

# The Data Accelerator

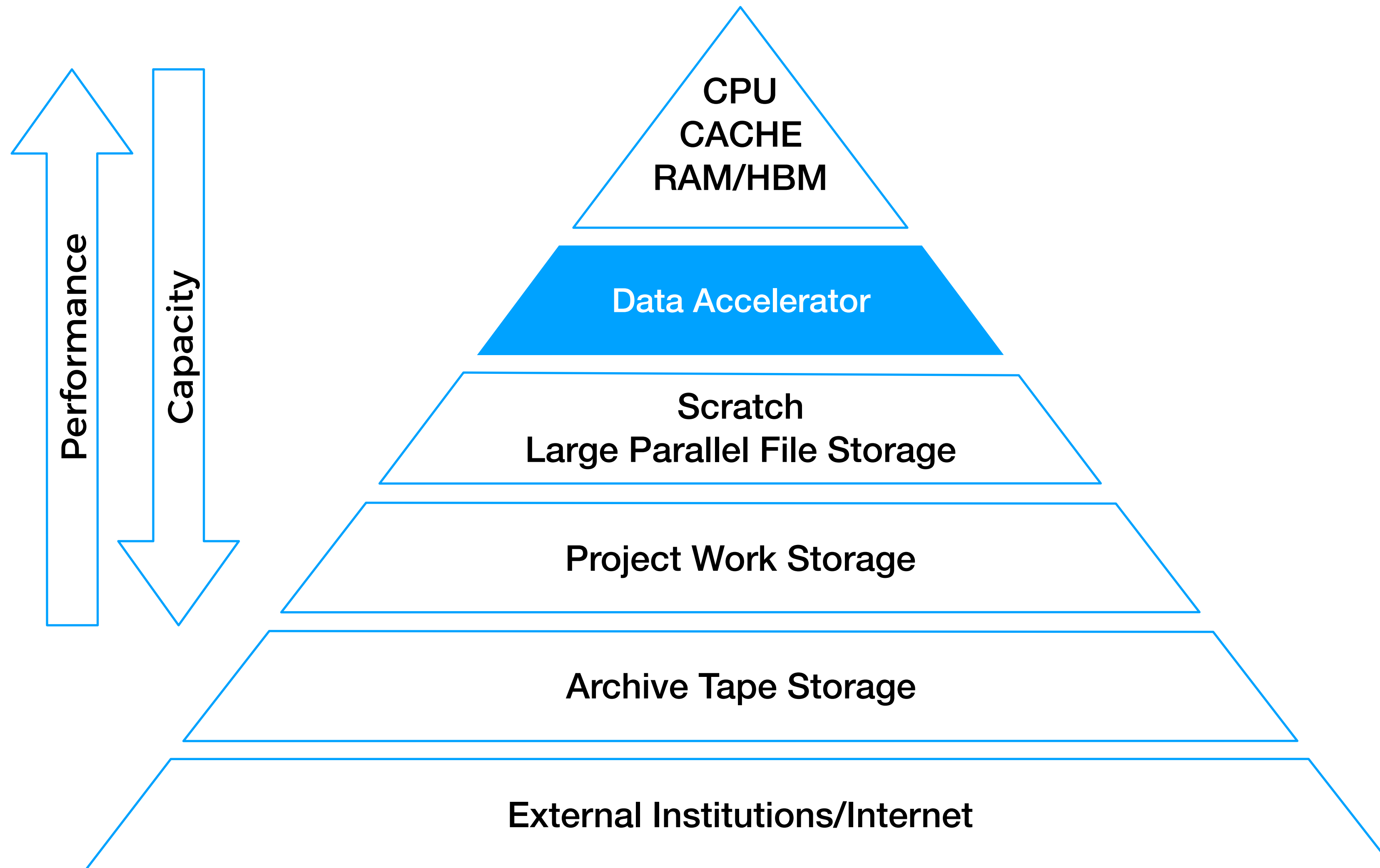
Lustre Administrator and Developer Workshop LAD '18  
Alasdair King [ajk203@cam.ac.uk](mailto:ajk203@cam.ac.uk)



UNIVERSITY OF  
CAMBRIDGE  
Research Computing Services

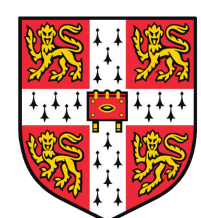








# Cambridge Lustre Estate



UNIVERSITY OF  
CAMBRIDGE  
Research Computing Services



# Lustre - Filesystem Hardware

Type	Quantity	Specs
MDS Server	2	Dell R630 - Dual E5-2667v3 3.2GHz 8 Core 128GB RAM FDR IB, 10GB Ethernet
MDT Storage	1	Dell PowerVault MD3420 20x 300GB SAS 15K HDDs Dual RAID controller with 8GB cache
OSS Server	6	Dell R630 - Dual E5-2623v3 3.0GHz 4 Core 64GB RAM FDR IB, 10GB Ethernet
OSS Storage	6	Dell PowerVault MD3460 60x 6TB NL-SAS HDDs Dual RAID controller with 8GB cache
Network	Infiniband	Mellanox SX6036 FDR IB switches 36x 56Gb FDR
	Ethernet	Mellanox SN2410 Spectrum switches 48x 10GbE + 8x 100GbE

\*Newer Lustre systems use Intel Omni-Path

## Dell MD3420

MDT - 20 disk Raid 10 with 4x SSD cache device



## Dell MD3460

Each chassis supports 6x OSTs

Each OST: 10-disk Raid-6, ~58TB per OST

~350TB per chassis





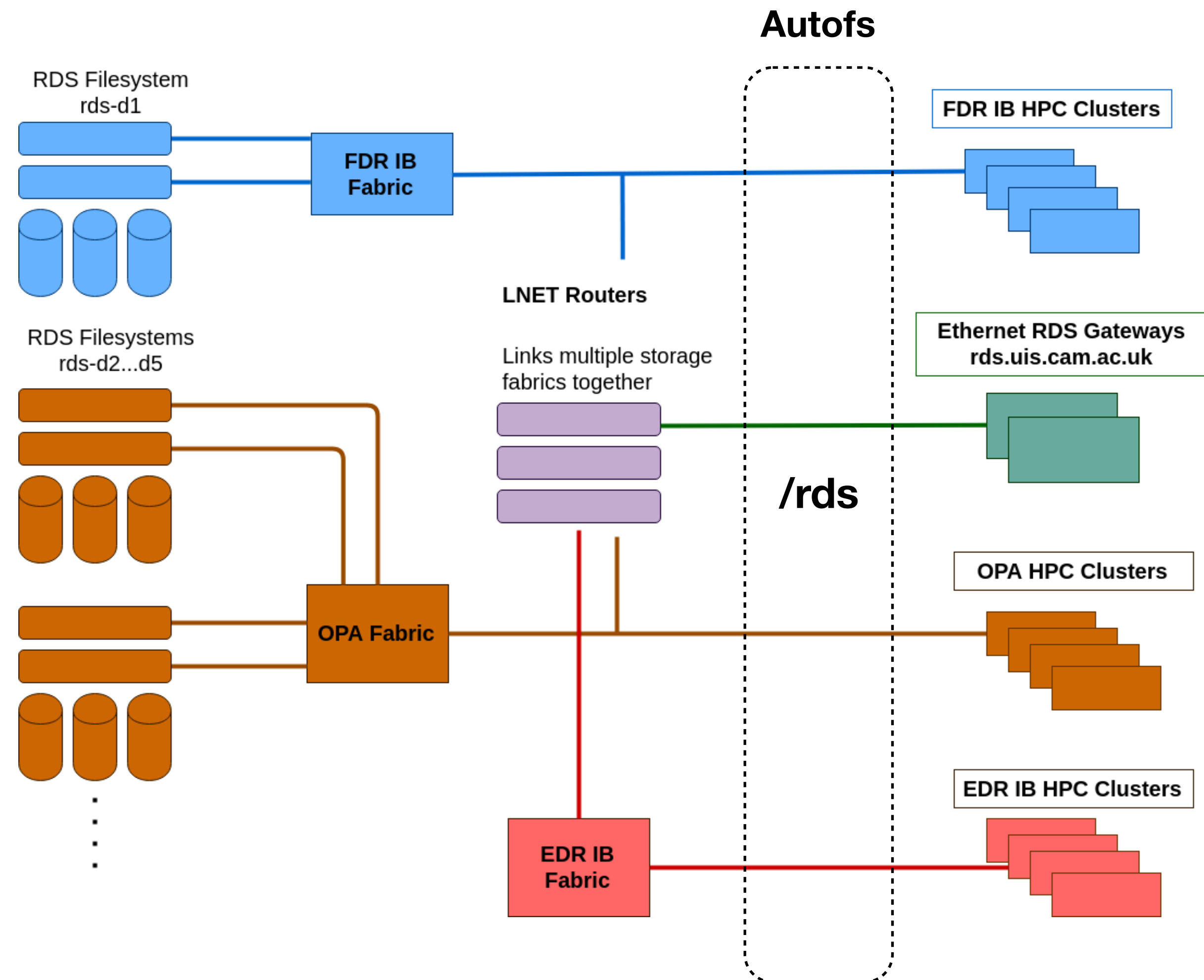
# Cambridge Research Data Store

Multiple filesystems abstracted into one global namespace.

Multiple network technologies across various access points.

LNET routing bridges networks to ensure access from all user gateways

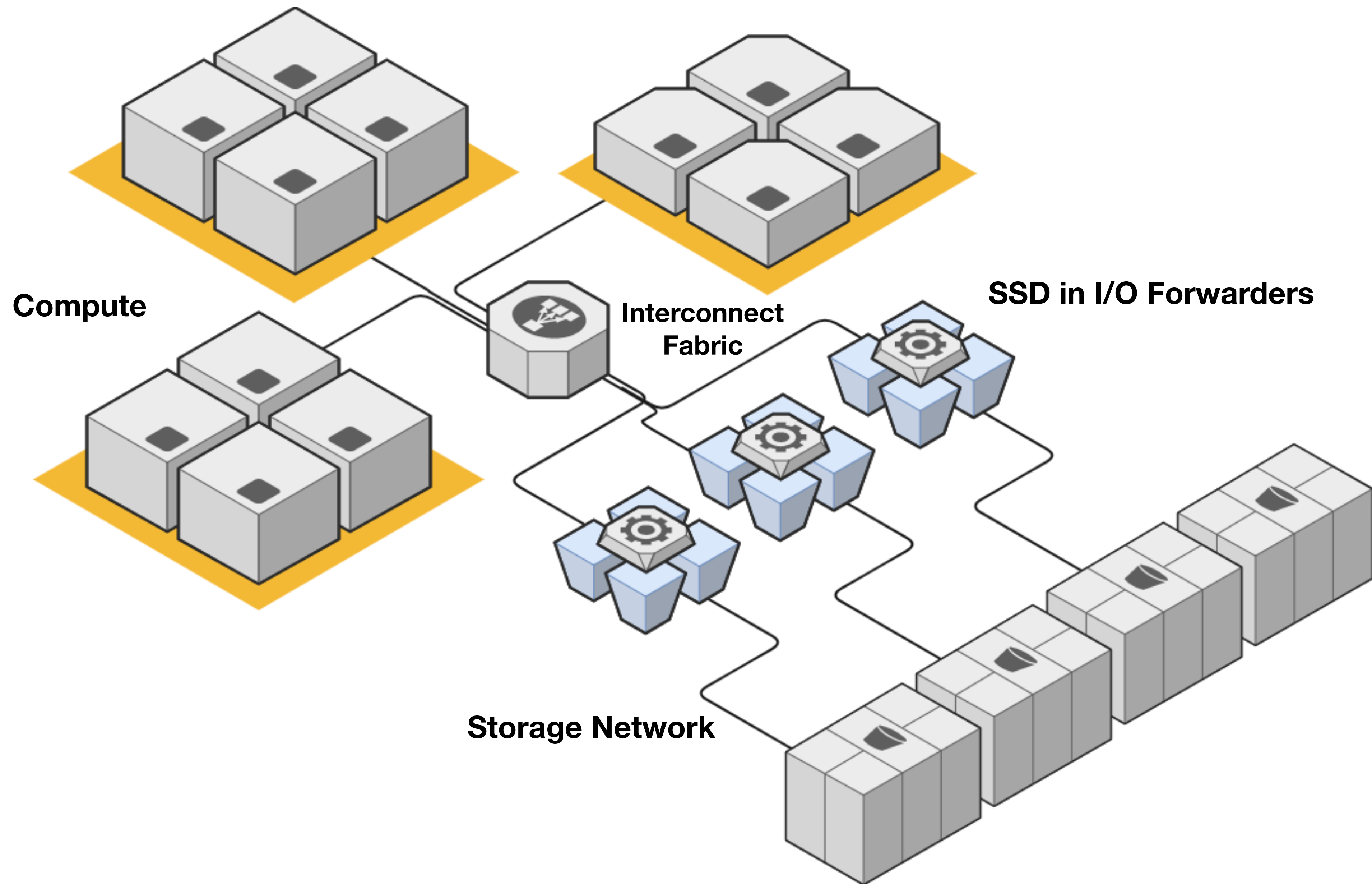
HSM Integrated with Robin-hood.



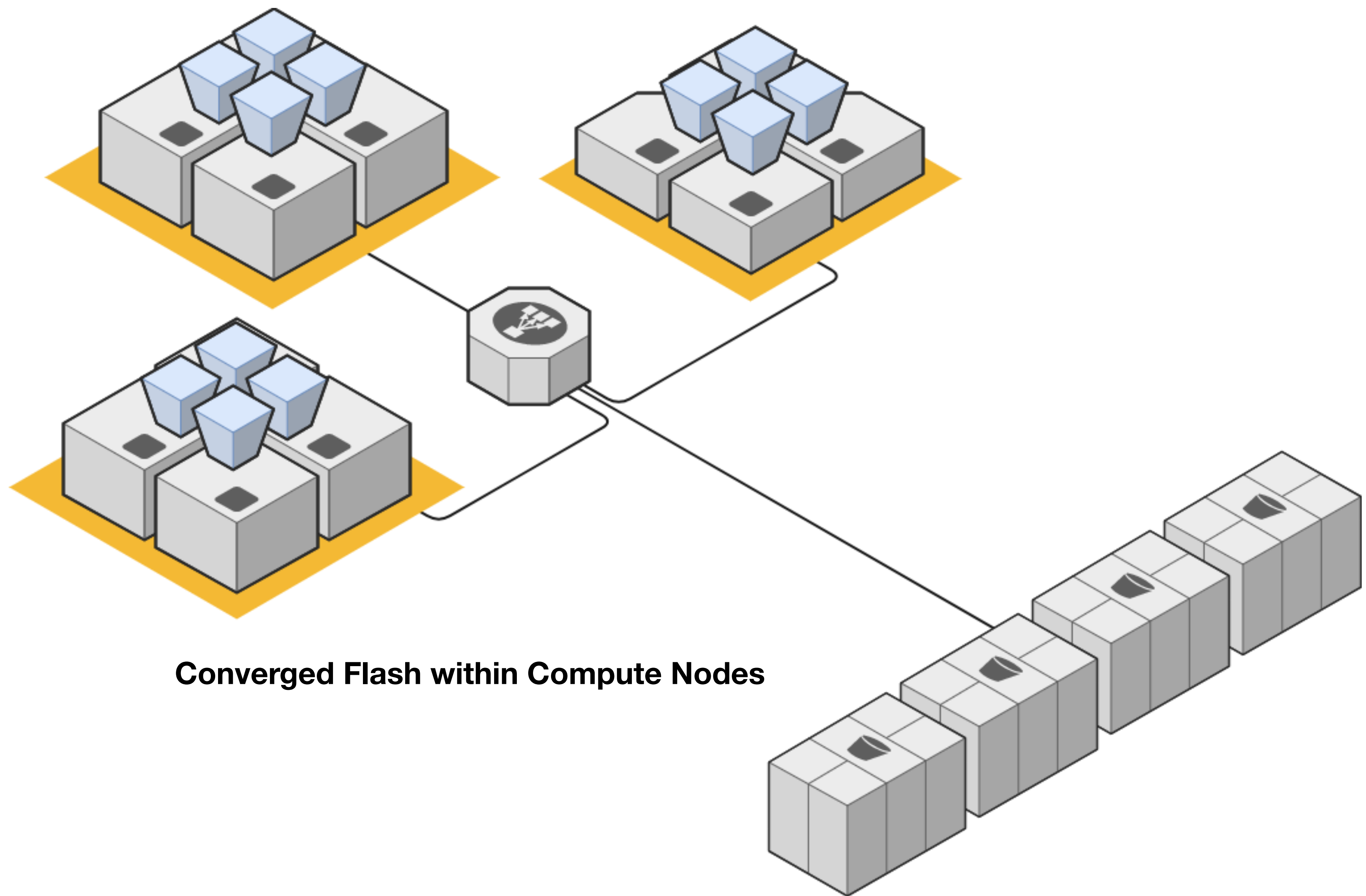


# IO Acceleration Architectures



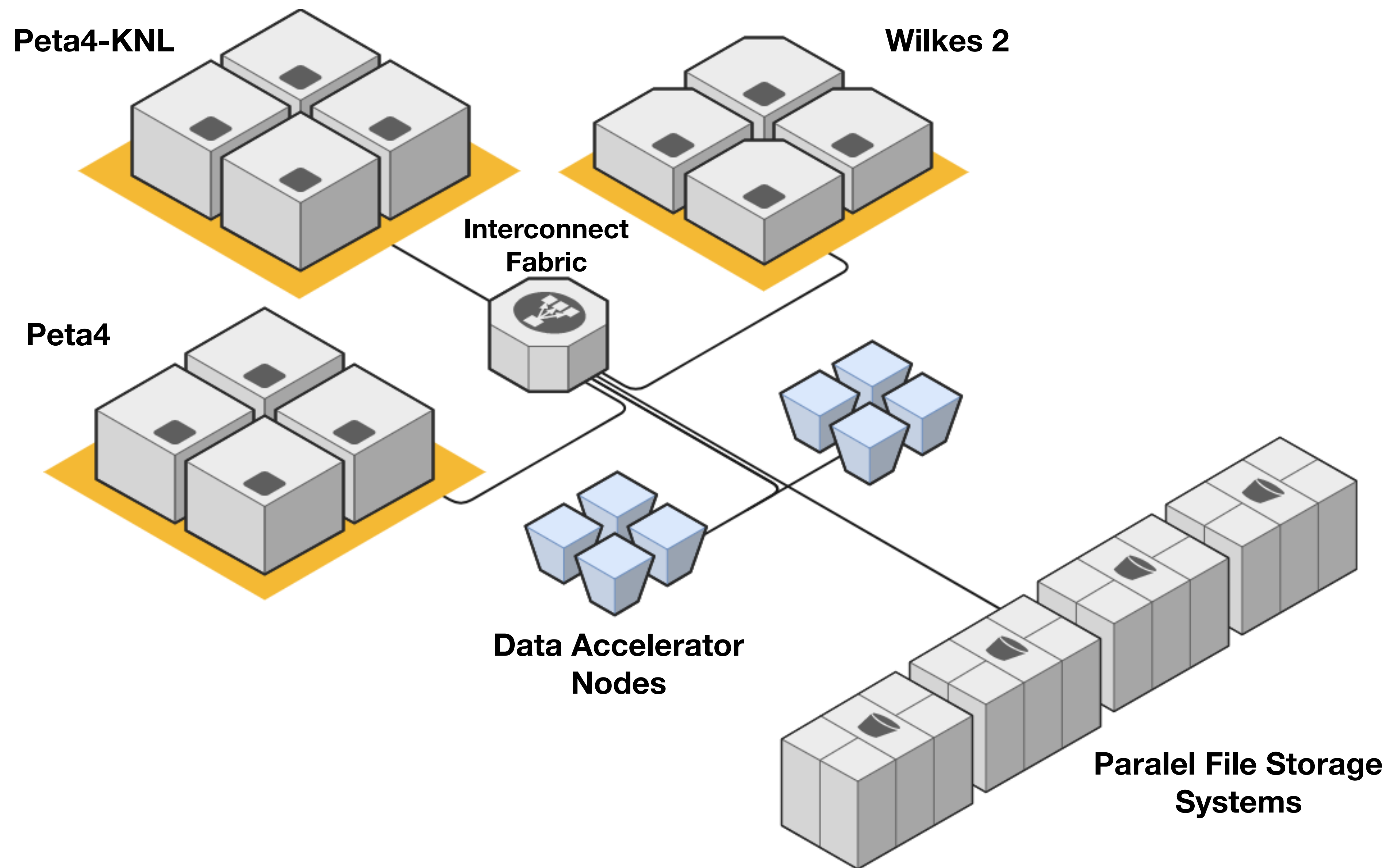






**Converged Flash within Compute Nodes**







# Data Accelerators

## Workflows and Features

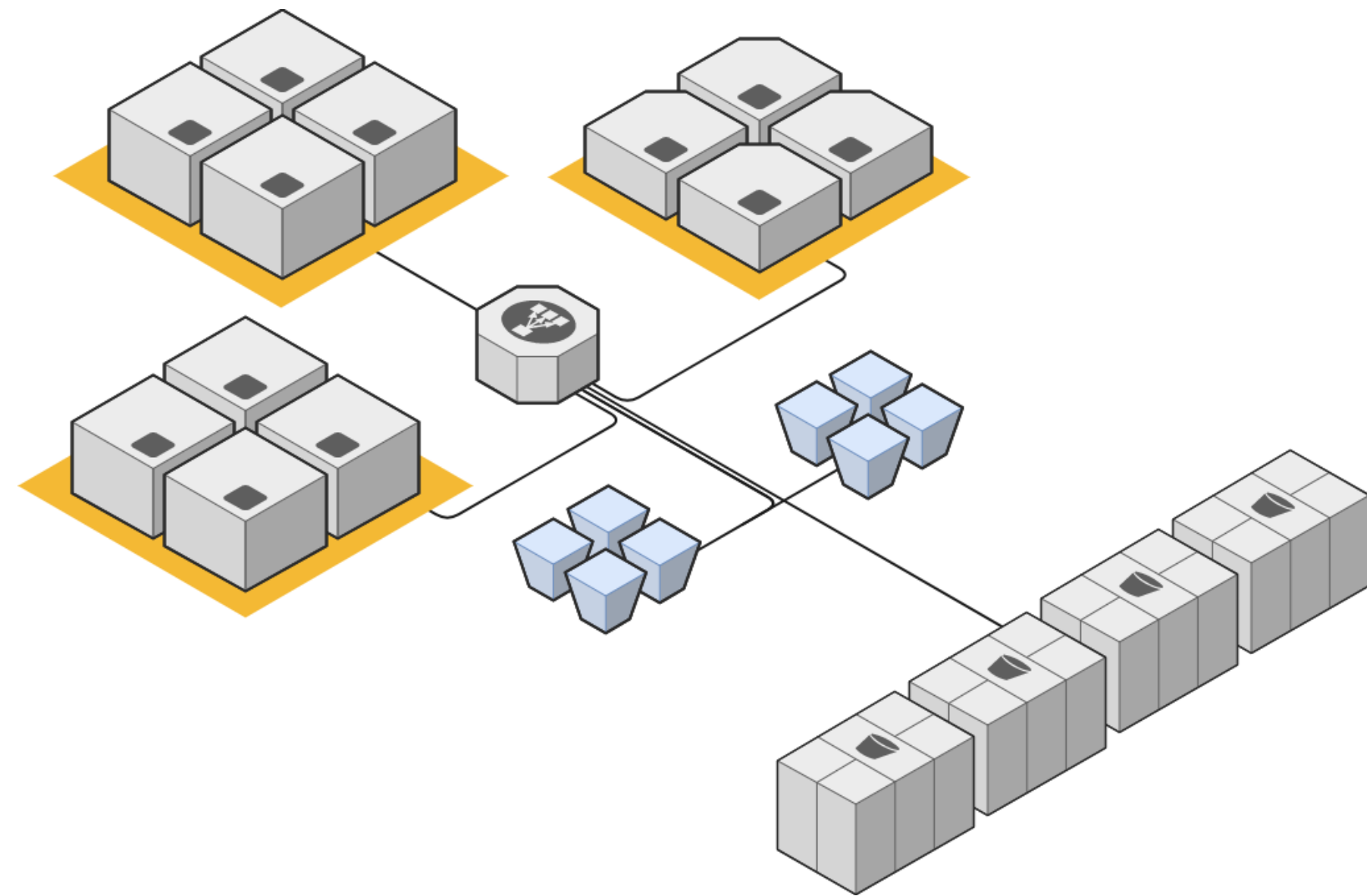
- Stage in/Stage out
- Transparent Caching
- Checkpoint
- Background data movement
- Journaling
- Swap memory

**Storage volumes - namespaces - can persist longer than the jobs and shared with multiple users, or private and ephemeral.**

**POSIX or Object ( this can also be at a flash block load/store interface )**



# The Data Accelerator Platform



**24 Dell EMC PowerEdge R740xd**  
**2 Intel Xeon Scalable Processors**  
**2 Intel Omni-Path Adaptors**  
**Each with 12 Intel SSD P4600**  
**1/2PB of Total Available Space**

Integration with SLURM and flexible storage  
orchestrator to provide maximum performance



# Integrating Lustre for the Data Accelerator



# Ansible Enabled Lustre Install

```
ansible-playbook test-dac-lustre.yml -i test-inventory-lustre --tag format --tag reformat_mgs
ansible-playbook test-dac-lustre.yml -i test-inventory-lustre --tag mount,create_mdt,create_mgs,create_osts,client_mount
ansible-playbook test-dac-lustre.yml -i test-inventory-lustre --tag stop_all,unmount,client_unmount
ansible-playbook test-dac-lustre.yml -i test-inventory-lustre --tag format

ansible-playbook test-dac-lustre.yml -i test-inventory-lustre --tag stop_mgs
ansible-playbook test-dac-lustre.yml -i test-inventory-lustre --tag reformat_mgs
```

```
*test-inventory-lustre

dac:
  children:
    fs1:
      hosts:
        dac1:
          fs1_mgs: nvme0n1
          fs1_mdt: nvme1n1
          fs1_osts: {nvme2n1: 2}
        dac2:
          fs1_osts: {nvme3n1: 1}
      vars:
        fs1_mgsnode: dac1
```

```
*test-dac-lustre.yml

---
- name: Setup buffer for fs1
  hosts: fs1
  become: yes
  roles:
    - role: lustre
  vars:
    fs_name: fs1
```



# Ansible Enabled Lustre Install

```
- name: Ensure MGS has been formatted
  command: /usr/sbin/mkfs.lustre --mgs /dev/{{ mgs }}
  register: command_result
  failed_when: "command_result.rc != 0 and ('was previously formatted for lustre' not in command_result.stderr)"
  changed_when: "command_result.rc == 0"
  when:
    - mgs is defined
  tags: [ 'never', 'format_mgs', 'format' ]

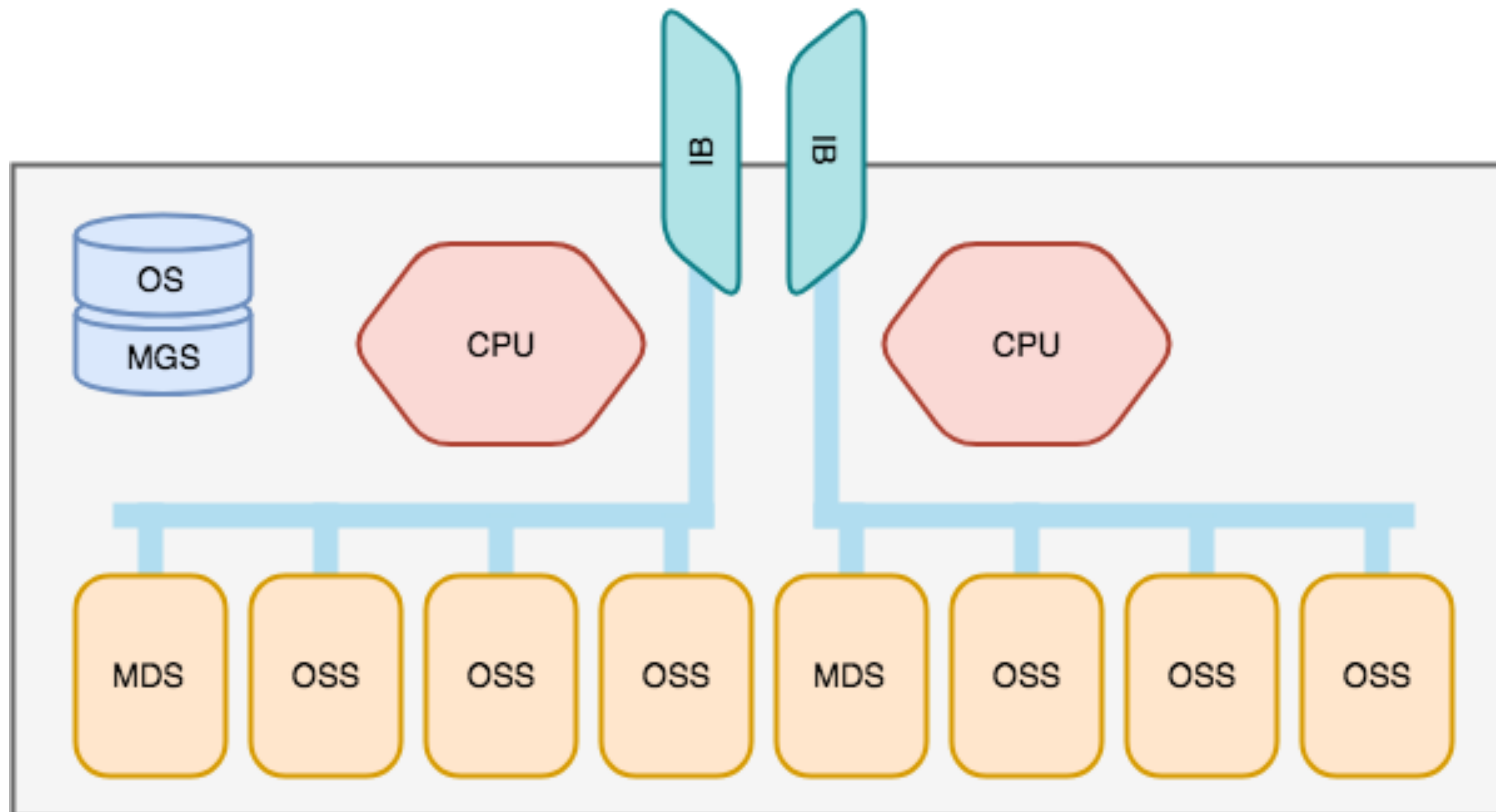
- name: Reformat MGS
  command: /usr/sbin/mkfs.lustre --mgs --reformat /dev/{{ mgs }}
  when:
    - mgs is defined
  tags: [ 'never', 'reformat_mgs' ]

- name: Reformat MDT
  command: /usr/sbin/mkfs.lustre --mdt --reformat --fsname={{ fs_name }} --index=0 --mgsnode={{ mgsnode }} /dev/{{ mdt }}
  when:
    - mdt is defined
  tags: [ 'never', 'reformat_mdts', 'format' ]

- name: Reformat OSTs
  command: /usr/sbin/mkfs.lustre --ost --reformat --fsname={{ fs_name }} --index={{ item.value }} --mgsnode={{ mgsnode }} /dev/{{ item.key }}
  loop: "{{ osts|dict2items }}"
  when:
    - osts is defined
  tags: [ 'never', 'reformat_osts', 'format' ]
```



# Lustre Server Composition



- Each DAC uses an internal SSD for the MGS should it be elected to run a file system.
- NVMeS then have an MDS or OSS applied. This arrangement can be changed as required.

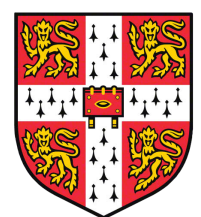


# Multirail Lustre

- Set up the ARP and Linux Kernel Routing before enabling multirail

**#Setting ARP so it doesn't broadcast  
(Do this for every IB interface)**

```
sysctl -w net.ipv4.conf.all.rp_filter=0  
sysctl -w net.ipv4.conf.ib0.arp_ignore=1  
sysctl -w net.ipv4.conf.ib0.arp_filter=0  
sysctl -w net.ipv4.conf.ib0.arp_announce=2  
sysctl -w net.ipv4.conf.ib0.rp_filter=0
```



UNIVERSITY OF  
CAMBRIDGE



Research Computing Services

# Multirail Lustre

- Set up the ARP and Linux Kernel Routing before enabling multirail

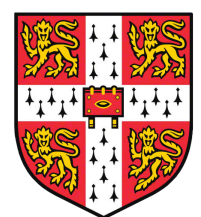
```
ip neigh flush dev ib0  
ip neigh flush dev ib1
```

```
echo 200 ib0 >> /etc/iproute2/rt_tables  
echo 201 ib1 >> /etc/iproute2/rt_tables
```

```
ip route add 192.168.0.0/16 dev ib0 proto kernel scope link src 192.168.1.1 table ib0  
ip route add 192.168.0.0/16 dev ib1 proto kernel scope link src 192.168.2.1 table ib1
```

```
ip rule add from 192.168.1.1 table ib0  
ip rule add from 192.168.2.1 table ib1
```

```
ip route flush cache
```





# Multirail Lustre

- Now set up your yaml or lnet.conf

```
#Example lnet.conf in /etc/modprobe.d
```

```
options lnet networks="o2ib1(ib0,ib1)"
```

```
#Example command to create yaml conf
```

```
lnetctl net add --net o2ib --if ib0,ib1
```

# System Performance



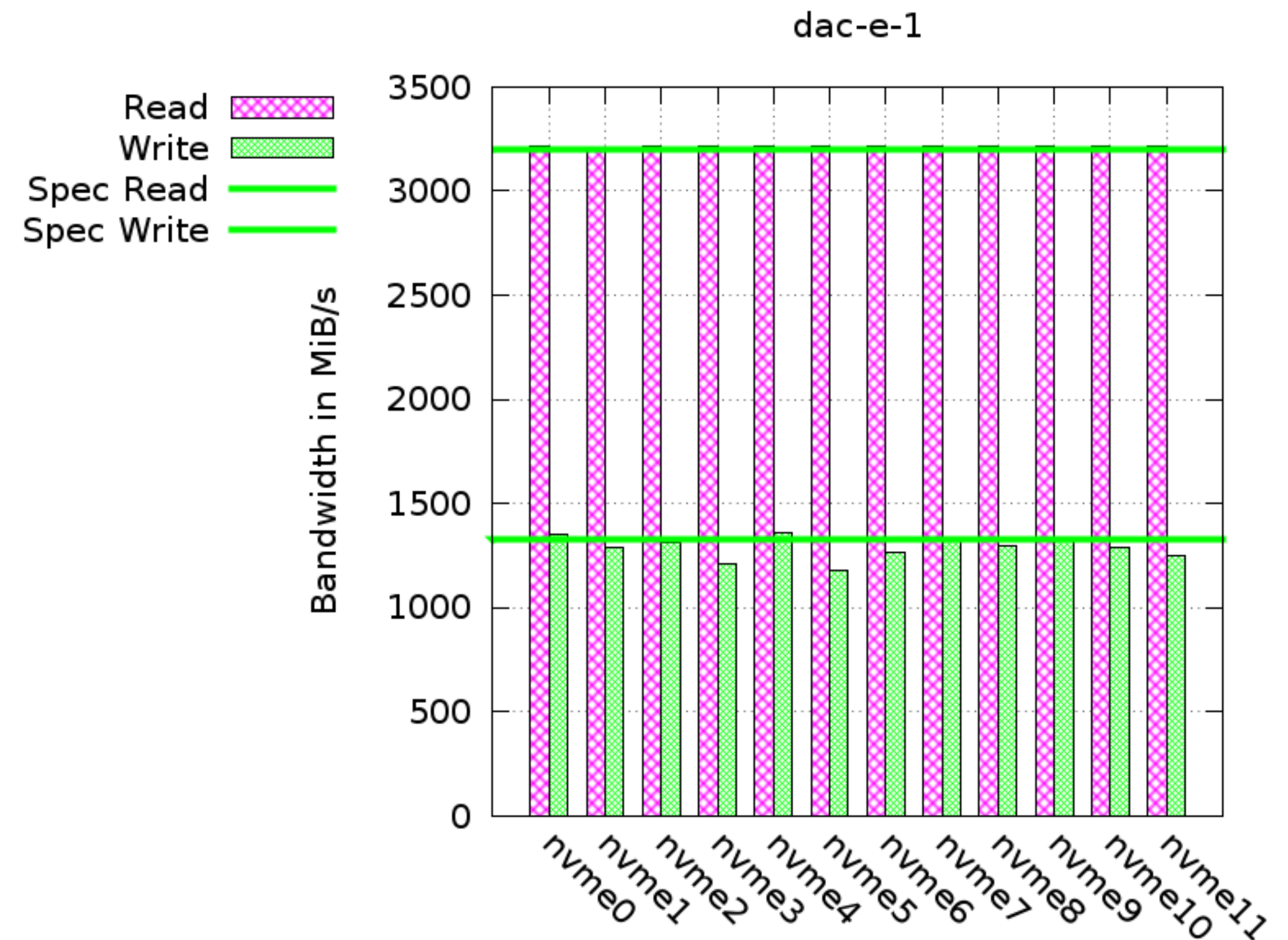
# Intel P4600 Performance

## Single FIO performance tests

```
fio --name=nvme${i} --readwrite=write -bs=128k
--invalidate=1 --group_reporting --direct=1 --time_based
--ramp_time=300 --runtime=300 --ioengine=libaio --numjobs=32
--iodepth=256 --norandommap --randrepeat=0 --stonewall
--size=100% --exitall --filename=/dev/nvme${i}n1p1
```

```
fio --name=nvme${i} --readwrite=randwrite -bs=4k
--invalidate=1 --group_reporting --direct=1 --time_based
--ramp_time=300 --runtime=300 --ioengine=libaio --numjobs=32
--iodepth=256 --norandommap --randrepeat=0 --stonewall
--size=100% --exitall --filename=/dev/nvme${i}n1p1
```

**\*for brevity similar settings for read and randread**



# Early Lustre Performance Results

## Single Rail IOR Performance

**IOR READ 264 GiB/s**

**IOR WRITE 225 GiB/s**

**DAC has 74x better efficiency of performance per TB than our scratch lustre!**

**To achieve the same performance we would be spending in the region of 10x more on traditional storage.**

**MultiRail Results should show over 250GiB/s and 600GiB/s for write and read respectively.**



# Network Performance

## LNET\_selftest Lustre 2.11.0 Clients Write

```
[LNet Bandwidth of servers]
[R] Avg: 22050.16 MiB/s Min: 22050.16 MiB/s Max: 22050.16 MiB/s
[W] Avg: 3.36 MiB/s Min: 3.36 MiB/s Max: 3.36 MiB/s
[LNet Rates of servers]
[R] Avg: 42669 RPC/s Min: 42669 RPC/s Max: 42669 RPC/s
[W] Avg: 42670 RPC/s Min: 42670 RPC/s Max: 42670 RPC/s
```

## LNET\_selftest Lustre 2.11.0 Clients Read

```
[LNet Bandwidth of servers]
[R] Avg: 3.62 MiB/s Min: 3.62 MiB/s Max: 3.62 MiB/s
[W] Avg: 23693.94 MiB/s Min: 23693.94 MiB/s Max: 23693.94 MiB/s
[LNet Rates of servers]
[R] Avg: 23697 RPC/s Min: 23697 RPC/s Max: 23697 RPC/s
[W] Avg: 47388 RPC/s Min: 47388 RPC/s Max: 47388 RPC/s
```

## LNET\_selftest.sh

```
#!/bin/bash
export LST_SESSION=$$
lst new_session read_write
lst add_group servers 10.47.18.1@o2ib1
lst add_group clients 10.47.0.[1-8]@o2ib1
lst add_batch bulk_rw
lst add_test --batch bulk_rw --concurrency 12 --from clients --to servers \
brw read check=simple size=1M
# start running
lst run bulk_rw
# display server stats for 30 seconds
lst stat servers & sleep 30; kill $!
# tear down
lst end_session
```

# Disk Performance

obdfilter test 8NVMe Lustre 2.11.0

ost 8 sz 1702887424K rsz 1024K obj 4096 thr 8192  
write 10670.58 [1237.28, 1529.71]  
rewrite 10575.61 [1153.67, 1542.69]  
read 17581.66 [2178.53, 3084.63]

ost 8 sz 838860800K rsz 1024K obj 64 thr 512  
write 10671.99 [1276.94, 1534.92]  
rewrite 10692.21 [1268.93, 1542.90]  
read 17020.97 [1485.89, 2460.10]

OSS not translating read performance correctly requires further investigation

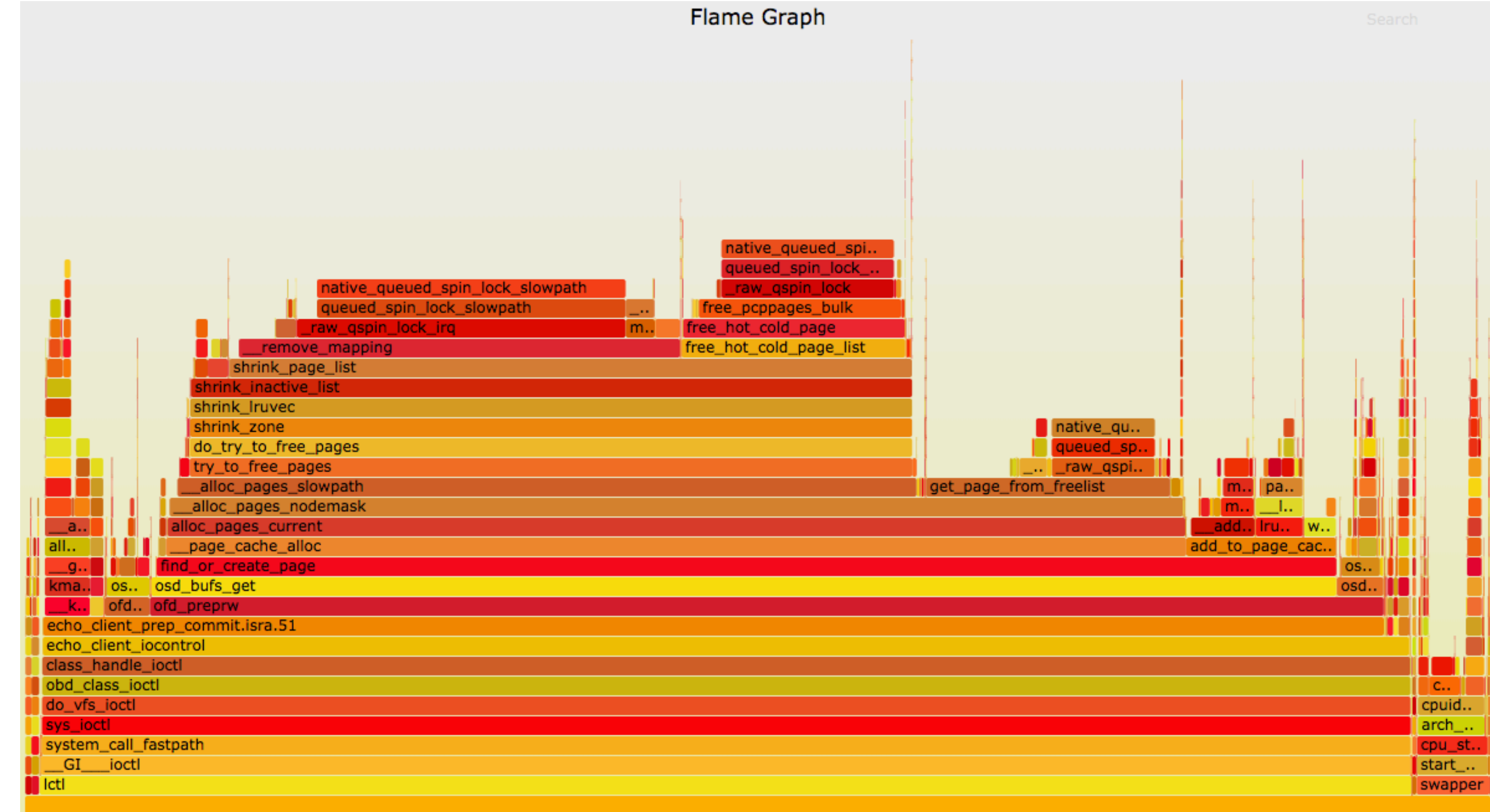
Can be made to reach ~19GiB/s but this is still lower than it should be (22-23GiB/s)



# Disk Performance



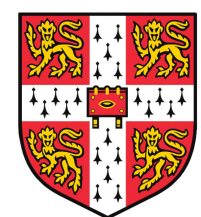
obdfilter result where threads=8192



obdfilter result where threads=512

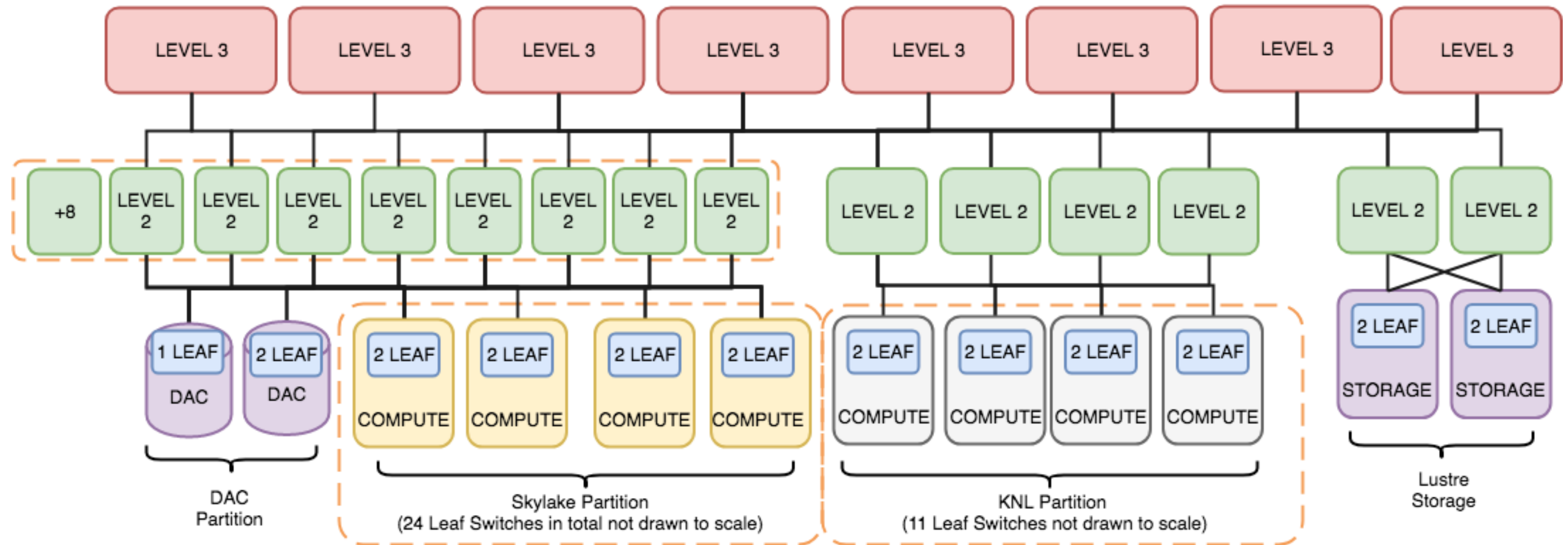
Flame Graphs show the OSS is stuck in in memory calls.  
/proc/buddyinfo shows the memory fragmentation is high  
This requires further investigation and submitted to jira.

# Disk Performance

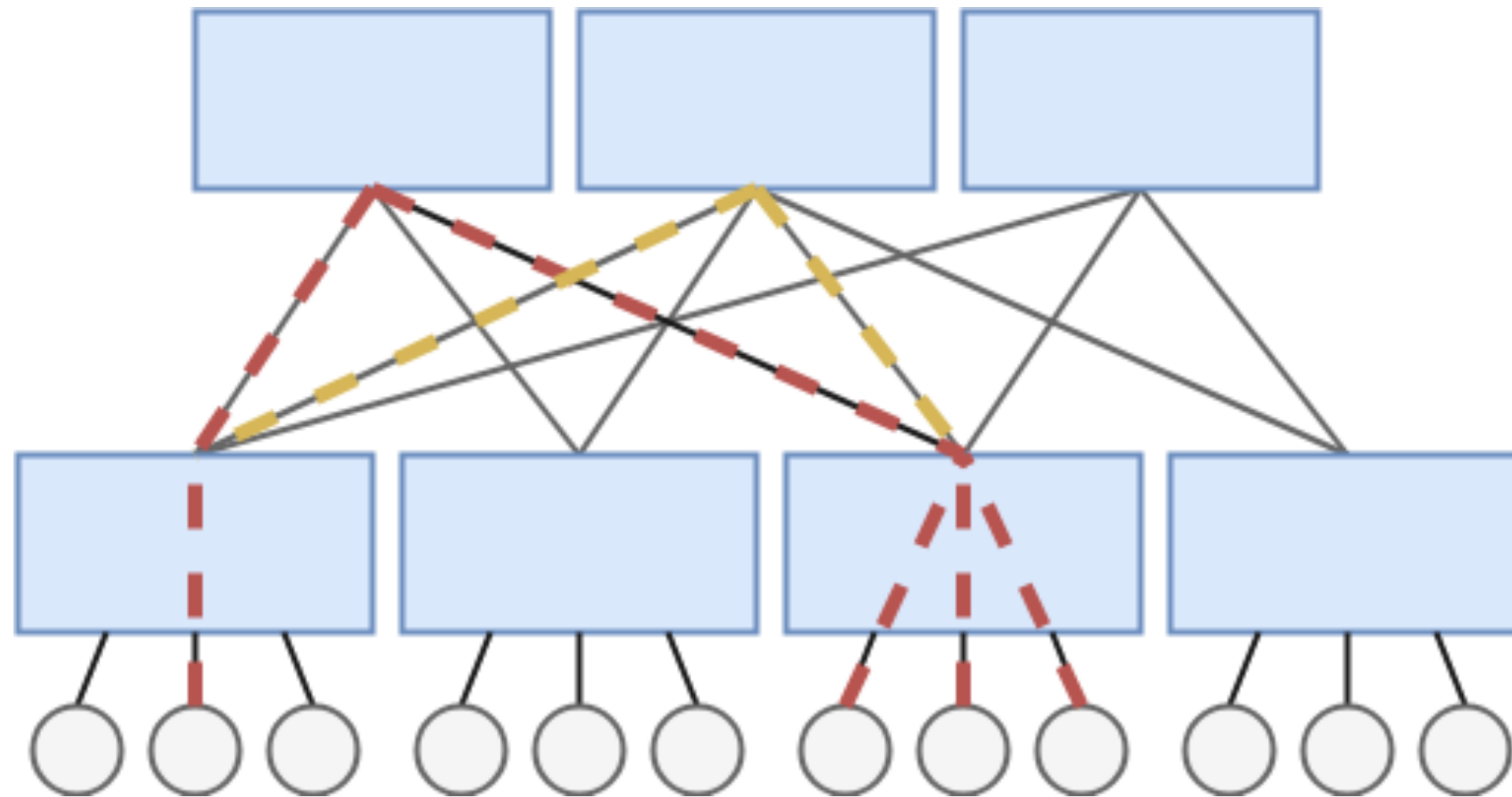




# Peta 4 Interconnect Topology



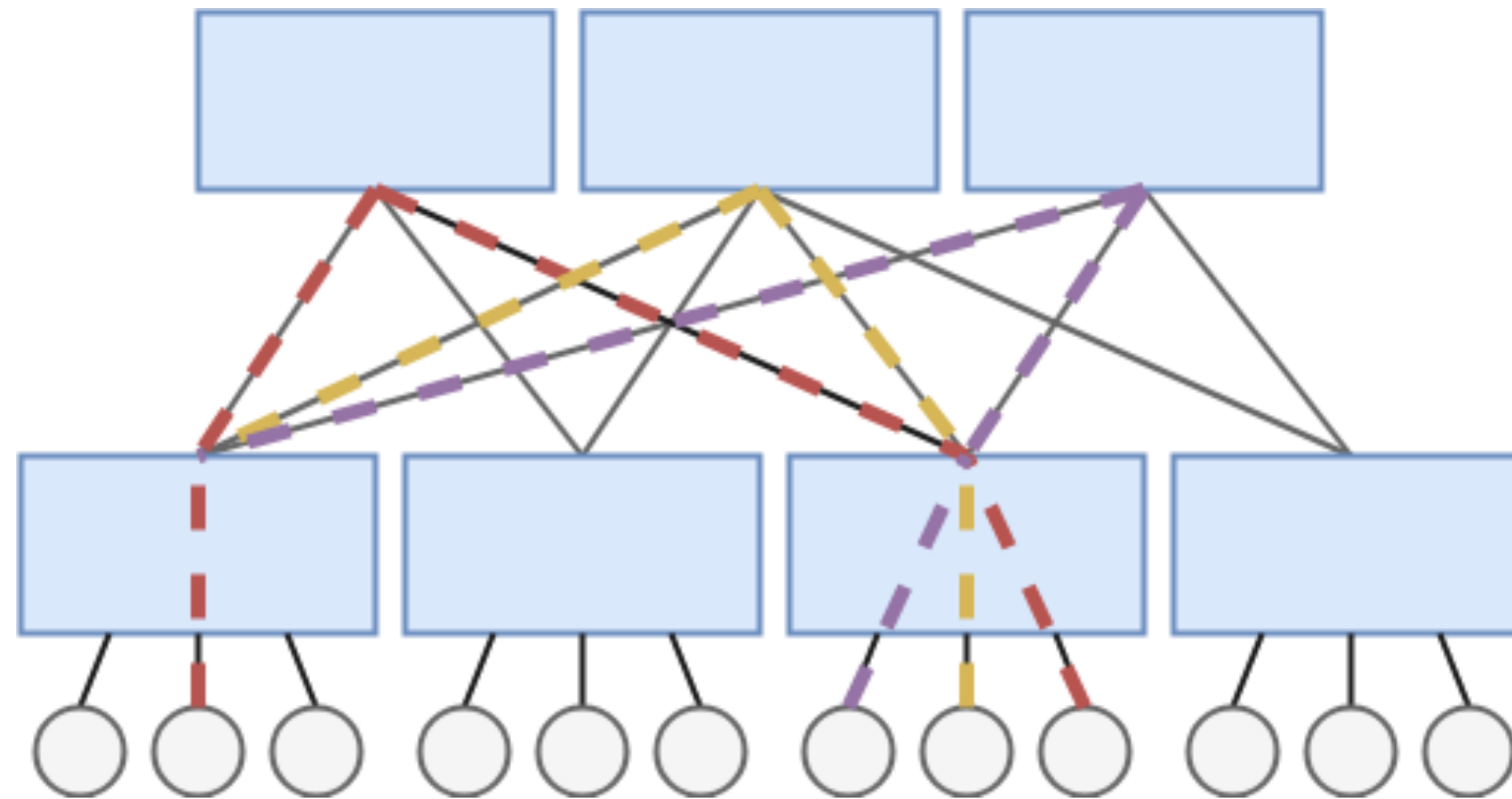
# Fat Tree Static Routing



- All nodes take the same Inter Switch Links(Red)
- Other Links are Possible(Gold)



# Adaptive Routing



- Nodes can now take alternate routes (Gold,Purple)
- Utilisation of Inter switch links improved

# Diagnosing in Intel Omni-Path

```
Group Focus: ALL   GrpNumPorts: 4015   NumPorts: 10   Number: 10
Ix  Util-High  LIDx  Port  Node GUID 0x  NodeDesc
0    0.0 0113  46  00117501020D8FAA opasw-fr16-u40
<->  84.6 04B3  3  00117501020F29B1 opasw-dr20-u35
1    0.0 0113  15  00117501020D8FAA opasw-fr16-u40
<->  74.8 0190  3  00117501020D8F8D opasw-dr20-u30
2    0.0 04B3  25  00117501020F29B1 opasw-dr20-u35
<->  70.0 04D2  19  00117501020D805E opasw-dr19-u42
3    0.1 0113  11  00117501020D8FAA opasw-fr16-u40
<->  67.4 04AE  3  00117501020F4147 opasw-dr20-u42
4    57.4 005A  3  00117501020C57AF opasw-dr20-u33
<->  0.1 0113  30  00117501020D8FAA opasw-fr16-u40
5    0.0 04D2  14  00117501020D805E opasw-dr19-u42
<->  51.9 04FB  1  00117501010DBA30 dac-e-13 hfi1_1
6    0.0 0190  41  00117501020D8F8D opasw-dr20-u30
<->  51.4 01A9  37  0011750102702B0F opasw-dr19-u41
7    0.0 04AE  9  00117501020F4147 opasw-dr20-u42
<->  49.9 04B6  33  00117501020C47F7 opasw-dr19-u42
8    0.0 018B  11  0011750102702978 opasw-fr16-u38
<->  49.8 04AE  36  00117501020F4147 opasw-dr20-u42

Quit up Live/rRev/fFwd/bookmrked Bookmrk Unbookmrk ?help | sS cC N0-n P0-n:
```

- Example of *opato* during a test. Can highlight oversubscribed links based on the percentage utilised.



# Intel Omni-Path Tuning

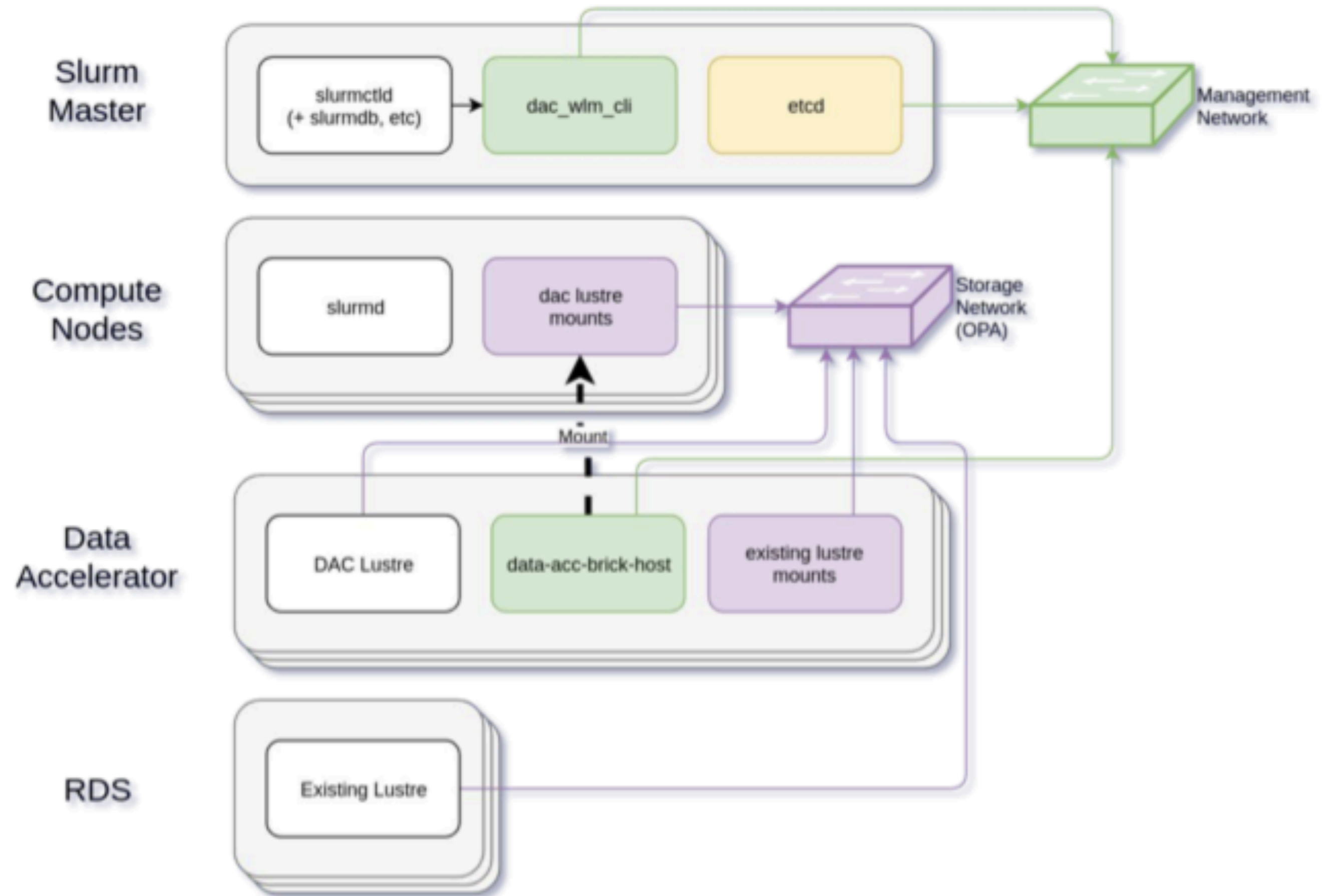
- /etc/modprob.d/hfi1.conf
- krcvqs=8, pcie\_caps=0x51,rcvhdrcnt=8192
- \*krcvqs can be varied based on node role

# Slurm Integration

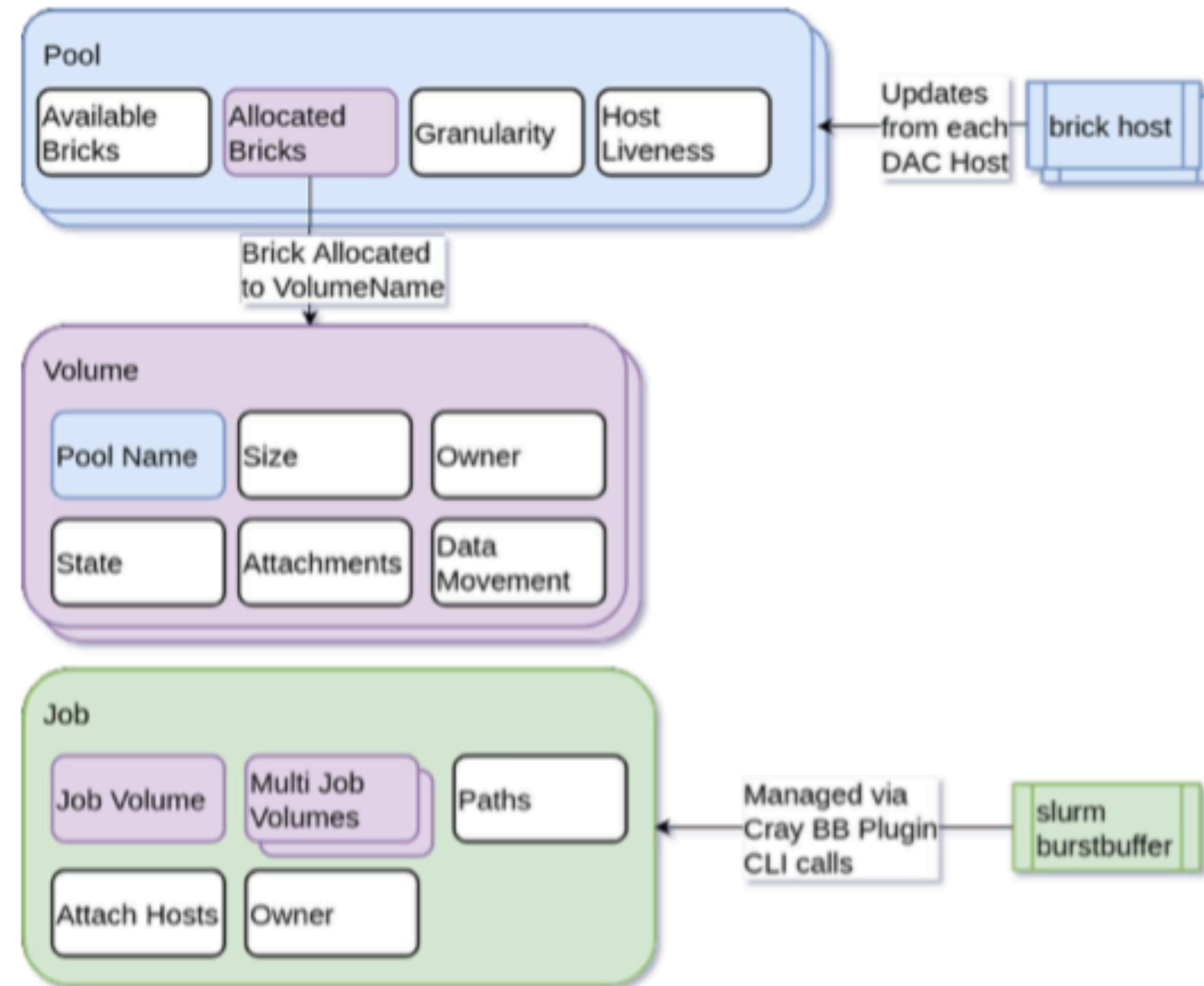


# DAC Plugin

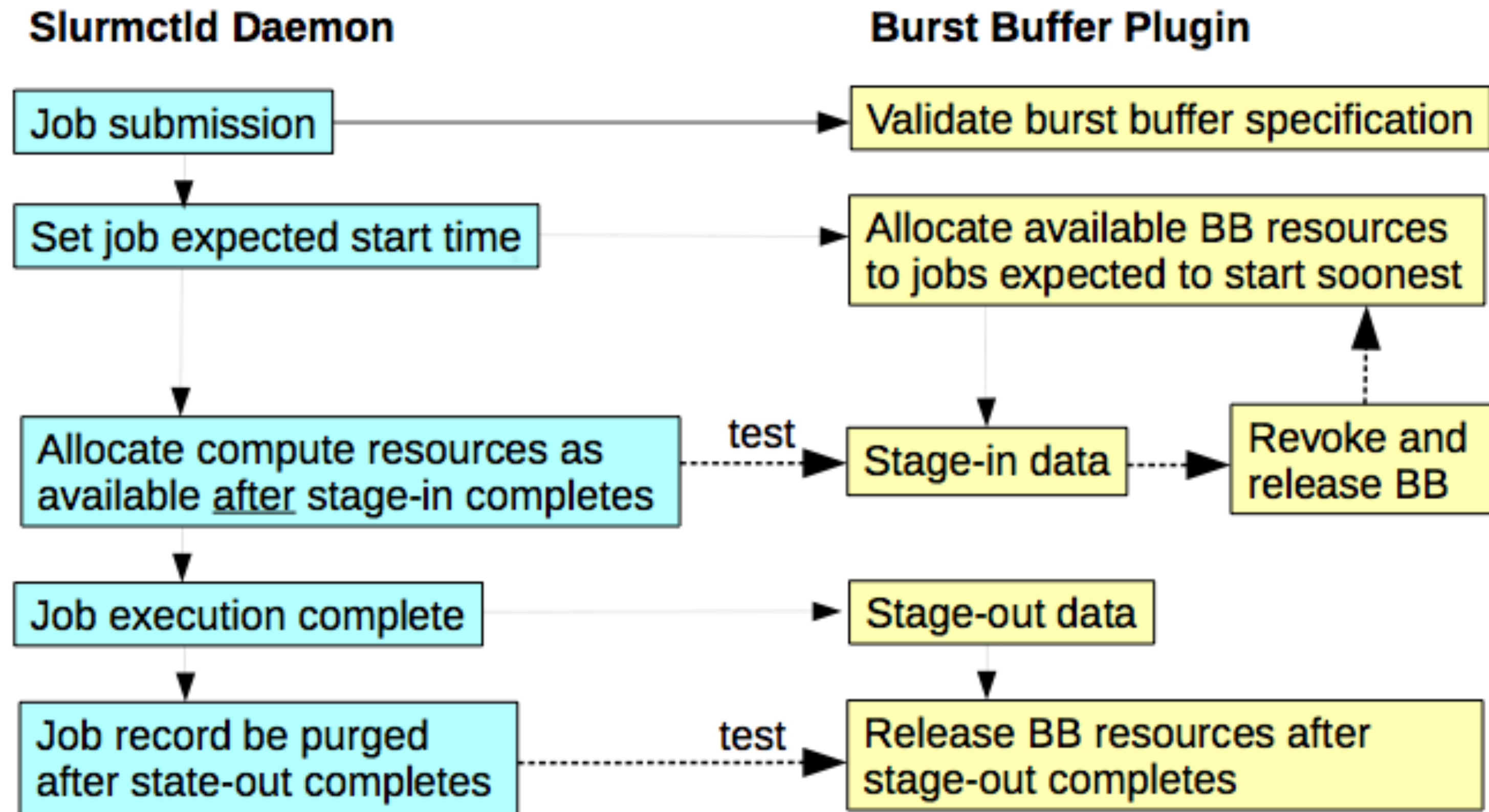
- Reuses Existing Cray plugin.
- Cambridge has implemented an orchestrator to manage the DAC nodes.
- CO-Lang project utilising ETCd and Ansible for dynamic automation
- To be released as an OpenSource project.



# SLURM DAC Plugin



# SLURM DAC Plugin



Copyright 2015 SchedMD LLC  
<http://www.schedmd.com>



# SLURM DAC Plugin

```
#!/bin/bash
#BB create_persistent name=alpha capacity=32GB access=striped type=scratch
#DW jobdw type=scratch capacity=1GB access_mode=striped
#DW stage_in type=file source=/home/alan/data.in destination=$DW_JOB_STRIPED/data
#DW stage_out type=file destination=/home/alan/data.out source=$DW_JOB_STRIPED/data
/home/alan/a.out
```

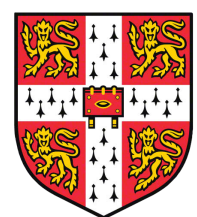
User specified burst buffer options within a slurm batch file.  
Example from [https://slurm.schedmd.com/burst\\_buffer.html](https://slurm.schedmd.com/burst_buffer.html)

# Special Thanks to...

- RCS Platforms Team
- [whamcloud.io](http://whamcloud.io)
- Intel Omni-Path Team
- Intel SSD Division
- Dell EMC Inovation Lab



# Questions and Comments?



UNIVERSITY OF  
CAMBRIDGE  
Research Computing Services

Alasdair King [ajk203@cam.ac.uk](mailto:ajk203@cam.ac.uk)



# Thanks for the Continued Support of :

The Dell EMC logo is displayed in a light blue color. It features the word "DELL" in a bold, sans-serif font, followed by a stylized icon consisting of three parallel, slanted lines. To the right of the icon is the word "EMC" in a lighter, all-caps sans-serif font.