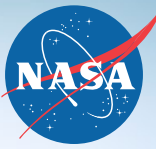


# HPCG Pleiades

[Bob.Ciotti@nasa.gov](mailto:Bob.Ciotti@nasa.gov)

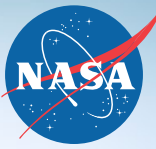
John Baron [jbaron@sgi.com](mailto:jbaron@sgi.com)

NASA Ames Research Center



# Pleiades HW Environment

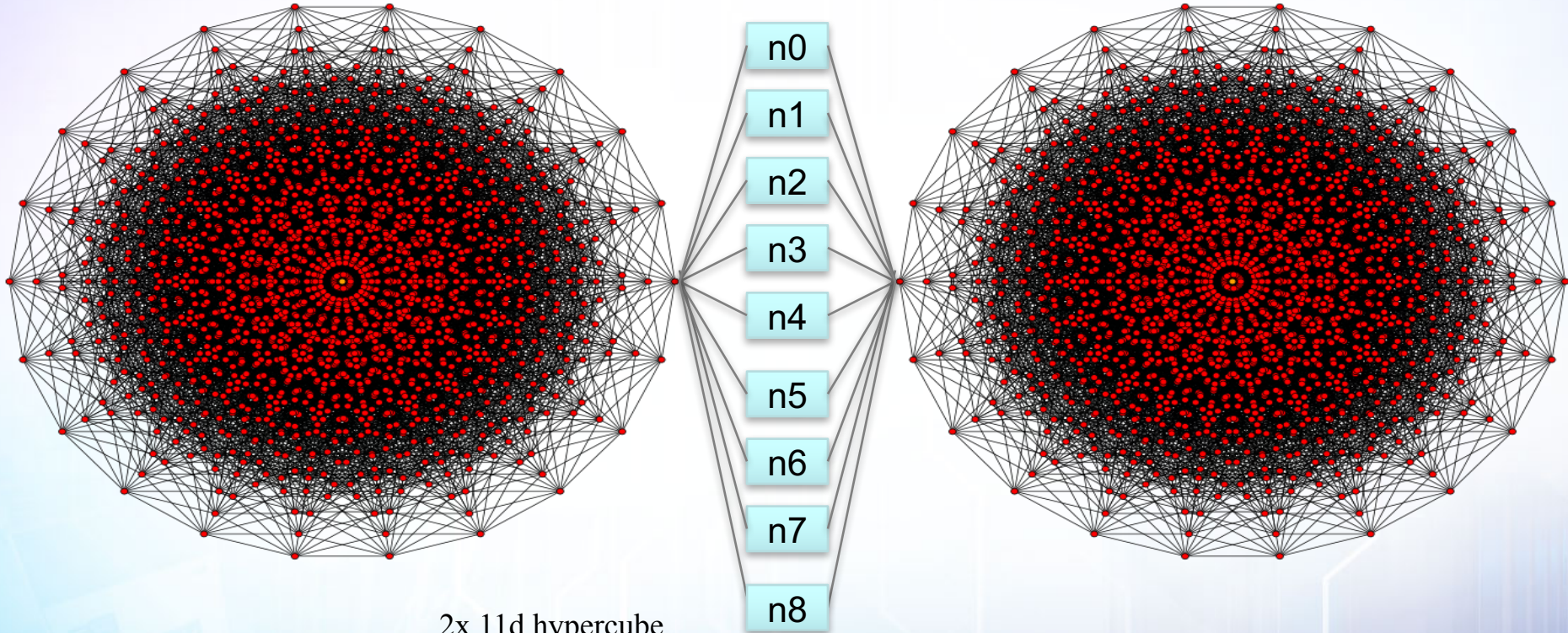
- 11, 472 compute nodes 246,048 x86 cores
  - 1,968 Sandybridge
  - 5,400 Ivybridge
  - 2,088 Haswell
  - 2,016 Broadwell
- 938TB Memory
- FDR Infiniband – dual rail hypercube
- Additional task specific nodes
  - GPU
  - Xeon Phi (KNC+KNL)
  - 1024/512 cpu large shared memory
  - Large memory data analysis nodes
  - Front Ends
  - hyperwall viz/data analysis
- + a couple hundred administration/management nodes of various types.



# Pleiades SW Environment

- LINUX
  - SLES11/12 (most user facing systems)
  - Red Hat/Centos (lustre servers)
- Lustre
- NFS
- Continuous Availability
- \*All\* software can be updated without full system dedicated outage
  - 'rolling updates for compute nodes
  - Suspend/Resume for service nodes (lustre/NFS servers, rack leaders)
- Compute nodes added/removed without dedicated system down

# SGI ICE Dual Plane – Topology



**ib0**

2x 11d hypercube

full 11d == 2048 vertices

Pleiades – partial 11d - 1296 vertices (2592 across both cubes)

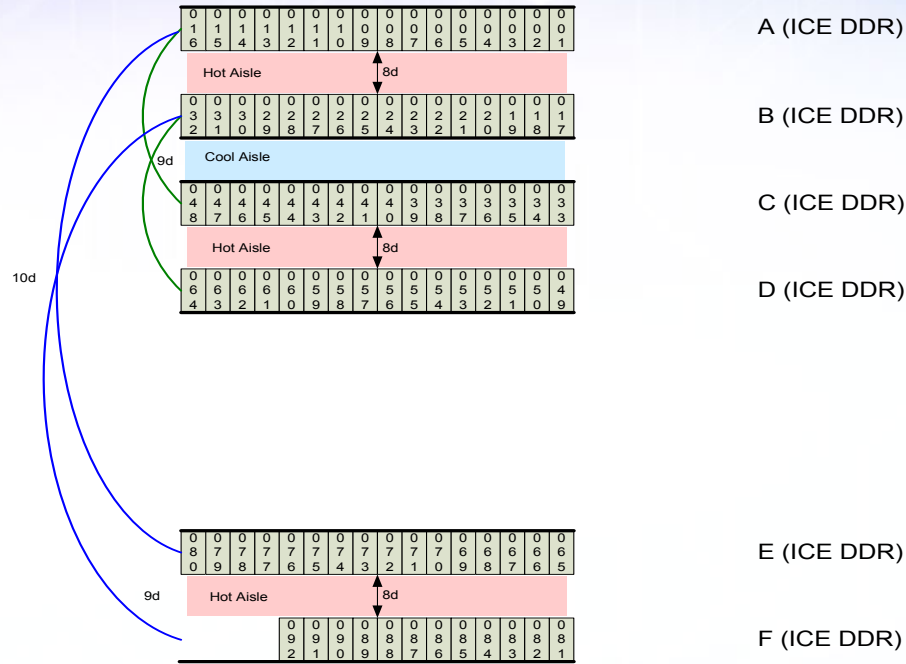
**ib1**

[http://en.wikipedia.org/wiki/User:Qef/Orthographic\\_hypercube\\_diagrams](http://en.wikipedia.org/wiki/User:Qef/Orthographic_hypercube_diagrams)





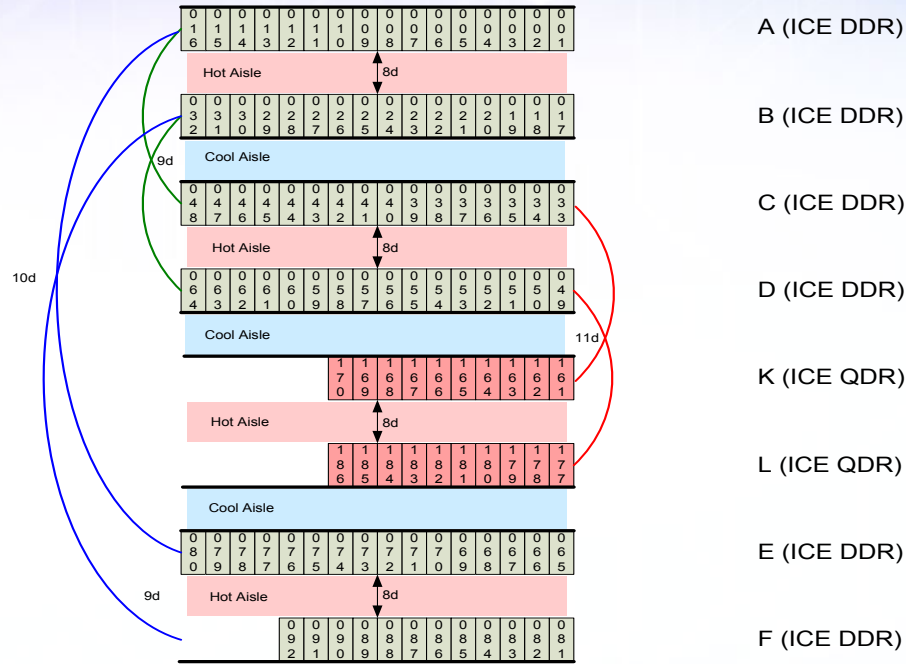
### NASA (Pleiades) Rack Layout



92 racks – 2008  
565 teraflops

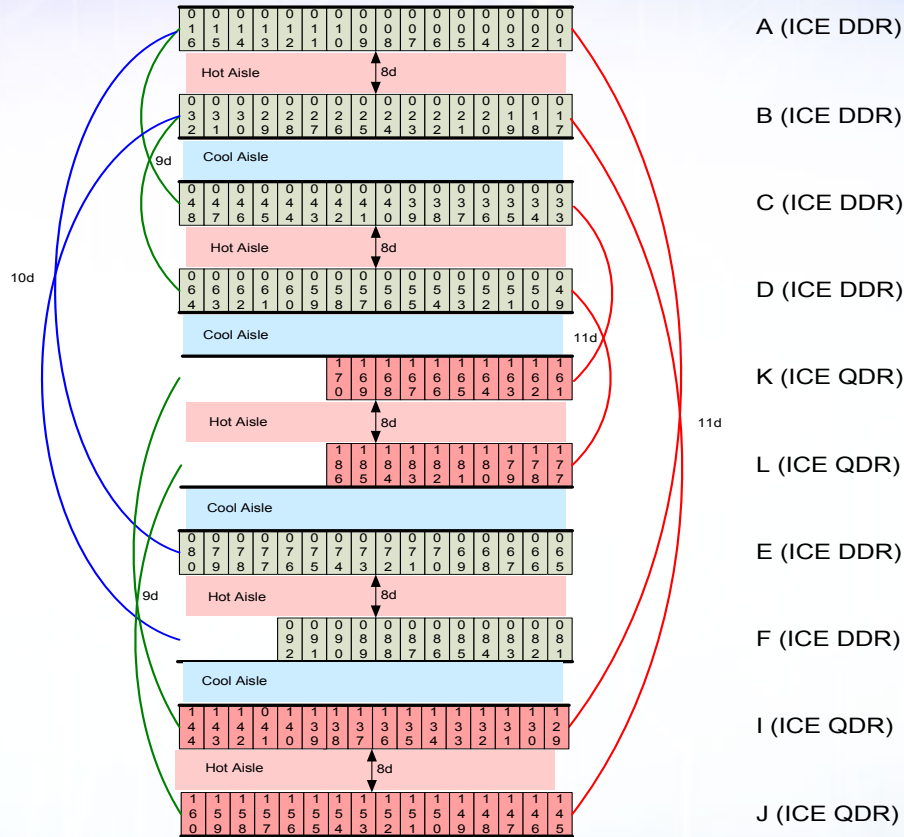
#3 Top500

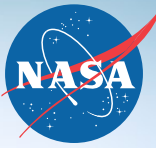
# NASA (Pleiades) Rack Layout



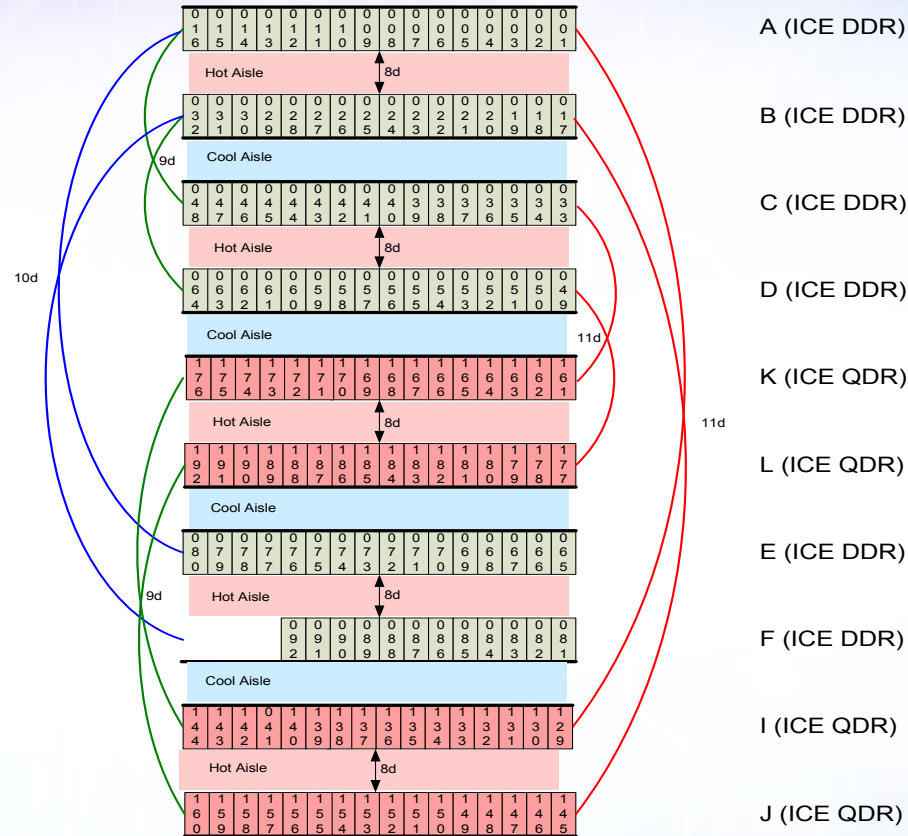
112 racks – 2009  
683 teraflops

# NASA (Pleiades) Rack Layout





# NASA (Pleiades) Rack Layout

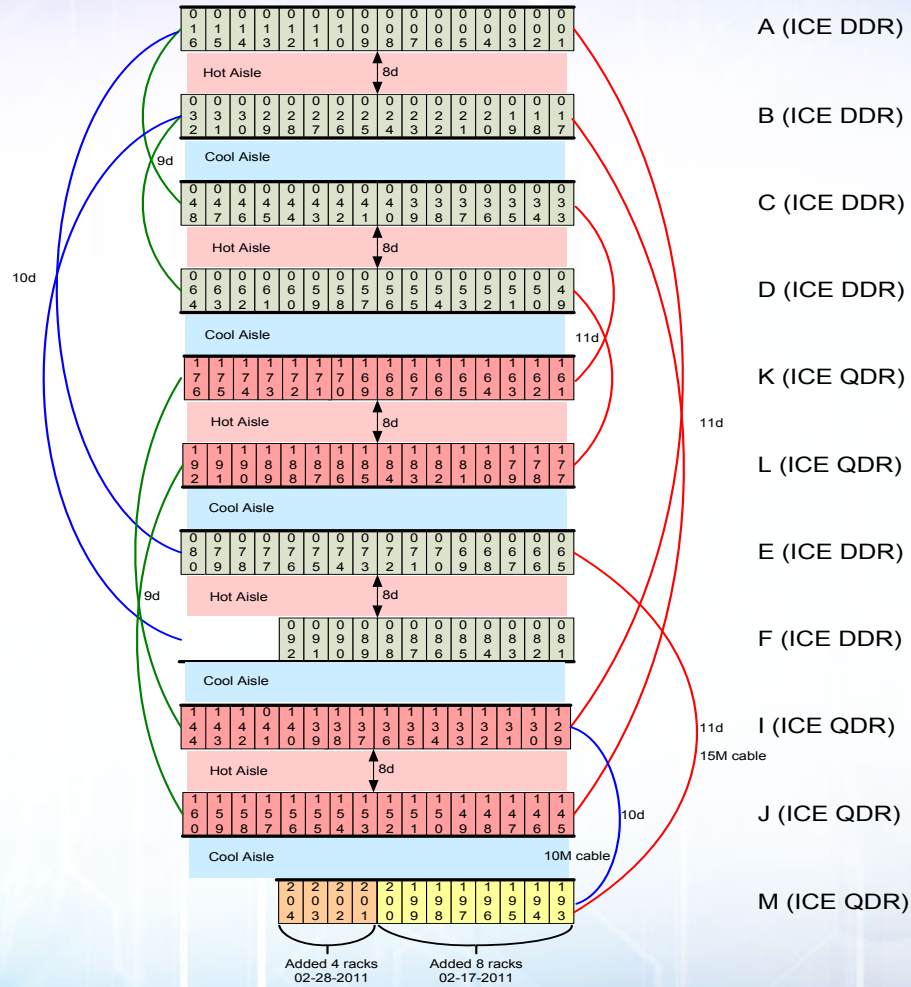


156 racks – 2010  
1.08 petaflops





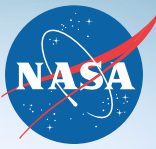
# NASA (Pleiades) Rack Layout



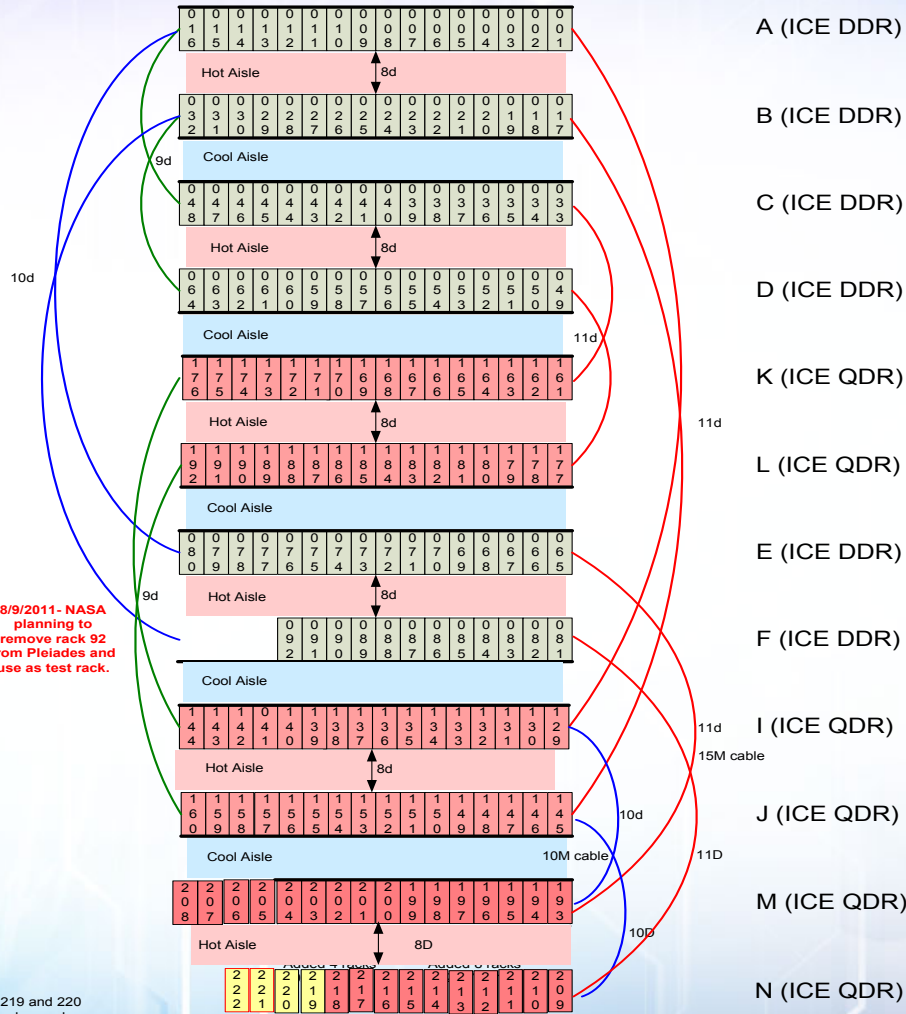
168 racks – 2011  
1.18 petaflops







# NASA (Pleiades) Rack Layout



186 racks – 2011  
1.33 petaflops

8/9/2011- NASA planning to remove rack 92 from Pleiades and use as test rack.

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.

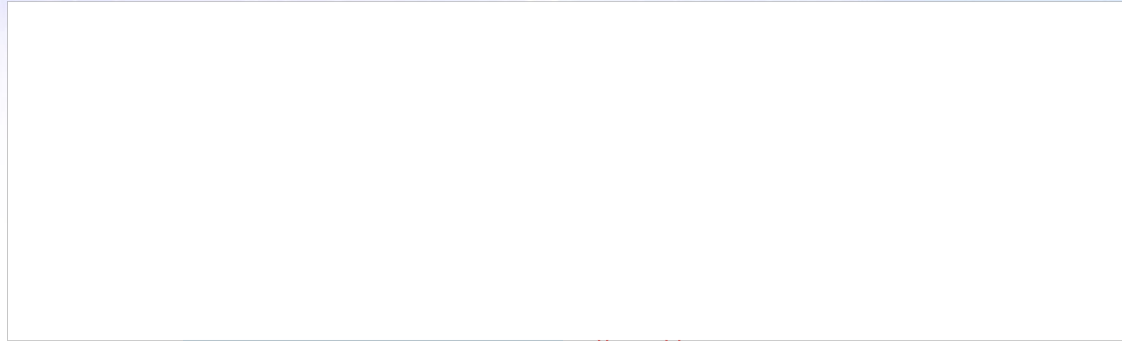
Note: Rack 221 will cable to on 11D to rack 92. There is no 11d for Rack 222, this is a problem. If we remove rack 92 then we have issue with racks 221 & 222.





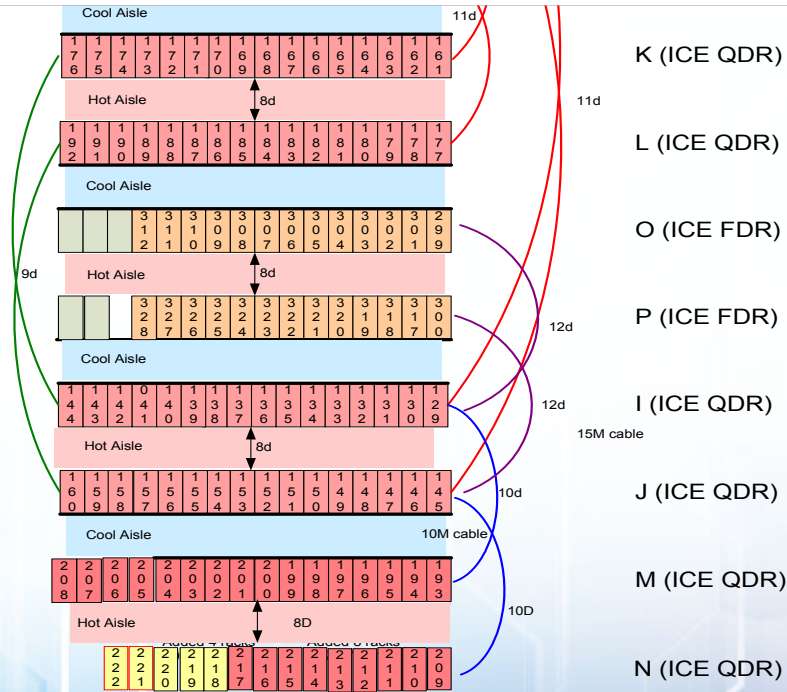


# NASA (Pleiades) Rack Layout



## 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



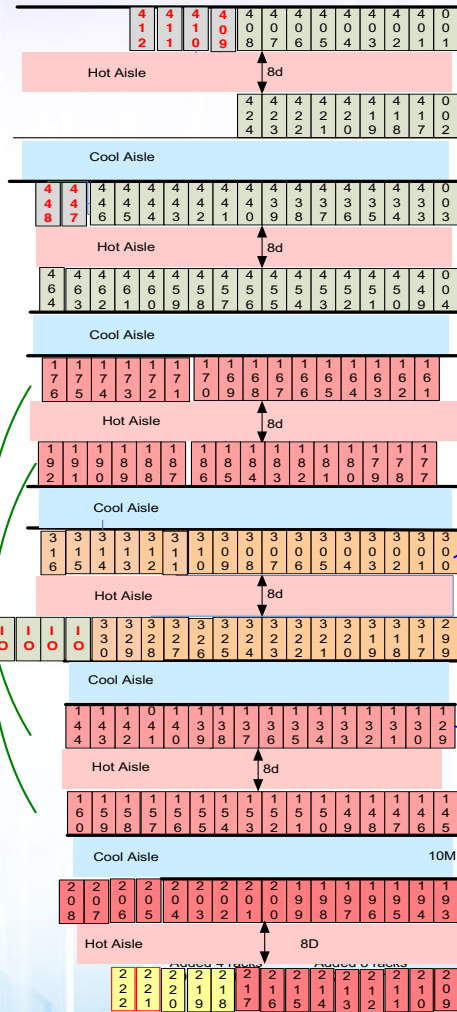
Gpgpu racks 219 and 220 but configured as rack 219, note switches on gpgpu are in rear of rack so cable lengths needs to be adjusted to reflect this.

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 12/30/2013



A (ICE FDR)  
RK 401-416 -SGI ICE X  
(ivybridge) Prem SW

B (ICE FDR)  
RK417-432 SGI ICE X  
(ivybridge) Prem SW

C (ICE FDR)  
RK 433-448 SGI ICE X  
(ivybridge) prem Sw

D (ICE FDR) Rk 449-464  
SGI ICE X  
(ivybridge) Prem SW

K (ICE QDR) rks 161-170  
Altix ICE 8200  
Nehalem + ICE QDR)  
171-176 Altix ICE 8400  
EX (Westmere)

L (ICE QDR) rks 177-186  
Altix ICE 8200  
Nehalem +(ICE QDR)  
187-17192 Altix ICE  
8400 EX (Westmere)

O (ICE FDR) RK 300-312  
SGI ICE X  
Sandybridge Prem SW/  
RK 313-316 128 node  
Pyramid in hypercube  
topology

P (ICE FDR) RK 317-330  
SGI ICE X SNB  
Prem SW

I (ICE QDR) RKS 129-144  
Altix Ice 8400 EX

J (ICE QDR) RKS 145-160  
ALTIX ICE 8400  
EX

M (ICE QDR) RKS 193-208  
ALTIX ICE 8400  
EX

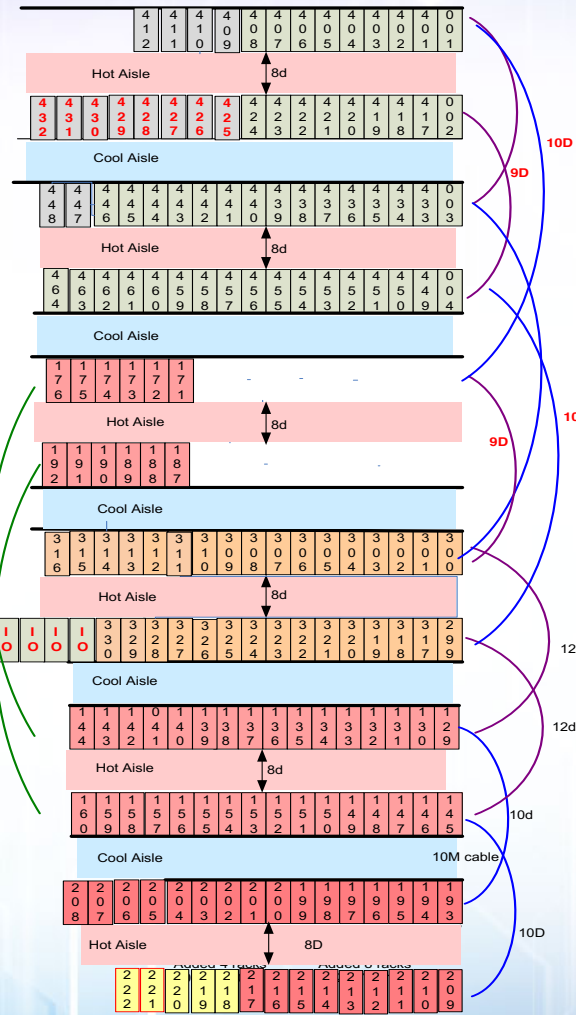
N (ICE QDR) RKS 209-218  
ALTIX ICE 8400 EX/  
RK 219-220 Coyote based  
Westmere with GPGPU  
M2090 in hypercube. RK  
221-222 Altix ICE 8400 EX

160 racks – 2013  
3.1 petaflops

This is the switch rack



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



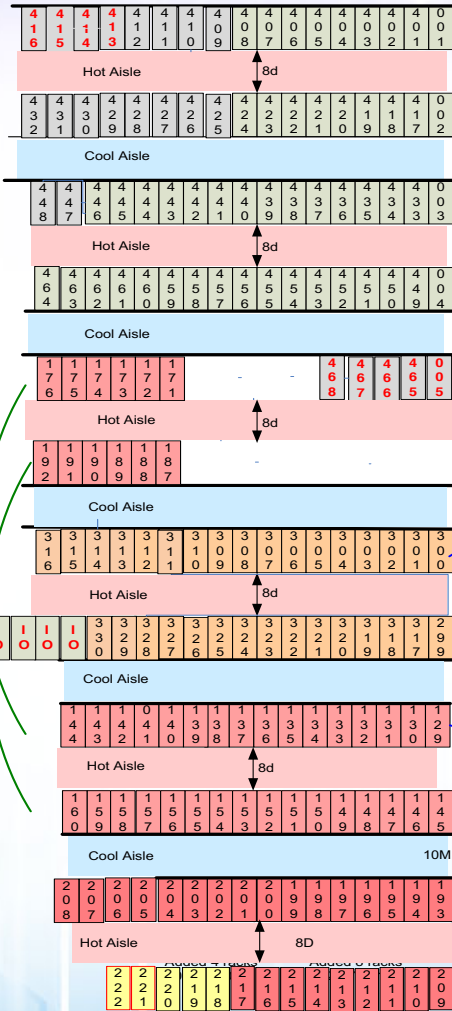
- A (ICE FDR)  
RK 401-412 -SGI ICE X  
(ivybridge) Prem SW
- B (ICE FDR)  
RK417-432 SGI ICE X  
(ivybridge) Prem SW
- C (ICE FDR)  
RK 433-448 SGI ICE X  
(ivybridge) prem Sw
- D (ICE FDR) Rk 449-464  
SGI ICE X  
(ivybridge) Prem SW
- E (ICE FDR) rks 465-468  
SGI ICE X  
Ivybridge Prem SW – K  
ICE QDR) 171-176 Altix  
ICE 8400 EX
- F (ICE FDR) rks 481-483  
SGI ICE X  
Ivybrodge Prem SW L-  
ICE QDR) 187-17192  
Altix ICE 8400 EX
- O (ICE FDR) RK 300-312  
SGI ICE X  
Sandybridge Prem SW/  
RK 313-316 128 node  
Pyramid in hypercube  
topology
- P (ICE FDR) RK 317-330  
SGI ICE X SNB  
Prem SW
- I (ICE QDR) RKS 129-144  
Altix Ice 8400 EX
- J (ICE QDR) RKS 145-160  
ALTIX ICE 8400 EX
- M (ICE QDR) RKS 193-208  
ALTIX ICE 8400 EX
- N (ICE QDR) RKS 209-218  
ALTIX ICE 8400 EX/  
RK 219-220 Coyote based  
Westmere with GPGPU  
M2090 in hypercube. RK  
221-222 Altix ICE 8400 EX

168 racks – 2013  
3.2 petaflops





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

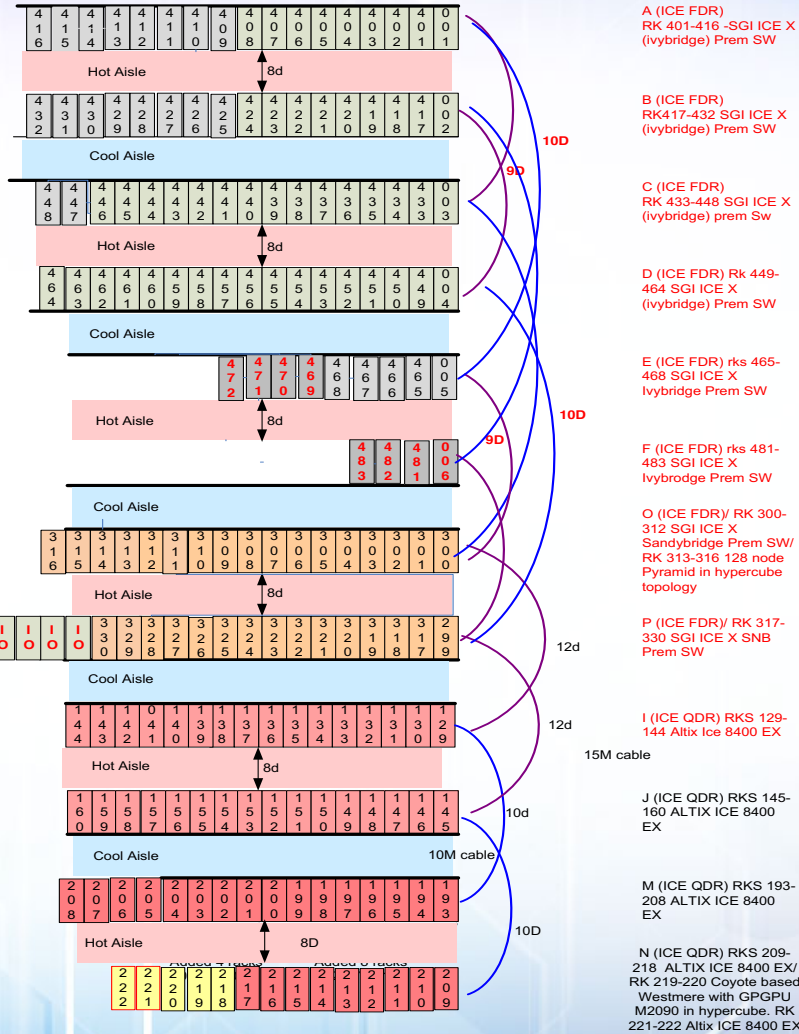


- A (ICE FDR)  
RK 401-416 -SGI ICE X  
(ivybridge) Prem SW
- B (ICE FDR)  
RK 417-432 SGI ICE X  
(ivybridge) Prem SW
- C (ICE FDR)  
RK 433-448 SGI ICE X  
(ivybridge) prem Sw
- D (ICE FDR) Rk 449-464  
SGI ICE X  
(ivybridge) Prem SW
- E (ICE FDR) rks 465-468  
SGI ICE X  
Ivybridge Prem SW – K  
ICE QDR) 171-176 Altix  
ICE 8400 EX
- F (ICE FDR) rks 481-483  
SGI ICE X  
Ivybridge Prem SW L-  
ICE QDR) 187-17192  
Altix ICE 8400 EX
- O (ICE FDR) RK 300-312  
SGI ICE X  
Sandybridge Prem SW/  
RK 313-316 128 node  
Pyramid in hypercube  
topology
- P (ICE FDR) RK 317-330  
SGI ICE X SNB  
Prem SW
- I (ICE QDR) RKS 129-144  
Altix Ice 8400 EX
- J (ICE QDR) RKS 145-160  
ALTIX ICE 8400 EX
- M (ICE QDR) RKS 193-208  
ALTIX ICE 8400 EX
- N (ICE QDR) RKS 209-218  
ALTIX ICE 8400 EX/  
RK 219-220 Coyote based  
Westmere with GPGPU  
M2090 in hypercube. RK  
221-222 Altix ICE 8400 EX

168 racks – 2014  
3.3 petaflops



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

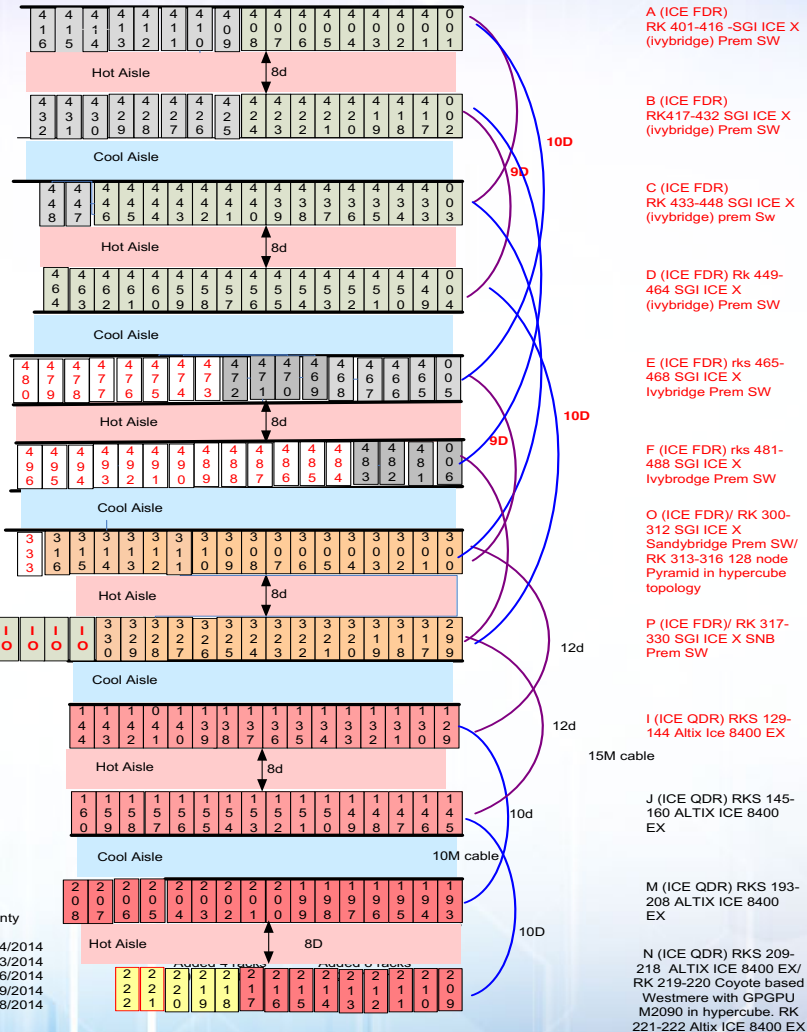


170 racks – 2014  
3.5 petaflops

This is the switch rack



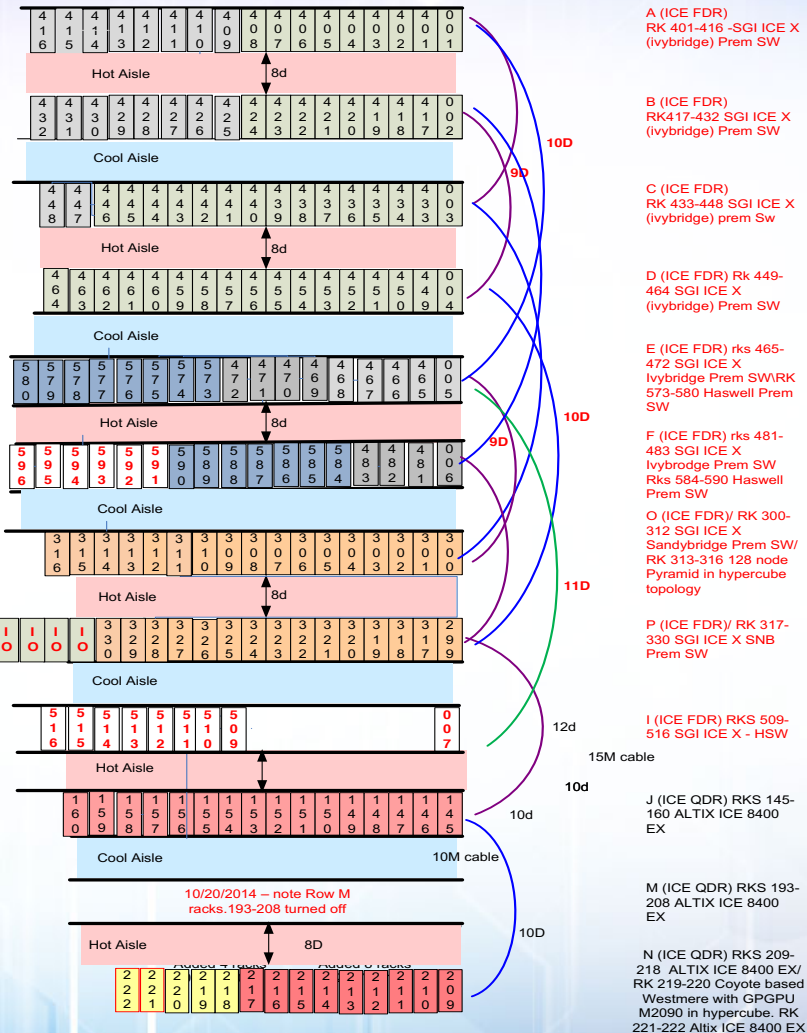
# NASA (Pleiades) Rack Layout



168 racks – 2014  
4.5 petaflops



# NASA (Pleiades) Rack Layout



168 racks – 2015  
5.4 petaflops

This is the switch rack

11/02/2014 – remove 16 racks of WSM – rks 129-144 and replace with racks 509-516 haswell

10/20/2014 – note Row M racks.193-208 turned off



# NASA (Pleiades) Rack Layout

4	4	4	4	4	4	4	4	4	4	4	4	4	4	0
1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
6	5	4	3	2	1	0	9	8	7	6	5	4	3	2

Hot Aisle

8d

4	4	4	4	4	4	4	4	4	4	4	4	4	4	0
3	3	3	2	2	2	2	2	2	2	2	2	1	1	1
2	1	0	9	8	7	6	5	4	3	2	1	0	9	8

Cool Aisle

4	4	4	4	4	4	4	4	4	4	4	4	4	4	0
4	4	4	4	4	4	4	4	4	3	3	3	3	3	3
8	7	6	5	4	3	2	1	0	9	8	7	6	5	4

Hot Aisle

8d

4	4	4	4	4	4	4	4	4	4	4	4	4	4	0
6	6	6	6	6	5	5	5	5	5	5	5	5	5	4
4	3	2	1	0	9	8	7	6	5	4	3	2	1	0

Cool Aisle

5	5	5	5	5	5	5	5	4	4	4	4	4	4	0
8	7	7	7	7	7	7	7	7	6	6	6	6	6	0
0	9	8	7	6	5	4	3	2	1	0	9	8	7	6

Hot Aisle

8d

5	5	5	5	5	5	5	5	5	5	5	5	4	4	0
9	9	9	9	9	9	8	8	8	8	8	8	8	8	0
6	5	4	3	2	1	0	9	8	7	6	5	4	3	2

Cool Aisle

3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
1	1	1	1	1	1	1	0	0	0	0	0	0	0	0
6	5	4	3	2	1	0	9	8	7	6	5	4	3	2

Hot Aisle

8d

This is the switch rack

I	O	I	O	I	O	I	O	I	O	I	O	I	O	I
0	9	8	7	6	5	4	3	2	1	0	9	8	7	6

Cool Aisle

11/02/2014 -- remove 16 racks of WSM -- rks 129-144 and replace with racks 509-516 haswell

5	5	5	5	5	5	5	5	6	6	6	6	6	6	0
1	1	1	1	1	1	1	0	0	0	0	0	0	0	0
6	5	4	3	2	1	0	9	8	7	6	5	4	3	2

Hot Aisle

8d

1	1	1	1	1	1	1	1	X	6	6	6	6	6	0
6	5	5	5	5	5	5	5	X	2	2	2	2	1	0
0	9	8	7	6	5	4	3	2	1	0	9	8	7	6

Cool Aisle

10M cable

10/20/2014 -- note Row M racks. 193-208 turned off

Hot Aisle

8d

2	2	2	2	2	2	2	2	X	X	X	X	X	X	X
2	2	2	1	1	1	1	1	X	X	X	X	X	X	X
2	1	0	9	8	7	6	5	X	X	X	X	X	X	X

A (ICE FDR) RK 401-416 -SGI ICE X (ivybridge) Prem SW

B (ICE FDR) RK417-432 SGI ICE X (ivybridge) Prem SW

C (ICE FDR) RK 433-448 SGI ICE X (ivybridge) Prem SW

D (ICE FDR) Rk 449-464 SGI ICE X (ivybridge) Prem SW

E (ICE FDR) rks 465-472 SGI ICE X Ivybridge Prem SW/RK 573-580 Haswell Prem SW

F (ICE FDR) rks 481-483 SGI ICE X Ivybridge Prem SW Rks 584-590 Haswell Prem SW

O (ICE FDR) RK 300-312 SGI ICE X Sandybridge Prem SW/ RK 313-316 128 node Pyramid in hypercube topology

P (ICE FDR) RK 317-330 SGI ICE X SNB Prem SW

I (ICE FDR) RKS 509-516 SGI ICE X -- HSW Racks 603-608 -BrdWL

J (ICE QDR) RKS 145-160 ALTIX ICE 8400 EX Turn off Rks 145-148 for RKS 603-608 Turn off RKS 149-152 for 601-602/617-622

M (ICE QDR) RKS 193-208 ALTIX ICE 8400 EX

N (ICE QDR) RKS 209-218 ALTIX ICE 8400 EX/ RK 219-220 Coyote based Westmere with GPGPU M2090 in hypercube. RK 221-222 Altix ICE 8400 EX Turn off 209-212 for Rks 603-608 Turn off 213-216 for RKS 601-602/617-620

162 racks – 2016  
7.1 petaflops

#13 Top500 Nov  
#9 HPCG ISC 16

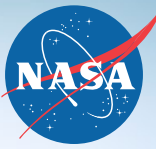
>20 major upgrades





# Highlights of SGI Optimized HPCG Code

- Lexicographical ordering for maximum data locality
- Left and right data structures for full matrix representation
- A variant of CSR storage format
- Pure MPI
- No overlapping of computation and communication
- Additional tuning for contiguous memory, setup time and combined computations



# Heterogeneous considerations

Load balancing via number of ranks per node

Broadwell E5-2680 v4 14-core 2.4 GHz

- 2015 nodes, 12 ranks/socket, **0.85 GF/rank**

Haswell E5-2680 v3 12-core 2.5 GHz

- 2080 nodes, 10 ranks/socket, **0.91 GF/rank**

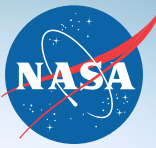
Ivy Bridge E5-2680 v2 10-core 2.8 GHz

- 5351 nodes, 9 ranks/socket, **0.86 GF/rank**

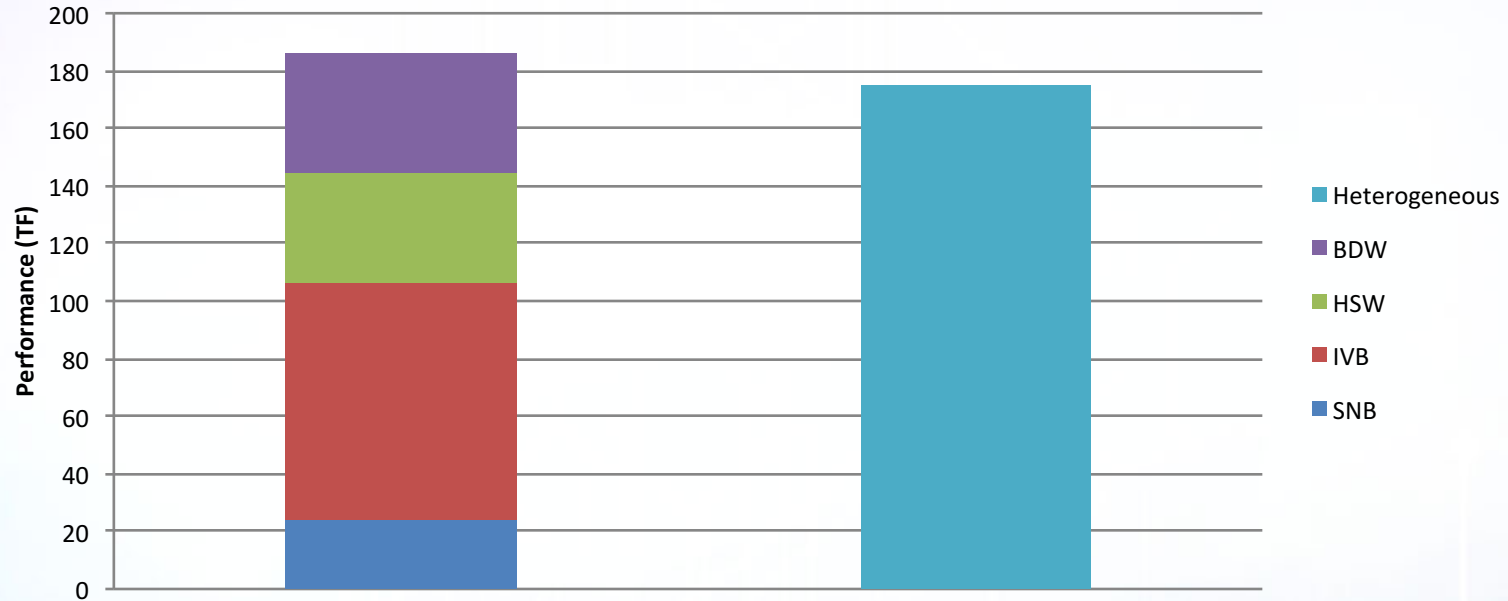
Sandy Bridge E5-2670 8-core 2.6 GHz

- 1853 nodes, 7 ranks/socket, **0.92 GF/rank**

# Pleiades results

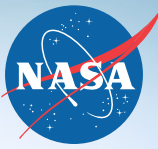


NASA Pleiades HPCG Performance



Aggregate performance is over 94% of the sum of the individual component results

# Credits to the Team



John Baron SGI

Cheng Laio SGI

Michael Raymond SGI

Jay Lan SGI

Scott Emery SGI

Jennifer Fung SGI

Jose Rodriguez SGI

Matt Lepp SGI

Jason Inoue SGI

Rich Davila SGI

John Dugan SGI

Davin Chan CSRA

Dale Talcott CSRA

Jim Karella CSRA

Greg Matthews CSRA

Herbert Yeung CSRA

Mahmoud Hanafi CSRA

Mike Hartman CSRA

Jeff Becker CSRA

Nathan Dauchy CSRA

Bill Thigpen NASA

Mark Tangney NASA

Bob Ciotti NASA