9   7.1   3.6   9.2   1.9   8.8   1.7  11.1   1.2   3.6 16847.1
   Max XmitData:   14.1 1393.7 6645.3 3405.3 1478.8 5506.1 13417.8 1675.2  2846.6 2498.5 1365.8 1210.5   8.8   2.0   6.9   3.8  10.4   1.2   8.9   2.1   4.7 15130.8

 Total  RcvData:    0.1   62.4    4.1    6.0    5.7   14.4   52.2   22.8  128.4   18.4  171.4  288.3   0.3   0.1   0.3   0.0   0.2   0.3   0.3   0.3   1.3  777.6
 Total XmitData:    0.1   17.7   11.2    6.4    6.3  105.0   15.0   15.0    8.9    9.8    2.8  301.8   0.3   0.1   0.3   0.1   0.3   0.3   0.2   0.4   1.3  502.7

r999i_mds         .     .  r41i0  r49i1  r57i1  r17i0  r25i0 r129i0 r137i0 r145i0 r153i0     .   r9i0  oss1  oss1  oss2  oss2  oss3  oss3  oss6  oss6 hwsw0    tot
r999i_mds RcvData:  0.0   0.2    0.6    0.4    0.2    0.1    0.7    2.0    0.3    1.2   0.0   8.5   0.1   0.1   0.1   0.0   0.2   0.1   0.2   0.1   0.0   15.1
r999i_mds XmitData: 0.0   1.9    1.3    0.9    0.2    0.1    1.2    2.2    0.3    2.1   0.0  11.2   0.0   0.0   0.0   0.0   0.0   0.0   0.0   0.0   0.8   22.2


r999i_oss1        .     .  r41i3  r49i3  r57i3  r17i3  r25i3 r129i3 r137i3 r145i3 r153i3   r1i3  r9i3  oss2  oss2   mds   mds  oss4  oss4  oss7  oss7 hwsw1    tot
r999i_oss1 RcvData:  0.0   5.2    0.5    0.3    0.8    2.9    4.9    2.0    1.9    5.6 170.4  37.2   0.0   0.0   0.0   0.0   0.0   0.0   0.0   0.0   0.1  231.9
r999i_oss1 XmitData: 0.0   3.3    1.3    0.5    0.8   12.0    1.8    1.7    1.1    0.9   1.9   4.2   0.0   0.0   0.1   0.0   0.0   0.0   0.0   0.0   0.0   29.7


r999i_oss2        .     .  r42i2  r50i2  r58i2  r18i2  r26i2 r130i2 r138i2 r146i2 r154i2   r2i2 r10i2   mds   mds  oss1  oss1  oss5  oss5  oss8  oss8 hwsw2    tot
r999i_oss2 RcvData:  0.0   7.3    0.5    0.3    0.7    2.8    7.6    2.0  115.3    1.8   0.2  46.3   0.0   0.0   0.0   0.0   0.0   0.0   0.1   0.0   0.3  185.1
r999i_oss2 XmitData: 0.0   1.8    1.3    0.5    0.7    0.9    1.8    1.7    2.2    0.9   0.2   1.3   0.1   0.0   0.0   0.0   0.1   0.1   0.0   0.0   0.2   13.6
```

# Lustre Metadata Caching

- Implemented a methodology to keep metadata cached
    - Identify sections of OST where metadata is stored.
        - Inodes, bitmaps, etc.
    - Open the raw block device and read those blocks every 5 minutes.
    - Read Caching Turned off on OSS

- Helps to limit the impacts of certain types or user behavior.
    - Vast improvement on certain operations.

- Thought we could turn off in 2.4, but returned to this after meltdown.

# What Do We Want from a Filesystem?

- Reliable
- Easy to Use
- Performance
- Free

# What Do We Want from a Filesystem?

- Reliable
  - Some things are surprisingly reliable
    - Suspend/lflush/reboot
    - LBUG in OSS doesn't kill everyone
      - Limited evictions
    - Recovery Works (sometimes)
  - Some things not
    - Cascasding failures
      - LBUG or KDB across all servers
      - 1000's of client evictions
    - *Always* hit already patched bugs

# What Do We Want from a Filesystem?

- Easy to Use
  - Generally – Very easy to use (POSIX compliant)
  - Maybe a few odd end cases
    - E.g. partial read or write

# What Do We Want from a Filesystem?

- Performance
  - Can get very good performance

  - Things you need to do to get performance doesn't always map easily to many applications.
    - ECCO

  - Large system
    - I/Os look random once they get to the back end

# What Do We Want from a Filesystem?

- Free
  - Yes – In the Stallman sense.

  - Still require high levels of support
    - Bug tracking/patching - steep curve here

# Issues

- Intel kept two maintenance releases 2.4 and 2.5, then dropped 2.4
- Got on 2.4 early, and then had problems moving to 2.5
- Hit many bugs that were already patched

- Bug tracking jira and Bug patching gerrit system need to talk. Missed some updated patch sets, resulted in more crashes.

# Issues

- Resiliance
  - Cascading failures.
  - Rebooting all 110 lustre servers
  - Commit on Share (help recovery?)

- Quiesce Filesystem for administrative work/upgrades
- Performance
  - Single user can drag down performance
  - Network Request Scheduler (LU-398) is on out list to test

- Single client performance
  -

# Issues

- Quotas seldom work. Moving to the OSTs made them more fragile

- We seem to always hit bugs that are already patched.
    - Over and over again. Since the beginning of time.

# What Does NASA Want from Lustre

- Increased Stability
    - Better Patch Management
- Better Workload Performance (500+ jobs).
- QoS – Limiting Damage of Creative Users
- Administrative Shutdown