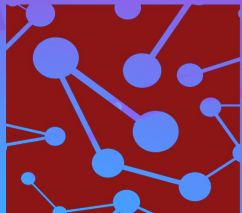


Overcoming HPC System Management Challenges: An Open Source Approach



SRCC

STANFORD RESEARCH COMPUTING CENTER

Michael Hartman (michael.hartman@stanford.edu)
- Stanford Research Computing Center (SRCC)
Stéphane Thiell (sthiell@stanford.edu)
- Stanford Research Computing Center (SRCC)
Kilian Cavalotti (kilian@stanford.edu)
- Stanford Research Computing Center (SRCC)

Outline

- Fir/Oak File Systems
- sasutils
- Oak Storage Lifecycle
- del_ost
- ibswinfo
- Slurm SPANK GPU Plugin



Lustre storage systems

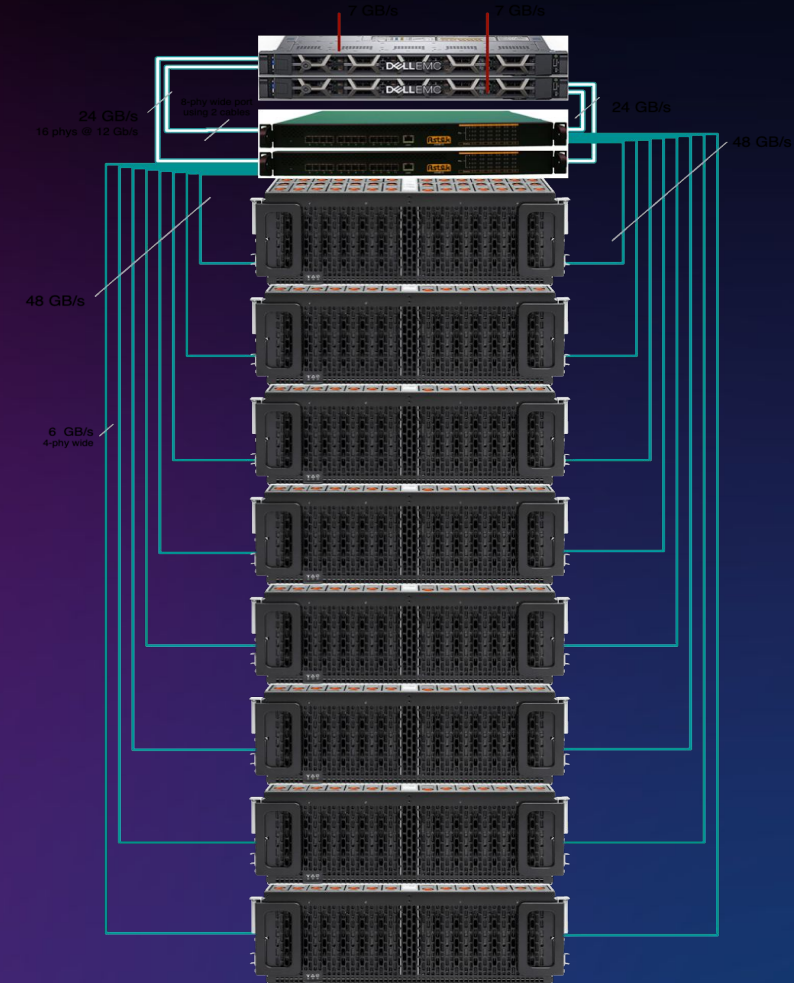
- Fir
 - Sherlock's scratch filesystem
 - High-performance storage system for temporary data
 - 4 MDTs, 96 OSTs, 960 hard drives in 8 I/O cells, 6 PB usable
- Oak
 - Global, capacity-oriented Lustre filesystem
 - 6 MDTs, 462 OSTs, 4,704 hard drives in 7 I/O cells, 52 PB usable total



sasutils

sasutils

- Display SAS fabric tree and provide aggregated view of devices
- `sas_discover`, `sas_devices`, `sas_counters`, `ses_report`
- Based on `sysfs` (and also `sg3_utils` and `smp_utils`)
- Support SES Enclosure Nickname
- Open source, available at <https://github.com/stanford-rc/sasutils>
- Published in EPEL 7, EPEL 8 and soon in EPEL 9
- Or use “`pip install sasutils`”



sasutils

- sas_devices

```
$ sas_devices
Found 2 SAS hosts
Found 4 SAS expanders
Found 1 enclosure groups
Enclosure group: [io1-jbod1-0][io1-jbod1-1]
NUM      VENDOR      MODEL  REV  PATHS
60 x     SEAGATE  ST8000NM0075  E002   2
Total: 60 block devices in enclosure group
```



sasutils

- sas_devices
- sas_discover

```
$ sas_discover -v
oak-io1-s1
|--host35 SAS9300-8e
|  |--8x--expander-35:0 ASTEK
|      |--1x--end_device-35:0:0
|          |--enclosure io1-sassw1 ASTEK
|              |--4x--expander-35:1 QCT
|                  |-- 60 x end_device -- disk
|                      |-- 1 x end_device -- enclosure io1-jbod1-0 QCT
`--host36 SAS9300-8e
    |--8x--expander-36:0 ASTEK
        |--1x--end_device-36:0:0
            |--enclosure io1-sassw2 ASTEK
                |--4x--expander-36:1 QCT
                    |-- 60 x end_device -- disk
                        |-- 1 x end_device -- enclosure io1-jbod1-1 QCT
```



sasutils

- sas_devices
- sas_discover
- sas_counters

```
$ sas_counters
```

```
...
```

```
oak-io1-s1.SAS9300-8e...Switch184.io1-sassw1.JB4602_SIM_0.io1-jbod1-0.bays.41.ST8000NM0075...ioerr_cnt 2 1487457378
```

```
oak-io1-s1.SAS9300-8e...Switch184.io1-sassw1.JB4602_SIM_0.io1-jbod1-0.bays.41.ST8000NM0075...iodone_cnt 7154904 1487457378
```

```
oak-io1-s1.SAS9300-8e...Switch184.io1-sassw1.JB4602_SIM_0.io1-jbod1-0.bays.41.ST8000NM0075...iorequest_cnt 7154906 1487457378
```

```
...
```

```
oak-io1-s1.SAS9300-8e.0x500605b00ab05678.Switch184.io1-sassw2.phys.15.invalid_dword_count 5 1487457378
```

```
oak-io1-s1.SAS9300-8e.0x500605b00ab05678.Switch184.io1-sassw2.phys.15.loss_of_dword_sync_count 1 1487457378
```

```
oak-io1-s1.SAS9300-8e.0x500605b00ab05678.Switch184.io1-sassw2.phys.15.phy_reset_problem_count 0 1487457378
```

```
oak-io1-s1.SAS9300-8e.0x500605b00ab05678.Switch184.io1-sassw2.phys.15.running_disparity_error_count 1 1487457378
```

```
...
```



sasutils

- sas_devices
- sas_discover
- sas_counters
- ses_report

```
$ ses_report --carbon --prefix=datacenter.stanford
datacenter.stanford.io1-sassw1.Cooling.Left_Fan.speed_rpm 19560 1476486766
datacenter.stanford.io1-sassw1.Cooling.Right_Fan.speed_rpm 19080 1476486766
datacenter.stanford.io1-sassw1.Cooling.Center_Fan.speed_rpm 19490 1476486766
```

```
# ses_report --status --prefix=datacenter.stanford | grep SIM
datacenter.stanford.io1-jbod1-0.Enclosure_services_controller_electronics.SIM_00 OK
datacenter.stanford.io1-jbod1-0.Enclosure_services_controller_electronics.SIM_01 OK
datacenter.stanford.io1-jbod1-0.SAS_expander.SAS_Expander_SIM_0 OK
datacenter.stanford.io1-jbod1-0.SAS_expander.SAS_Expander_ISIM_2 OK
datacenter.stanford.io1-jbod1-0.SAS_expander.SAS_Expander_ISIM_0 OK
datacenter.stanford.io1-jbod1-1.Enclosure_services_controller_electronics.SIM_00 OK
datacenter.stanford.io1-jbod1-1.Enclosure_services_controller_electronics.SIM_01 OK
datacenter.stanford.io1-jbod1-1.SAS_expander.SAS_Expander_SIM_1 OK
datacenter.stanford.io1-jbod1-1.SAS_expander.SAS_Expander_ISIM_3 OK
datacenter.stanford.io1-jbod1-1.SAS_expander.SAS_Expander_ISIM_1 OK
```

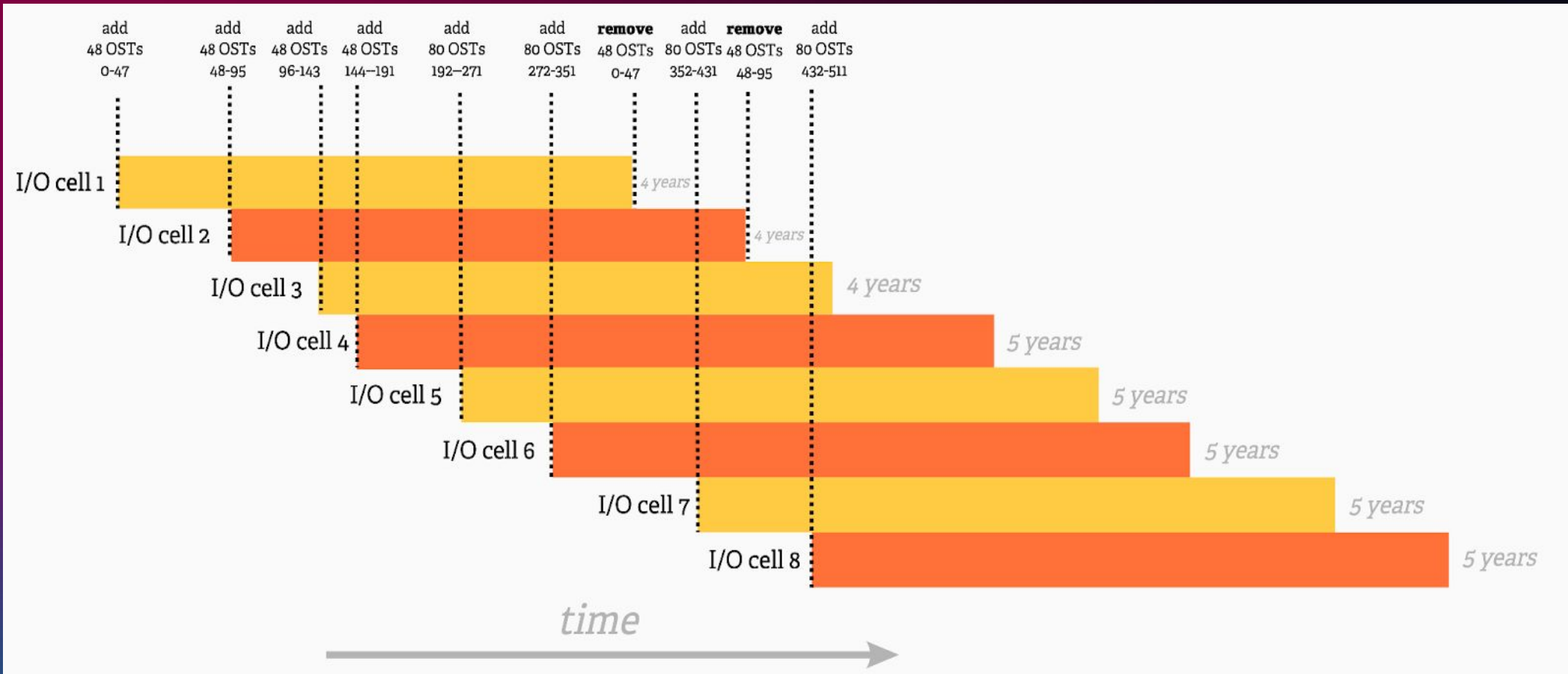


sasutils

- sas_counters
- sas_discover
- sas_devices
- ses_report
- udev scripts
 - sas_sd_snic_alias
 - sas_mpath_snic_alias
- Python library
 - Listing SAS Hosts, Listing Expanders, Listing Unique Expanders, SCSI Enclosure



Oak storage lifecycle overview



lctl del_ost

- Introduced in Lustre 2.16 (to be release mid-2023)
 - Documentation is located :
 - https://build.whamcloud.com/job/lustre-manual/lastSuccessfulBuild/artifact/lustre_manual.xhtml#lustremaint.remove_ost



ibswinfo

- Leverages the MFT (Mellanox Firmware Tools)
 - Normally used for Firmware updates on switches and IB cards
 - Can also probe the hardware and provide status data



ibswinfo

```
# ./ibswinfo.sh -d /dev/mst/<device>
```

```
=====
<node description>
=====
```

```
part number      | MQM8790-HS2F
serial number    | <redacted>
product name     | Jaguar Unmng IB
200
revision         | AC
ports            | 40
PSID             | MT_0000000063
GUID            | <redacted>
firmware version | 27.2000.1886
```

```
-----
uptime (d-h:m:s) | 196d-07:05:40
-----
```

```
PSU0 status      | OK
P/N              | MTEF-PSF-AC-C
S/N              | <redacted>
DC power         | OK
fan status       | OK
power (W)        | 64
PSU1 status      | OK
P/N              | MTEF-PSF-AC-C
S/N              | <redacted>
DC power         | OK
fan status       | OK
power (W)        | 47
```

```
-----
temperature (C)  | 34
max temp (C)     | 41
-----
```

```
fan status       | OK
fan#1 (rpm)      | 5426
fan#2 (rpm)      | 4746
fan#3 (rpm)      | 5426
fan#4 (rpm)      | 4798
fan#5 (rpm)      | 5426
fan#6 (rpm)      | 4815
fan#7 (rpm)      | 5382
fan#8 (rpm)      | 4868
fan#9 (rpm)      | 5471
-----
```



Slurm SPANK GPU Plugin

- Slurm Plug-in Architecture for Node and job (K)control - SPANK
- Written in Lua
- Requires that the slurm-spank-lua plugin is installed
- Flag:
 - `--gpu_cmode=[shared | exclusive | prohibited]`



Conclusion

sasutils: <https://github.com/stanford-rc/sasutils>

ibswinfo: <https://github.com/stanford-rc/ibswinfo>

Slurm

SPANK

GPU

Plugin:

https://github.com/stanford-rc/slurm-spank-gpu_cmode

