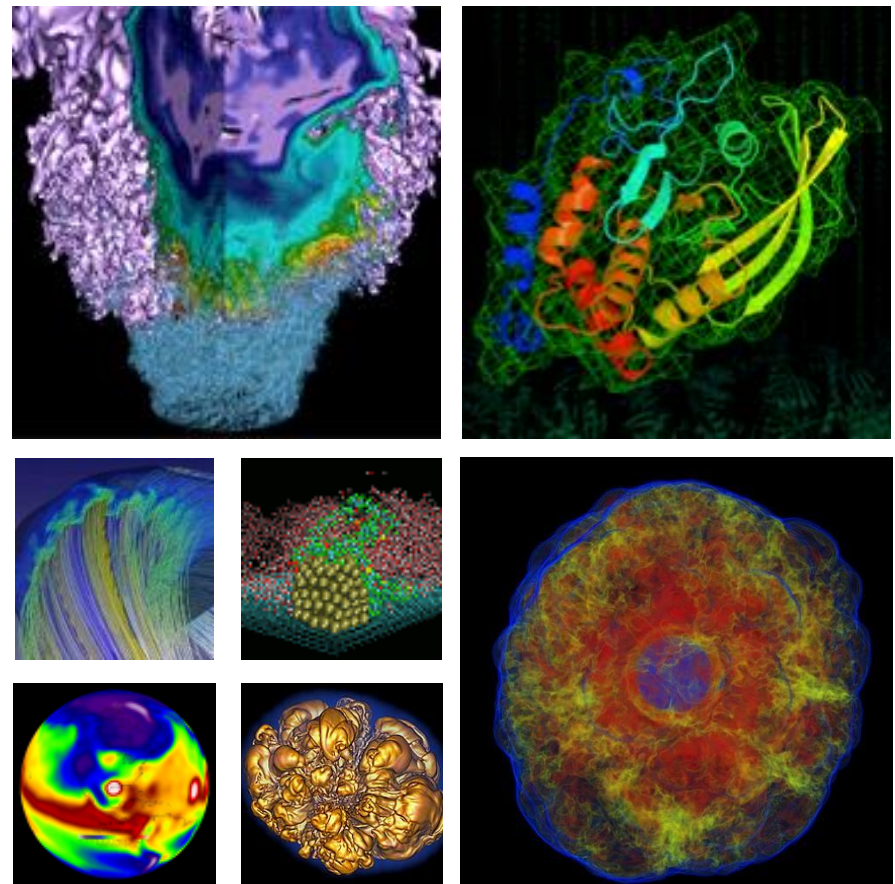


# Performance Variability on Xeon Phi



Brandon Cook, Thorsten Kurth,  
Brian Austin, Samuel Williams,  
Jack Deslippe

June 22, 2017

# Who cares?



- **Application developers**
  - understanding performance
  - reason effectively about optimizations
  - sound advice to application users
- **Users**
  - Efficient use of CPU allocations
  - Wasted cycles on terminated jobs
  - Correct estimates of campaign costs
- **Facilities**
  - System health
  - Advice for users
  - Utilization – scheduler efficiency

# Cori at NERSC



- **2388 Haswell**
  - 2x 16 core @ 2.3 GHz
  - 40 MB shared L3
  - 128 GB DDR
- **Cray Aries Interconnect**
  - dragonfly topology
- **9688 Xeon Phi (KNL) nodes**
  - 68 cores @ 1.4 GHz
  - 34 MB distributed L2
  - 96 GB DDR
  - 16 GB MCDRAM (on-package)



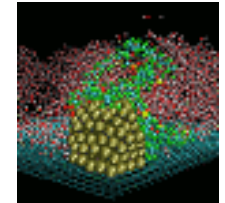
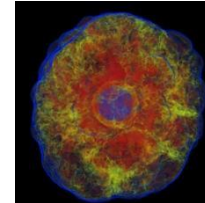
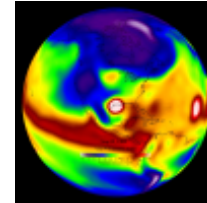
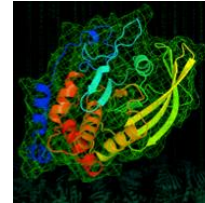
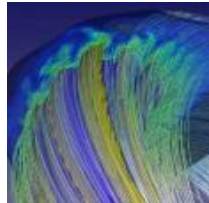
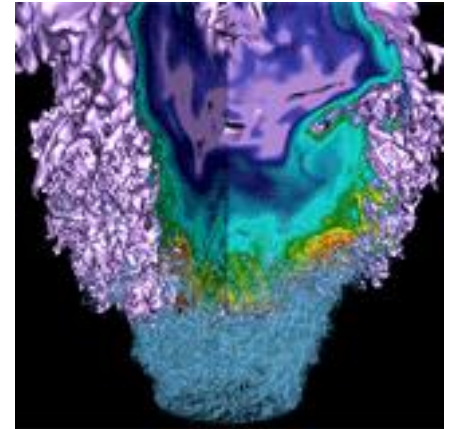
# Cori at NERSC



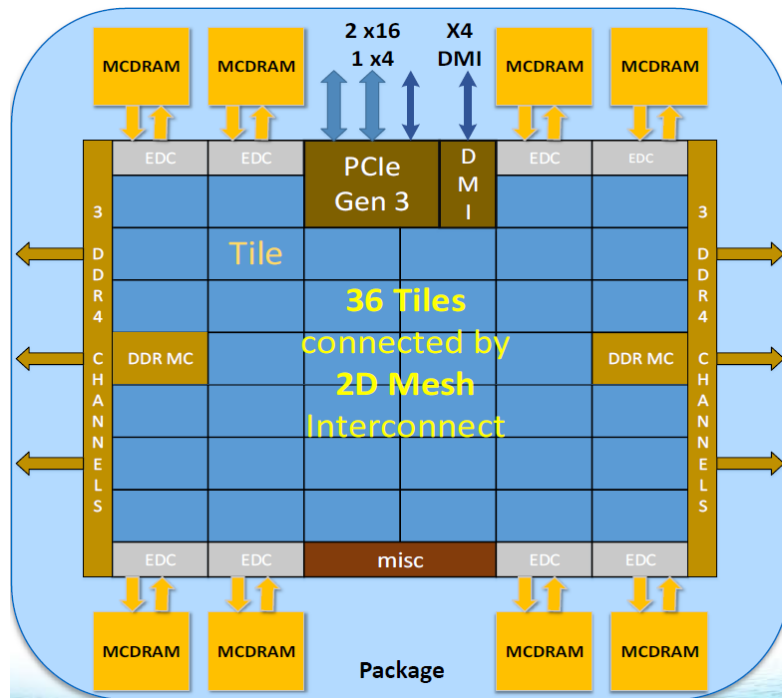
- **2388 Haswell**
  - 2x 16 core @ 2.3 GHz
  - 40 MB shared L3
  - 128 GB DDR
- **Cray Aries Interconnect**
  - dragonfly topology
- **9688 Xeon Phi (KNL) nodes**
  - 68 cores @ 1.4 GHz
  - 34 MB distributed L2
  - 96 GB DDR
  - **16 GB MCDRAM (on-package)**



# MCDRAM



# KNL is highly configurable



## Cluster modes

- all-to-all
- quadrant
- SNC2/4

## Memory modes

- flat
- cache
- hybrid



# MCDRAM cache mode



- 16GB MCDRAM cache
- single NUMA
- No code modification
- No NUMA programming or affinity issues (e.g. numactl)
- but?

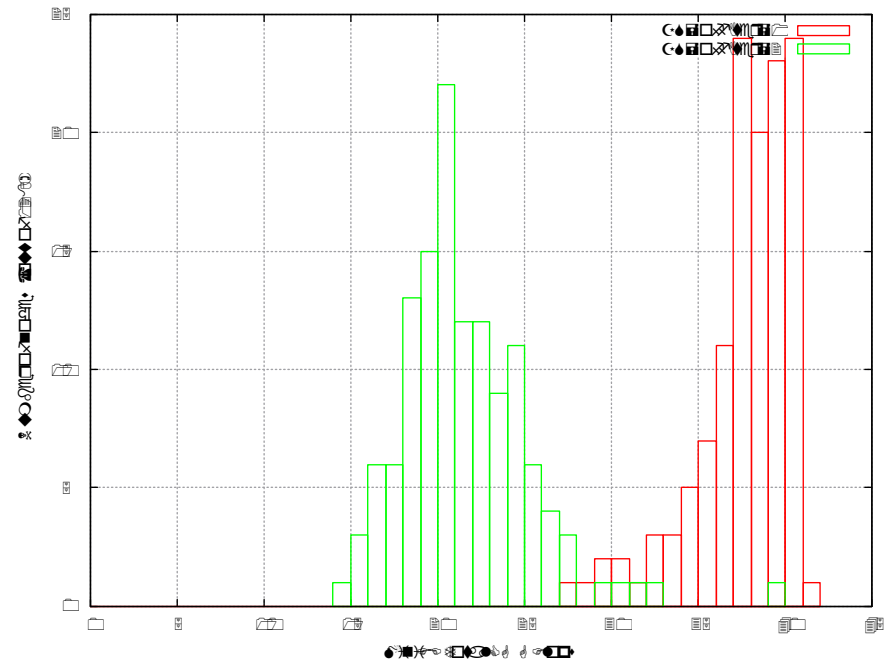
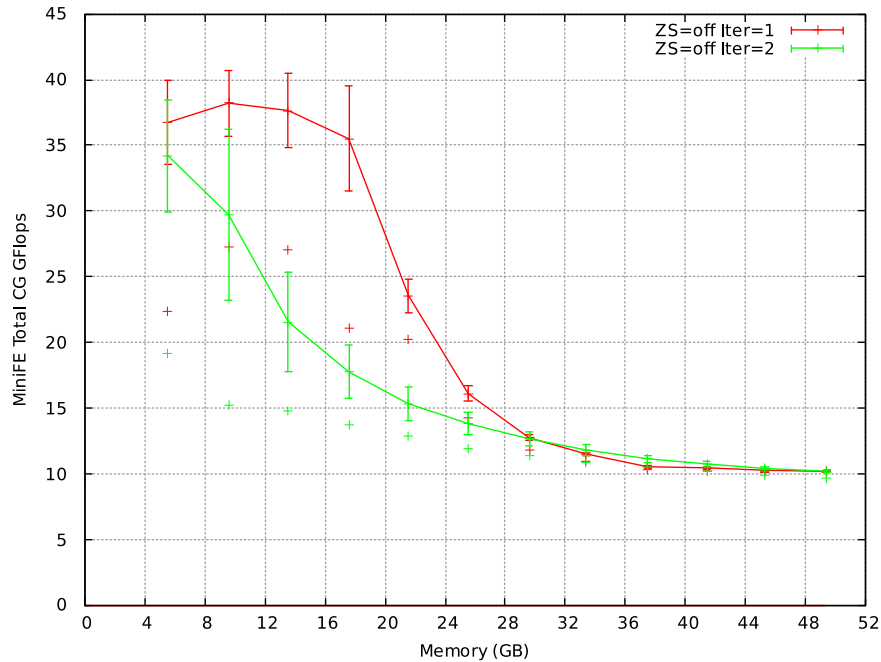
## Cluster modes

- all-to-all
- **quadrant**
- SNC2/4

## Memory modes

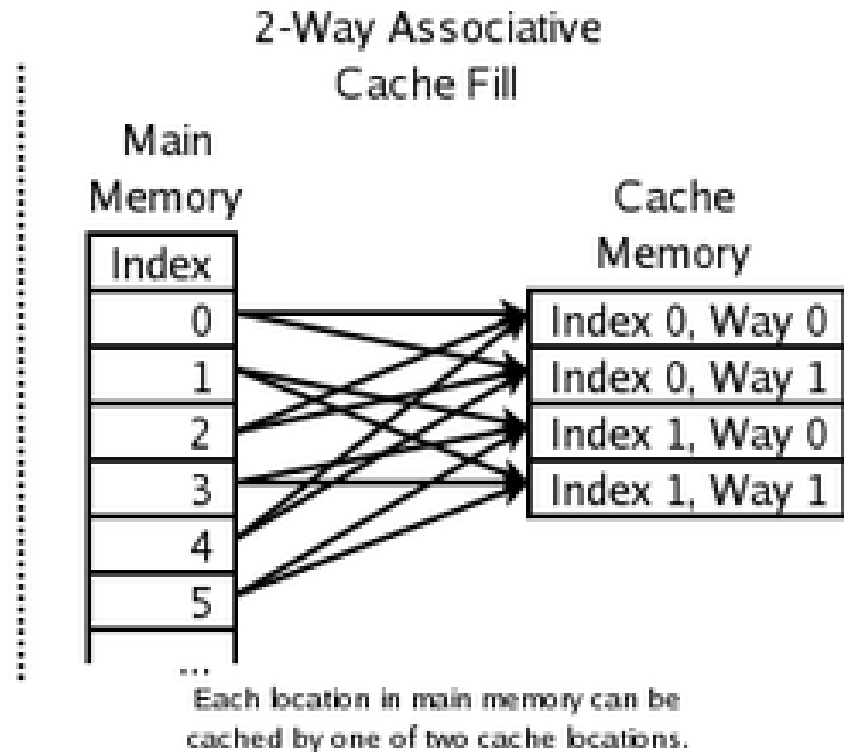
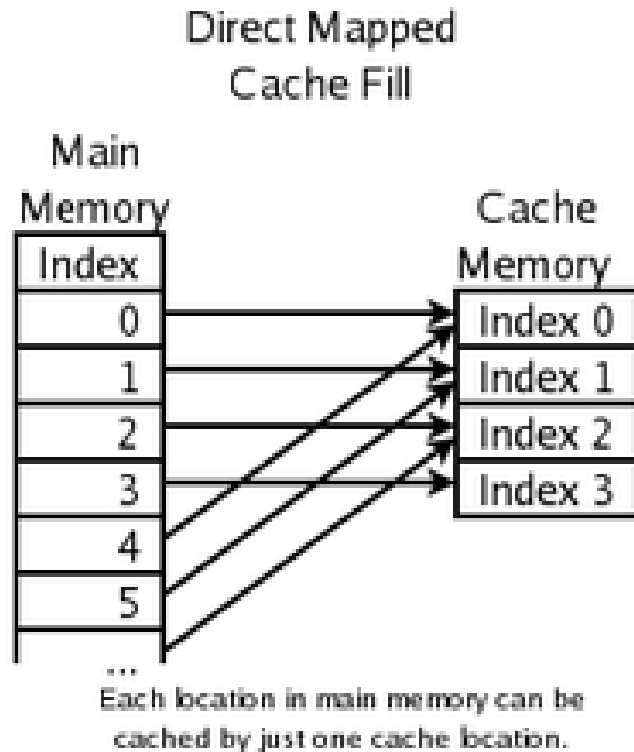
- flat
- **cache**
- hybrid

# Variability in cache mode





# Brief introduction to caches



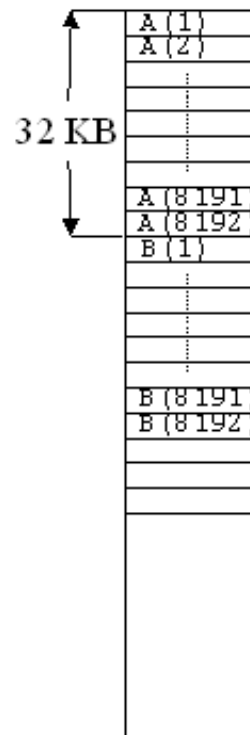
[https://en.wikipedia.org/wiki/CPU\\_cache](https://en.wikipedia.org/wiki/CPU_cache)

KNL's MCDRAM cache is direct-mapped.

# Direct-Mapped Caches: Thrashing

```
COMMON //A(8192), B(8192)
DO I=1,N
    PROD = PROD + A(I)*B(I)
ENDDO
```

(Virtual)  
memory



*Direct mapped* cache (32 KB)

Cache line: 4 words			
A(1)	A(2)	A(3)	A(4)
A(5)	A(6)	A(7)	A(8)
A(8185)	A(8186)	A(8187)	A(8188)
A(8189)	A(8190)	A(8191)	A(8192)

Registers  
in the CPU

*Thrashing: every memory reference results in a cache miss*

Location in the cache:  
 $(\text{memory-address}) \bmod (\text{cache-size})$   
 in this case  $\text{loc}(A(1)) \bmod 32\text{KB} = \text{loc}(B(1)) \bmod 32\text{KB}$   
 [because  $B(1) = A(1) + 8192$ ;  $8192 \cdot 4 \bmod 32\text{KB} = 0$ ]

<http://sc.tamu.edu/help/power/powerlearn/html/ScalarOptnw/sld015.htm>

# Misses depend on free page list



- OS stores a list of free memory pages.
- Allocations are made from the top of the list.
- The free page list gets scrambled if memory is not freed in the order it was allocated.
- Consider a 16kB cache...

//Initial page list is sorted

A = malloc(16 kB) //4 pages

B = malloc( 4 kB) //Capacity conflicts

free(A)

free(B) //Scrambled!

A = malloc( 4 kB)

B = malloc( 4 kB) //Conflict!

OS View							
Free page list							
0	1	2	3	4	5	6	7
				4	5	6	7
					5	6	7
	0	1	2	3	5	6	7
4	0	1	2	3	5	6	7
	0	1	2	3	5	6	7
		1	2	3	5	6	7

Application View							
A				B			
0	1	2	3				
0	1	2	3	4			
				4			
4							
4				0			

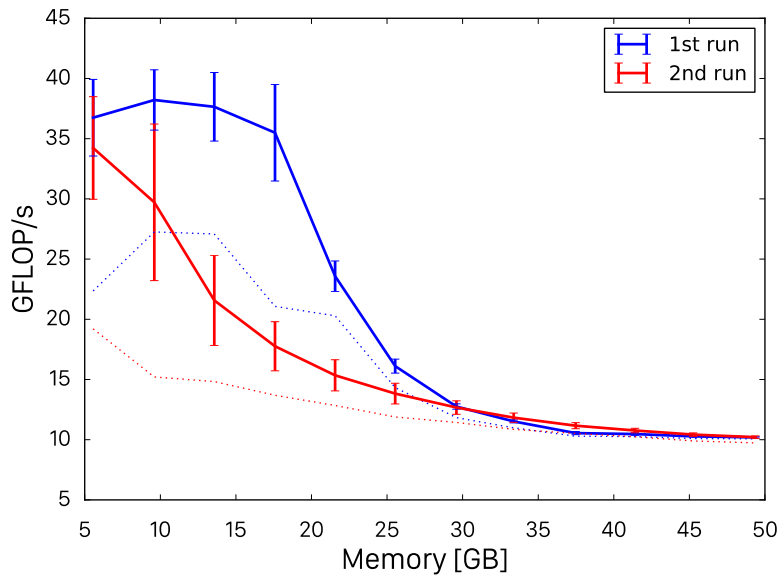
## **Solution: sort the free page list\***

- **zonesort: kernel module provided by Intel**
- **At NERSC**
  - **called immediately before application launch**

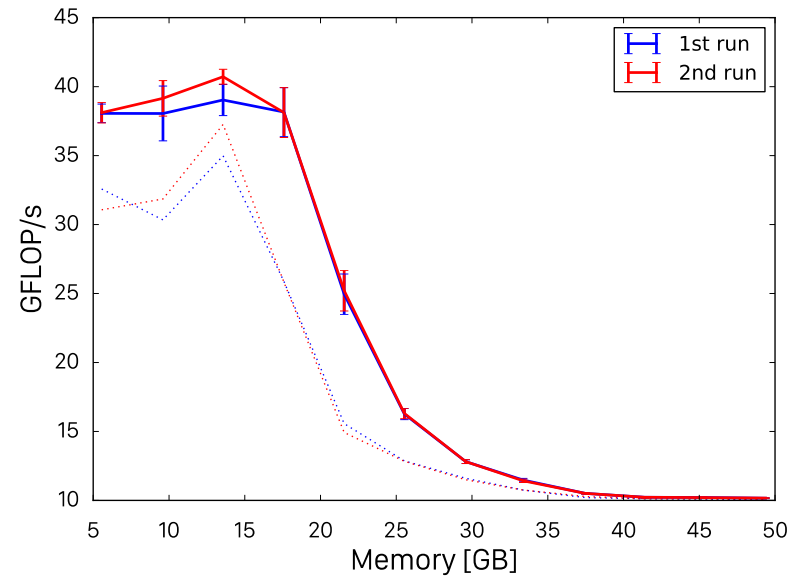
# benefit of zonesort for MiniFE



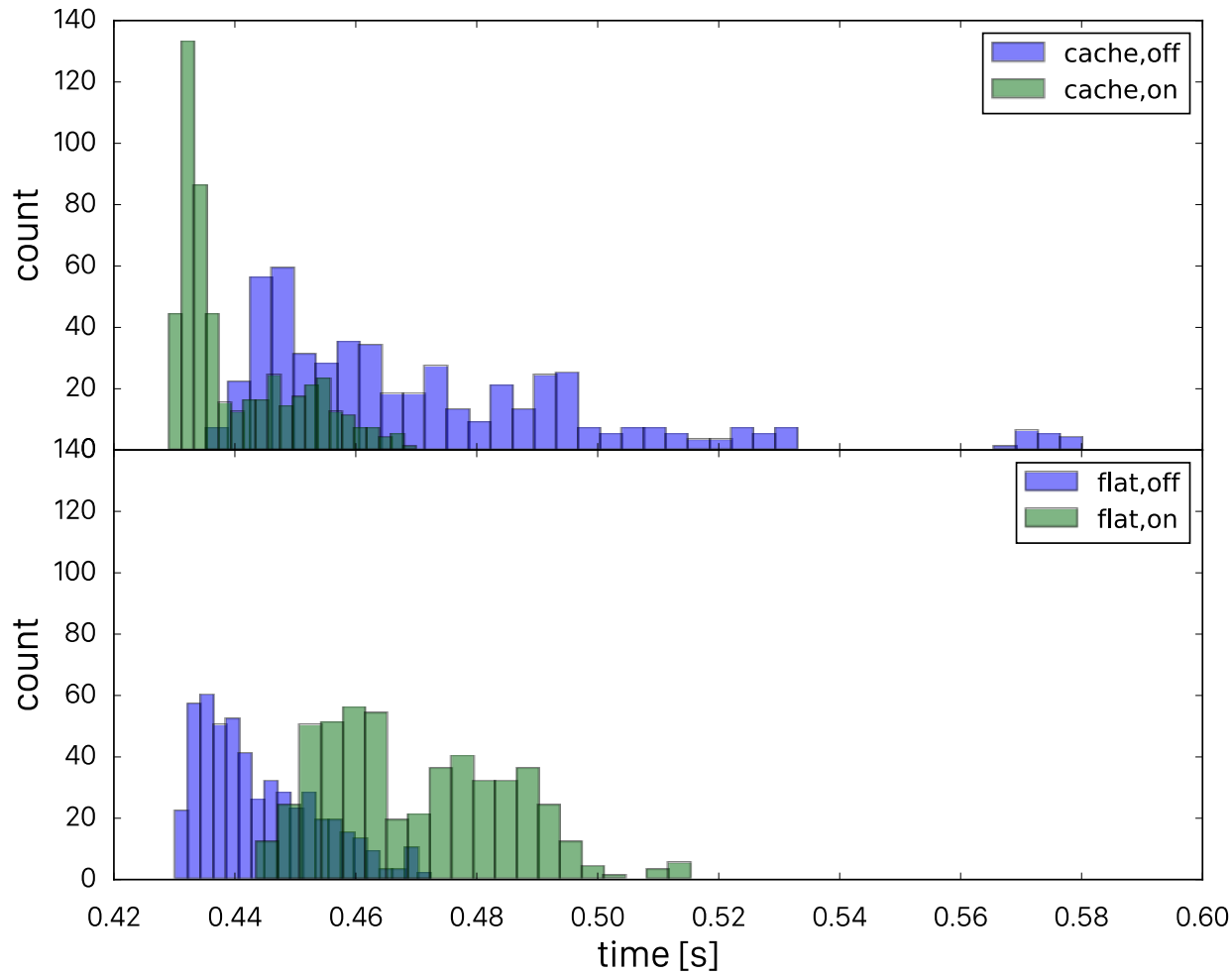
## zonesort off



## zonesort on



# effect of zonesort for HPGMG



High Performance  
Geometric Multi-Grid

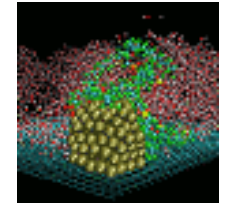
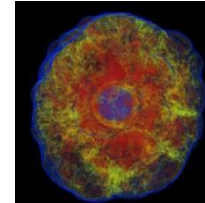
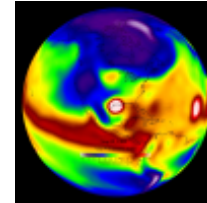
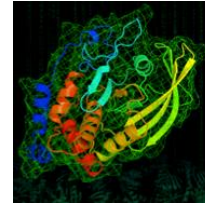
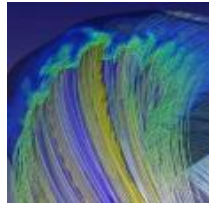
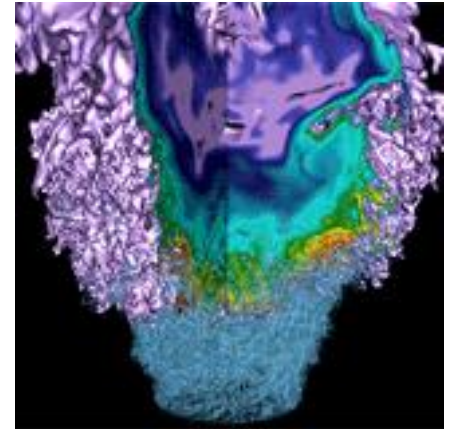
Highly instrumented

Perfectly load balanced  
problem

“smooth time”

$256^3$  grid per rank

# Job placement





# A Xeon Phi issue?

---



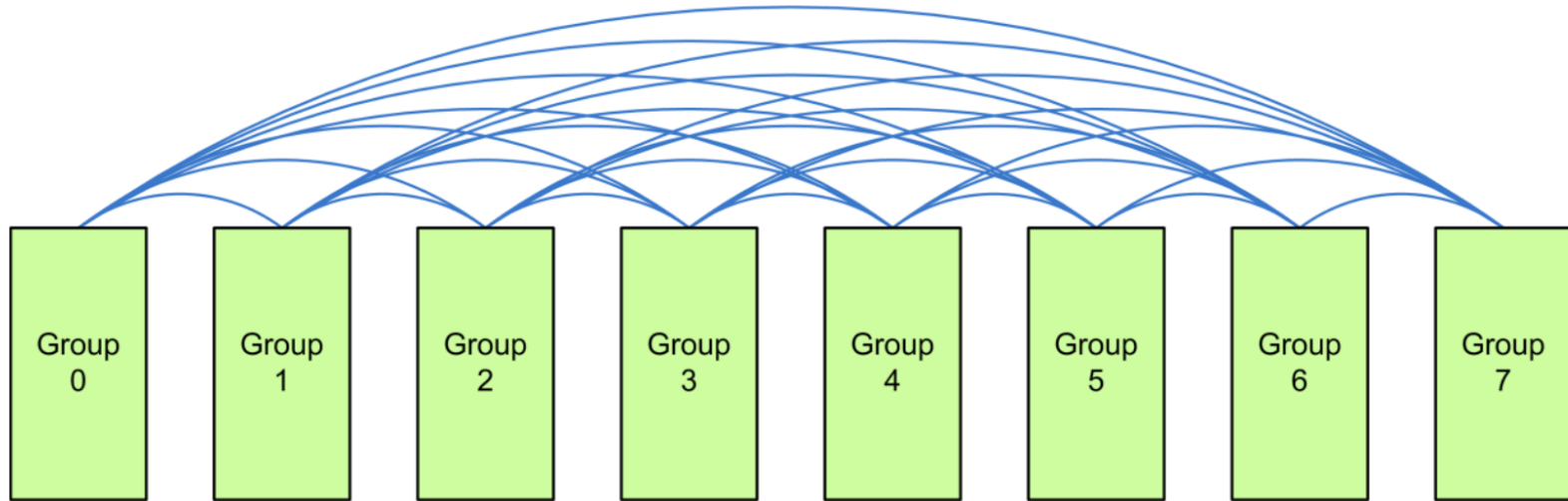
# A Xeon Phi issue?



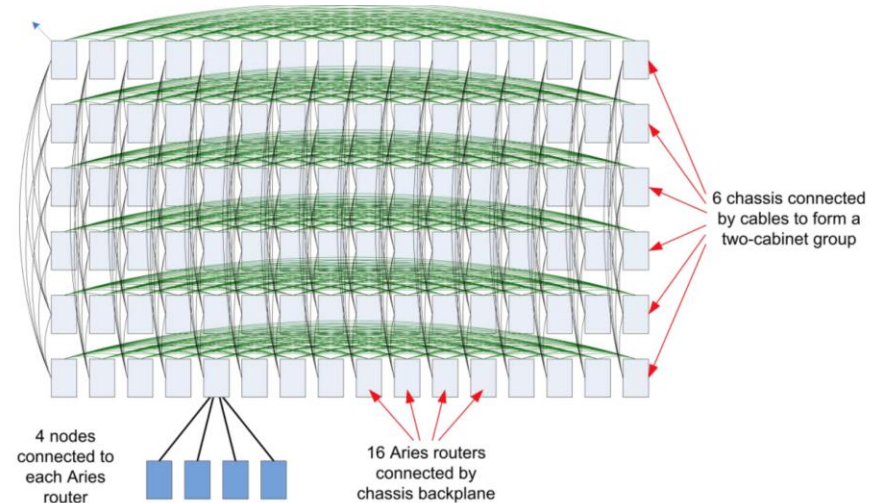
	Haswell	Xeon Phi
<b>Flops</b>	1.2 TFlops	3 TFlops
<b>Memory Bandwidth</b>	~100 GB/s	~400 GB/s
<b>Memory Capacity</b>	128 GB	96 GB
<b>Capacity / bandwidth</b>	1.28 s	0.24 s

**More flops & lower memory capacity / bandwidth  
+  
same network  
=  
more pressure on network!**

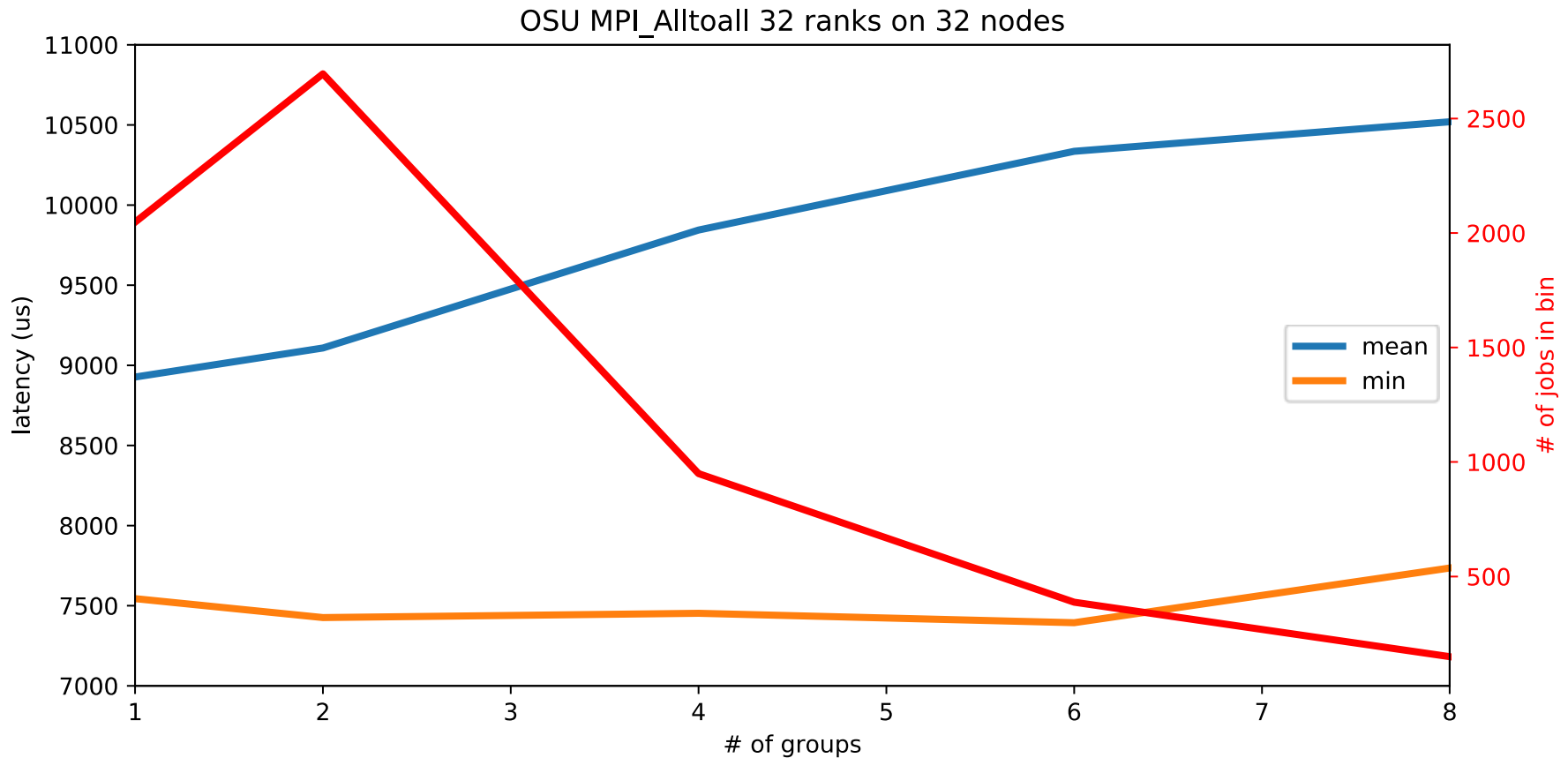
# Aries topology



**~386 nodes per group**



# Impact of # of groups



# scheduler topology awareness



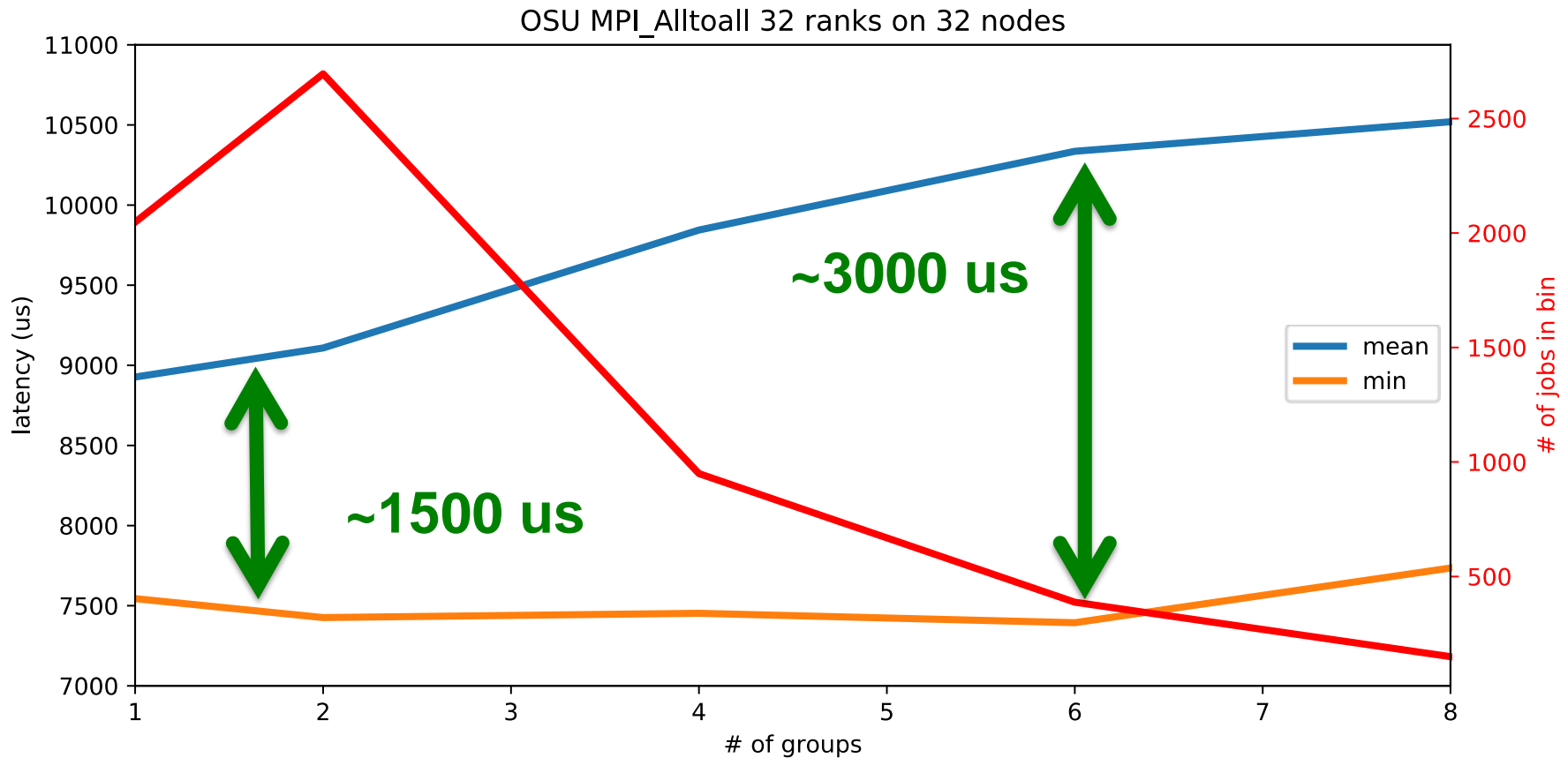
```
sbatch --switches=<count>[@<max-time>]
```



<count> = # of groups

<max-time> = time to wait for  
constraint

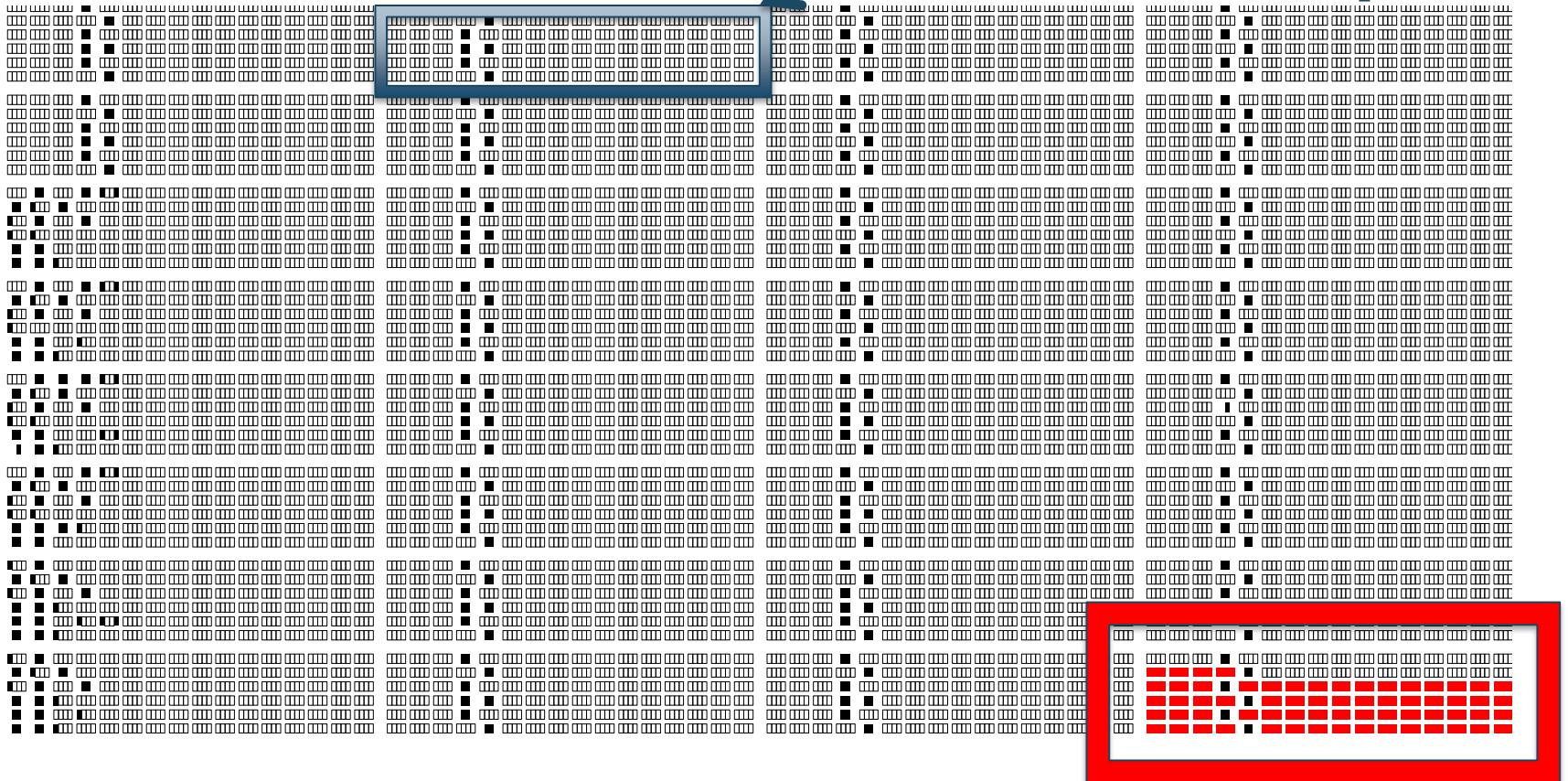
# Impact of # of groups



# Job Placement



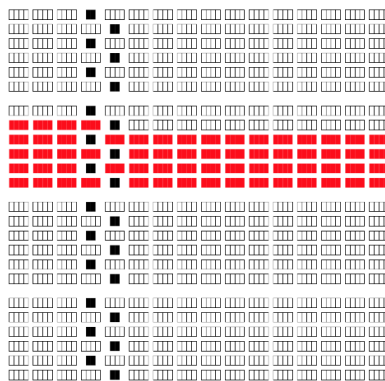
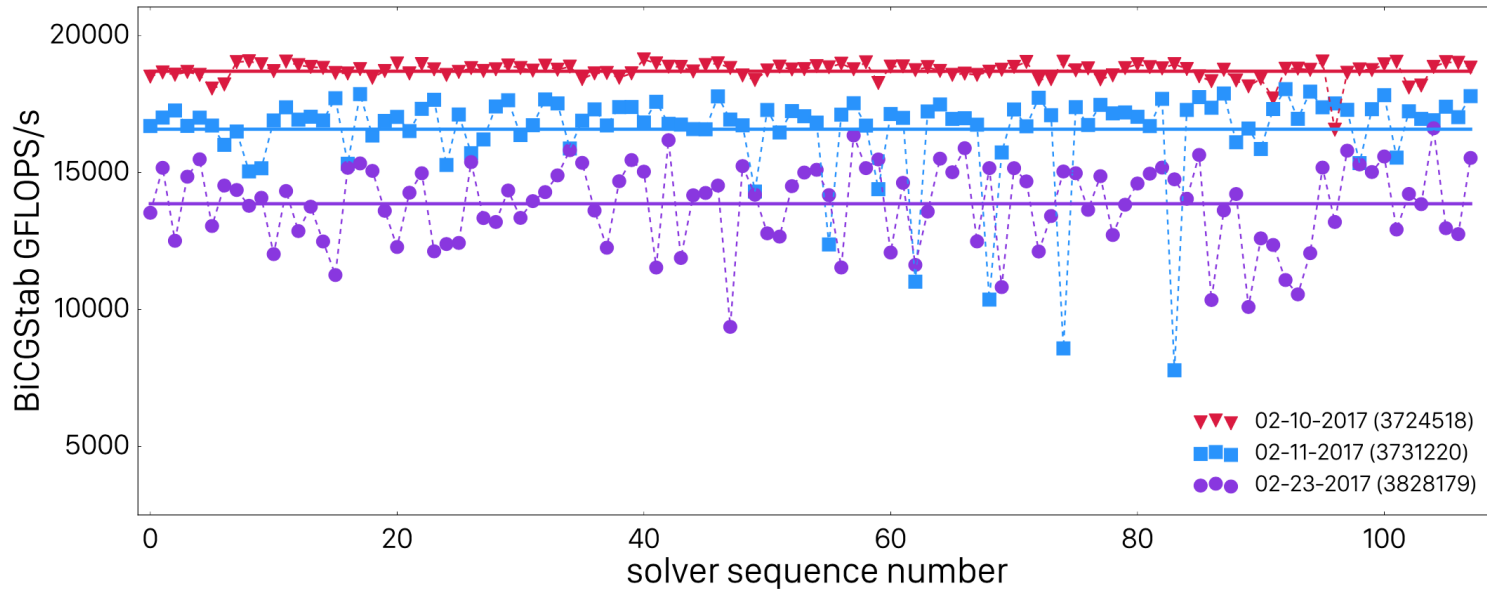
1 Aries Group



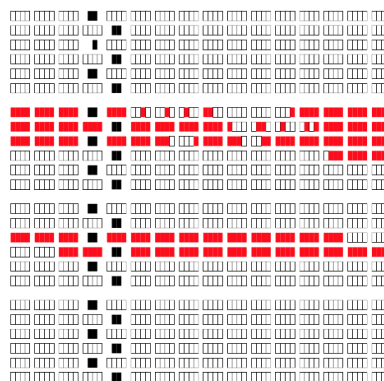
nodes allocated to job



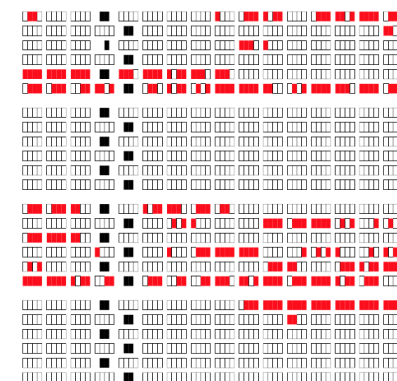
# Chroma HMC – 256 nodes



(b) job 3724518



(c) job 3731220



(d) job 3828179

- **MCDRAM Cache**
  - direct map cache
  - leads to cache conflicts
  - Intel zonesort
- **Job placement**
  - Bad placement introduces extra hops for data
  - Bad placement increases potential for interference
  - SLURM topology control helps (# of nodes < 350)
- **Not covered in this talk**
  - **IO! (burst buffer on compute fabric)**
  - Identification of network “Aggressors”
  - frequency scaling (DVFS)



# National Energy Research Scientific Computing Center