

Computing & Storage for on-site experiments - Petra3, FLASH and EuXFEL

Martin Gasthuber, ASAP³, FS-EC and Eu-XFEL Team
BNL, September 24, 2018



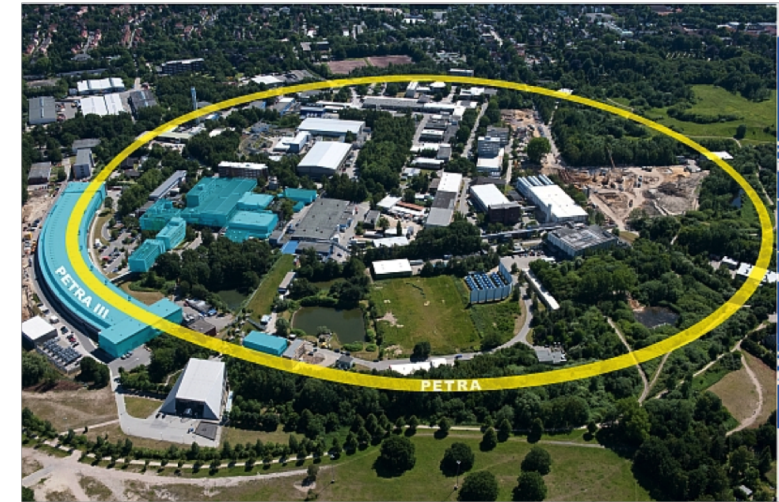
the bigger instruments...

- currently three on-site accelerators are in operation
 - Petra 3
 - FLASH
 - EuXFEL
- all three share some of the computing infrastructure and have much of the architecture and 'mode of operations' in common
 - same team at central IT complemented by teams at each instrument (and vice versa ;-)

Petra 3

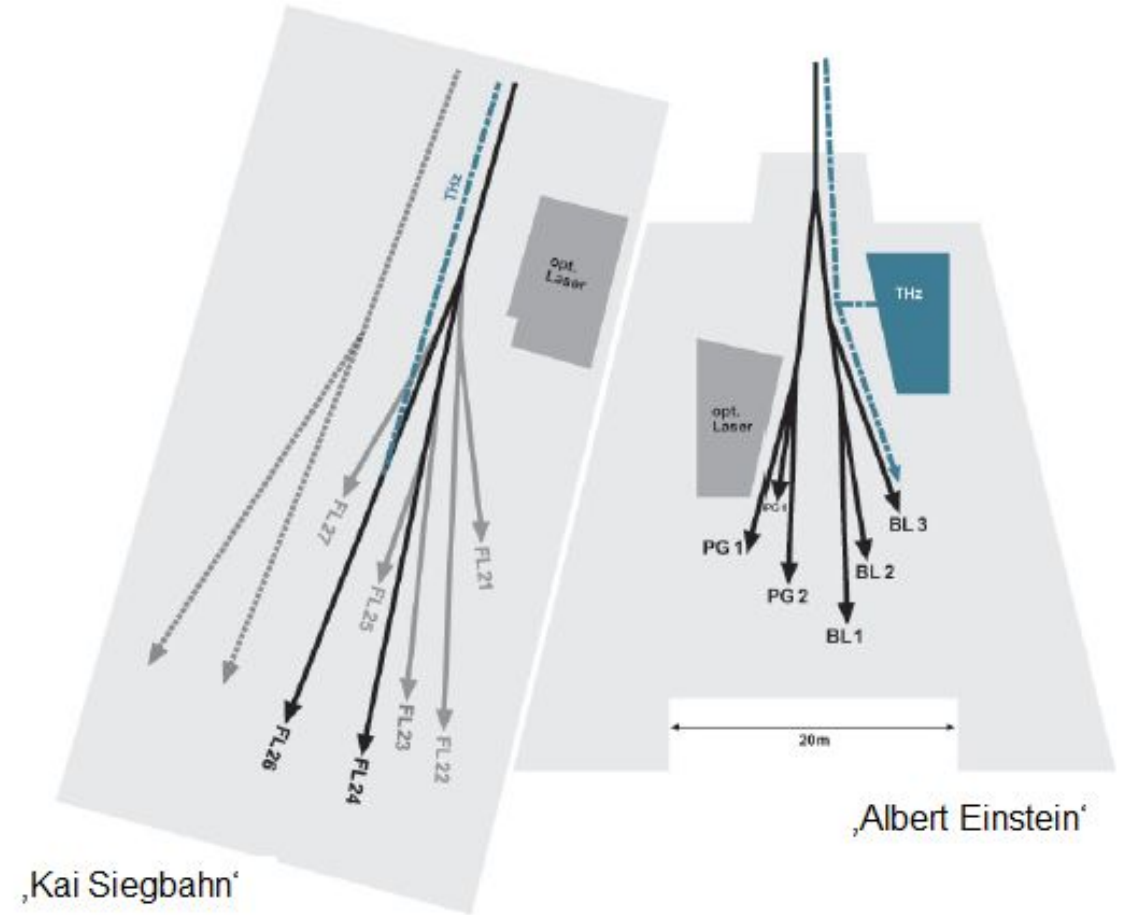
oldest workhorse, recently extended, inherited from particle physics era

- > built 1978 for HEP Experiments
- > Since 2009: 14 beamlines in operation
- > Since 2016: 10 additional beamlines



FLASH 1+2

- Linear accelerator
- started as test facility for TESLA technology
- since 2016 using ASAP³ for storage & analysis
- 12 beamlines (2 undulator lines) – one active per undulator



FLASH experimental halls

DOOR

DESY Online Office for Research with Photons

DOOR

- Workflow system for experiments processes at DESY Photon Science (PETRA III, FLASH)
 - proposal submission, internal and external review, beam time scheduling etc.
- Web application
 - php-based
 - Permanent further development of DUO clone (originated at PSI, Switzerland)
 - CentOS Linux Server, Apache
- Data stored in central Oracle Database Server
- Role-based System
- Other DUO clones are in use at EXFEL (UPEX), Hamburg and MAX IV (DUO), Lund, Sweden

A screenshot of the DOOR website. The header includes the DESY logo and navigation links for 'ACCELERATORS | PHOTON SCIENCE | PARTICLE PHYSICS'. The main content area features a 'WELCOME' message, 'NEXT DEADLINES' for FLASH and PETRA III, and sections for 'Registered DOOR user' and 'New DOOR user'. The footer contains the Helmholtz logo and contact information.

ACCELERATORS | PHOTON SCIENCE | PARTICLE PHYSICS
Deutsches Elektronen-Synchrotron
A Research Centre of the Helmholtz Association

DESY PHOTON SCIENCE »

DOOR HOME | CONTACT

DOOR

DESY Online Office for Research with Photons

HOME

- New User
- Lost Password
- Registered User

WELCOME

Welcome to DOOR, the DESY Online Office for Research with Photons.
After registration, you may use this system to submit research proposals or beamtime applications, and to complete all administrative steps required prior and after your experiment (e.g. online safety training, submission of declaration of substances, registration of participants, travel reimbursement reports, registration of publications).
Please do not hesitate to [contact us](#) in case you have further questions or if you encounter any problems using DOOR.

NEXT DEADLINES

Proposal Submission FLASH	01-Oct-2018
Proposal Submission PETRA III	No open call at present. For future deadlines see DESY Photon Science webpage .

Registered DOOR user
Log on using your DOOR user name and password or your Umbrella credentials.

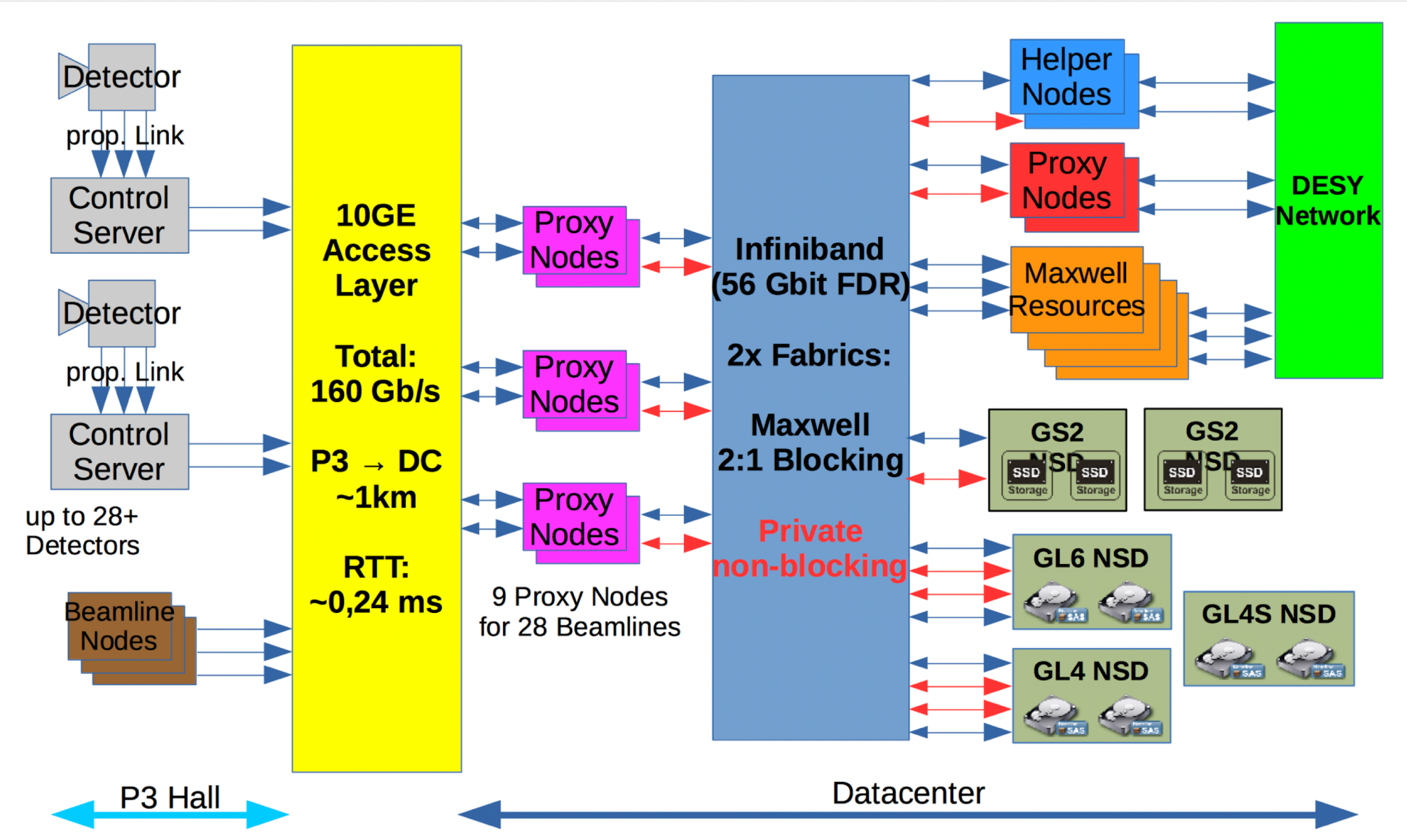
Forgotten password
If you do not remember your DOOR user name and/or password, your log on information will be sent to your previously registered e-mail address.

New DOOR user
To obtain a DOOR user name and password, please register here. Users with an existing Umbrella account might first log on at Umbrella [here](#).
Or you might set up an Umbrella account before registering [here](#).

HELMHOLTZ RESEARCH FOR GRAND CHALLENGES

Contact | DESY Data Privacy Policy | DOOR Data Privacy Policy | Imprint
© 2018 Deutsches Elektronen-Synchrotron DESY

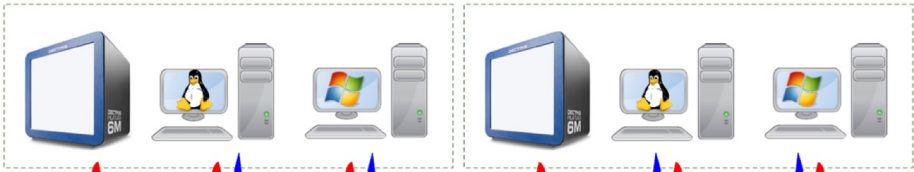
layout – networks, hardware



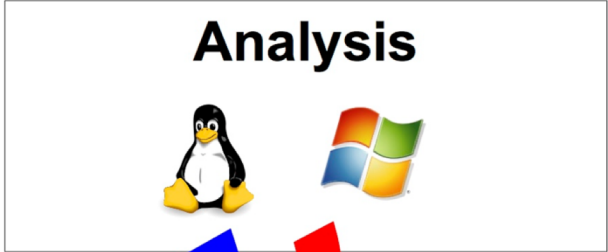
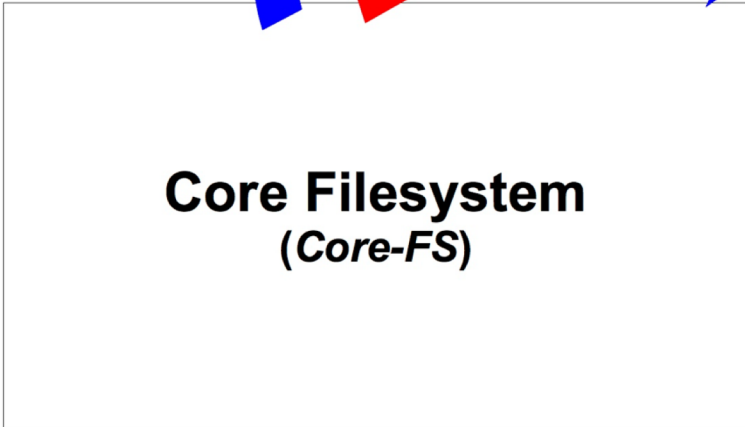
915 disk drives
 96 SSDs
 ~540 Nodes

Overview – dataflow & services

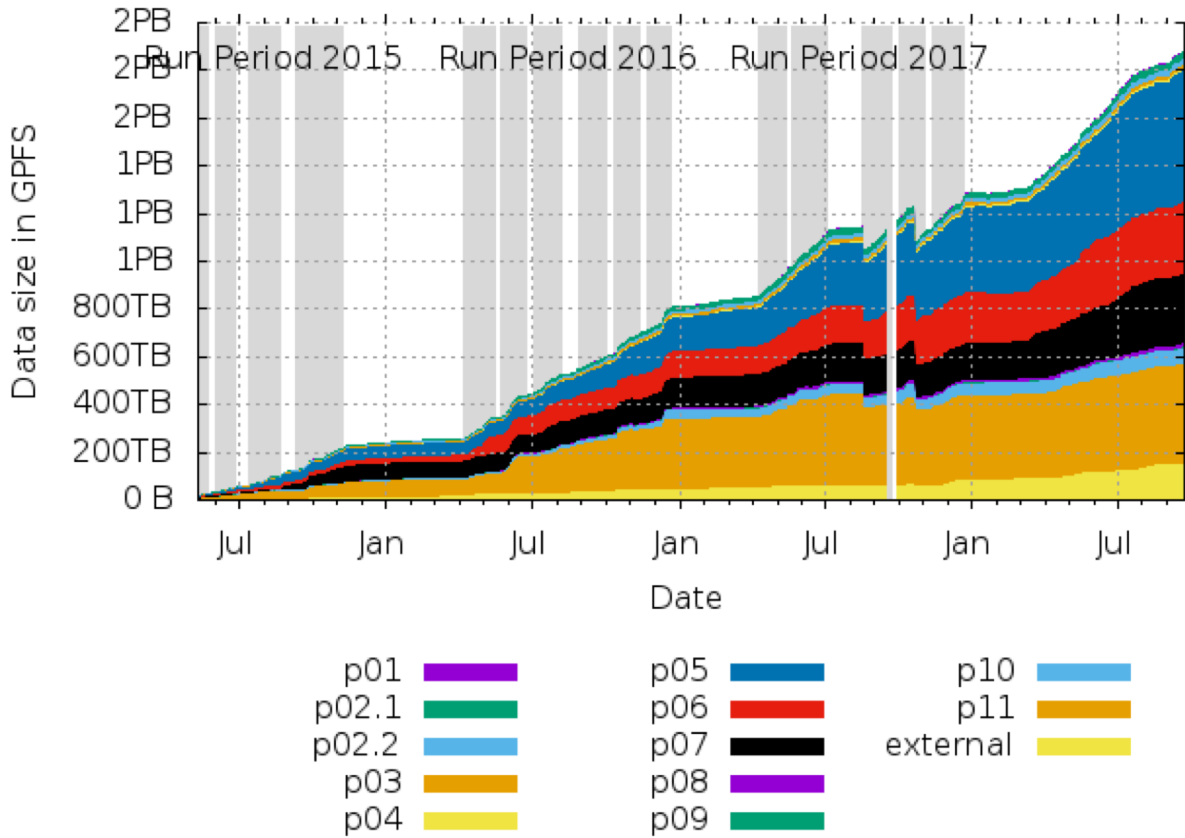
Sandbox per Beamline



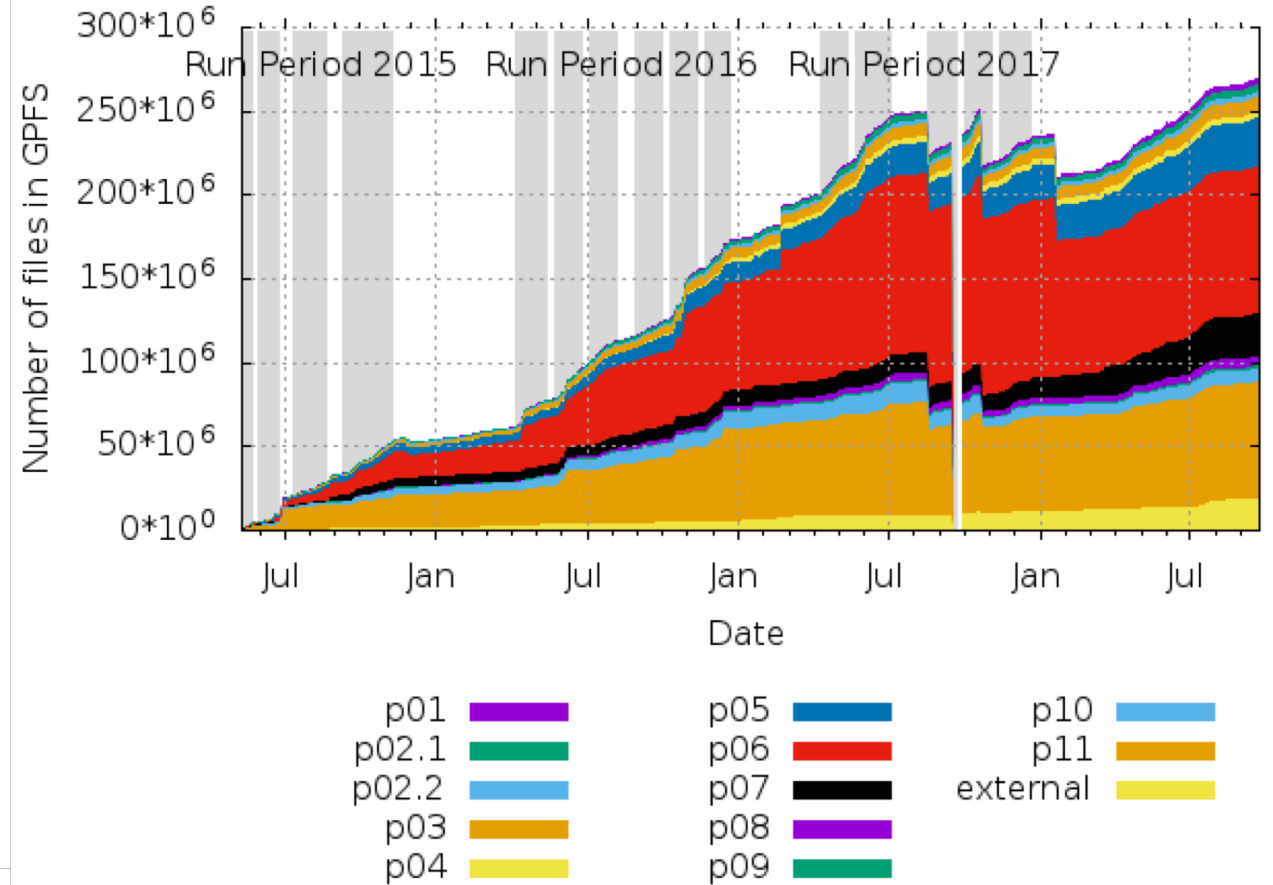
+HiDRA



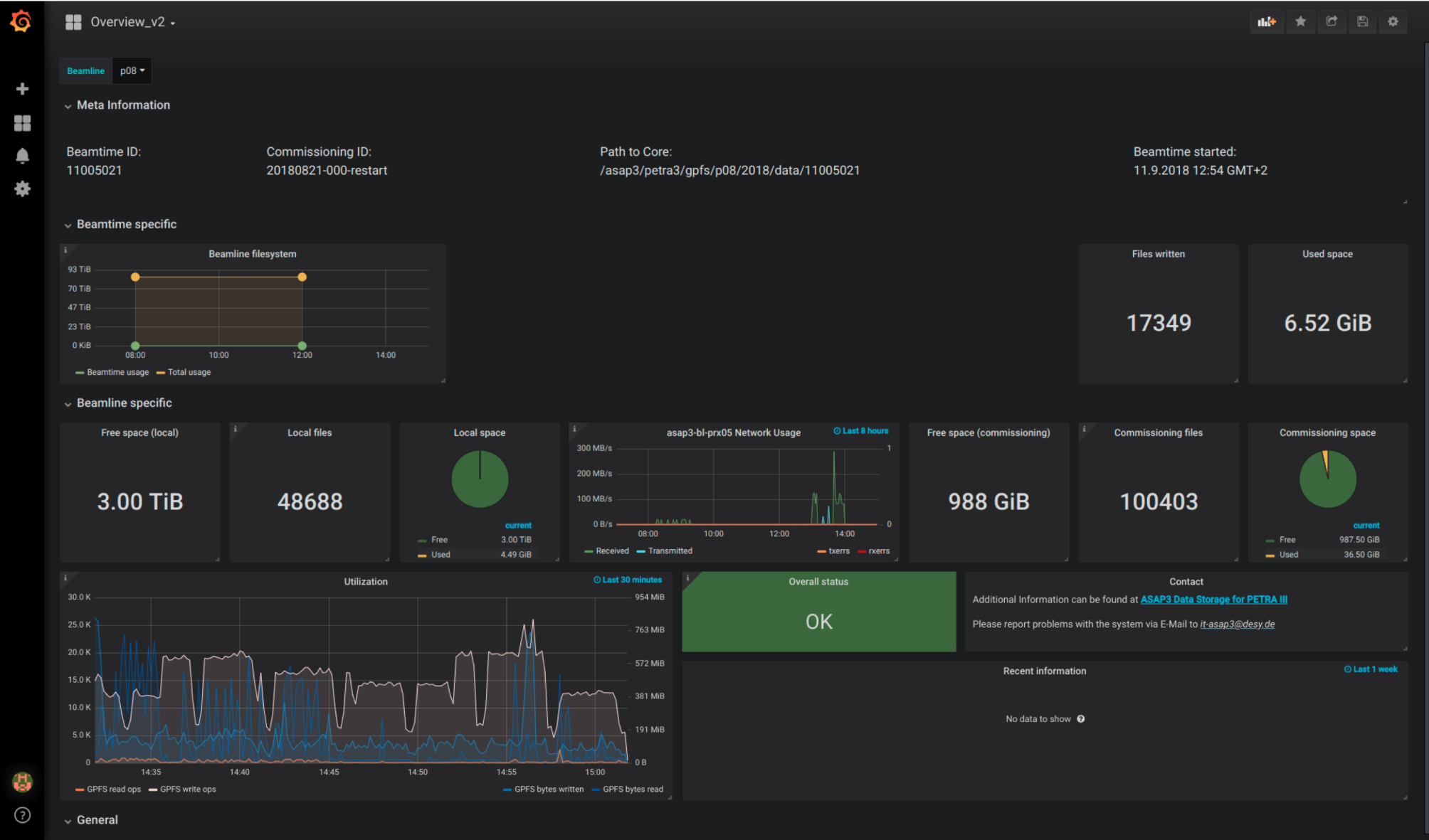
Storage consumption in size (per Beamline)



Storage consumption in number of files (per Beamline)

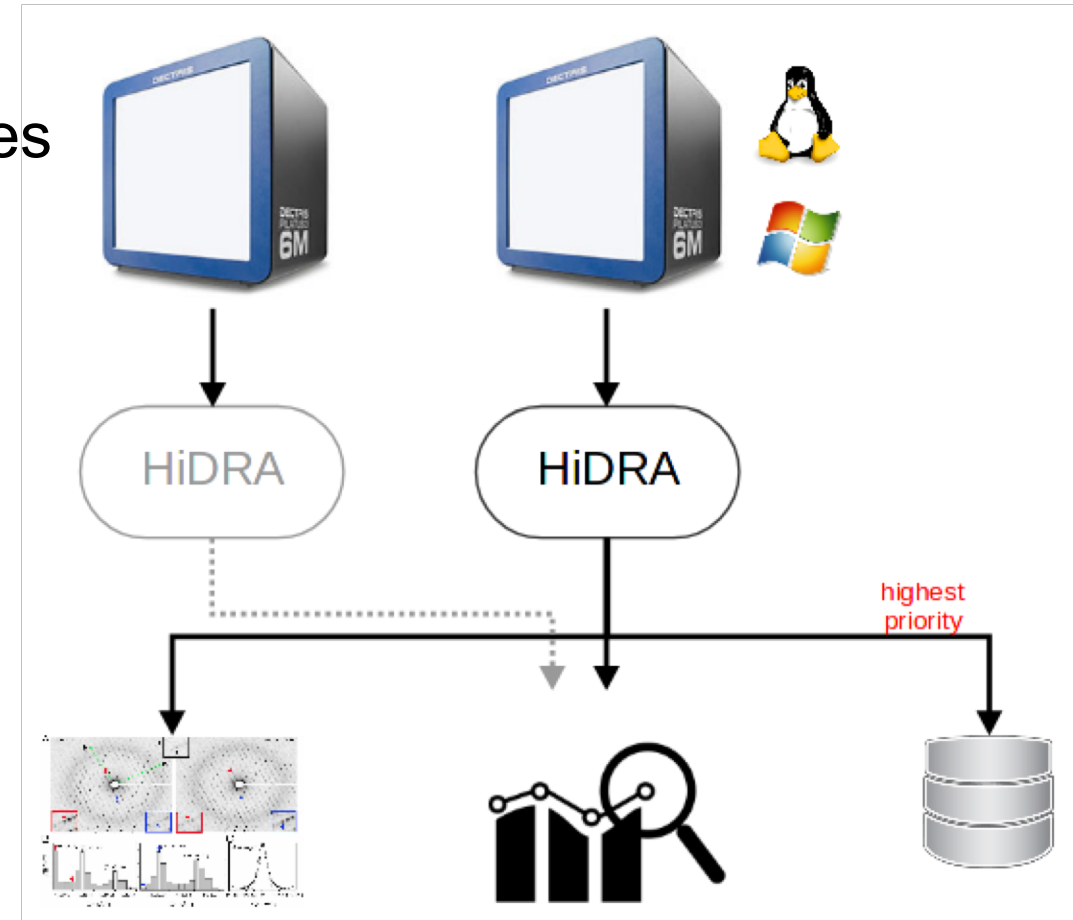


Dashboard – per beamline



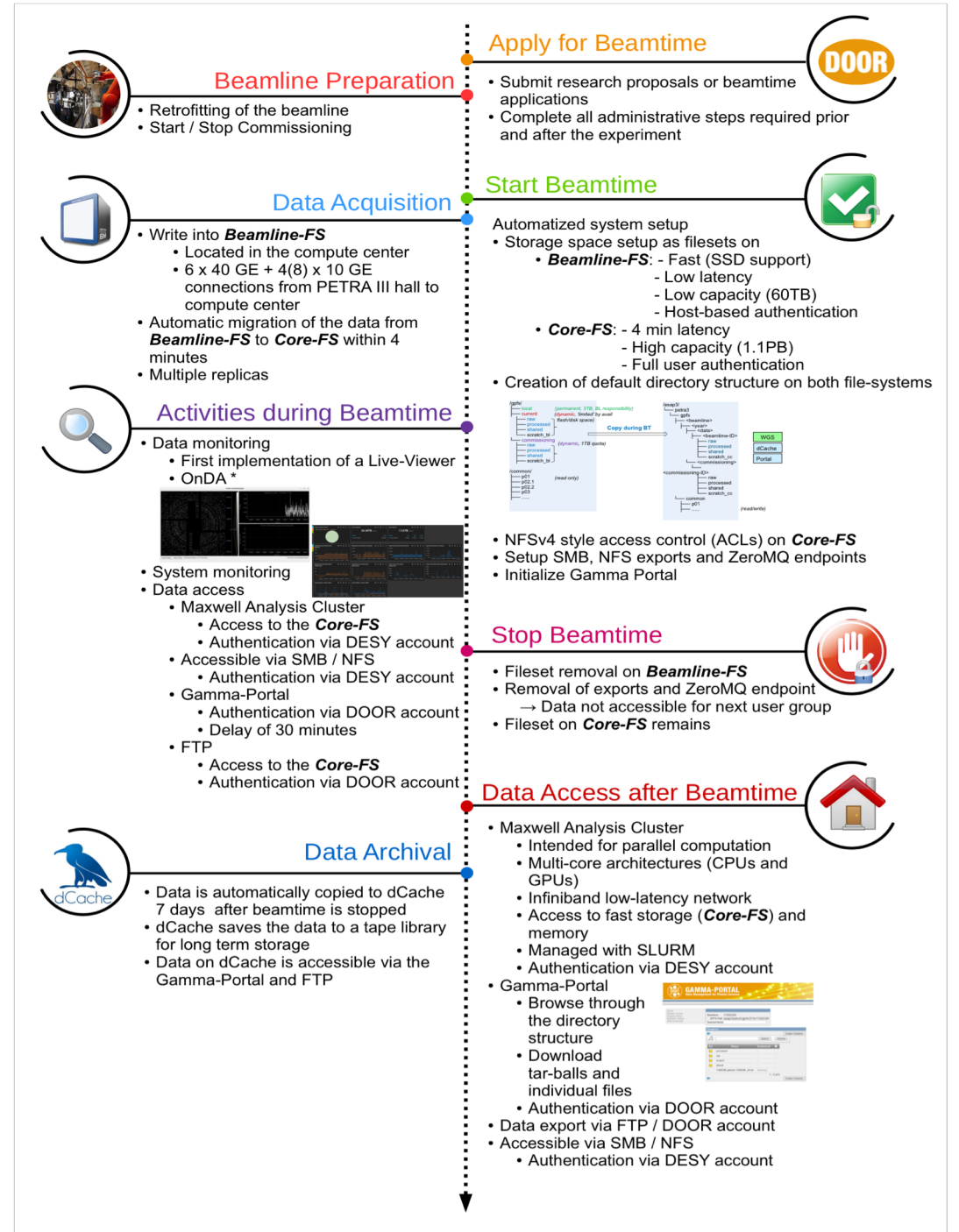
HiDRA

- High Data Rate Access
- Generic tool for high speed data multiplexing based on Python and ZeroMQ
- Actively used by FLASH and Petra 3 beamlines
 - data transfer from detector to GPFS
 - online monitoring & analysis
 - i.e. CFELs ONDA
- next generation in development
 - online & offline access – one API
 - query – include user key/value
 - scalable, more efficient

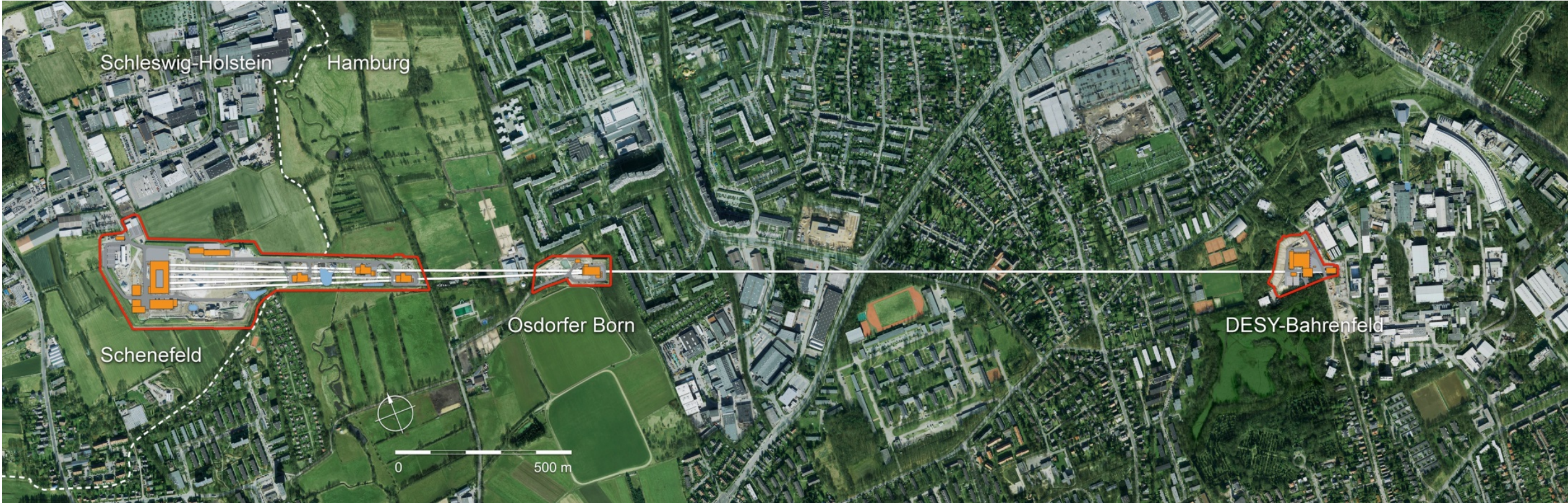


ASAP³ – glued services

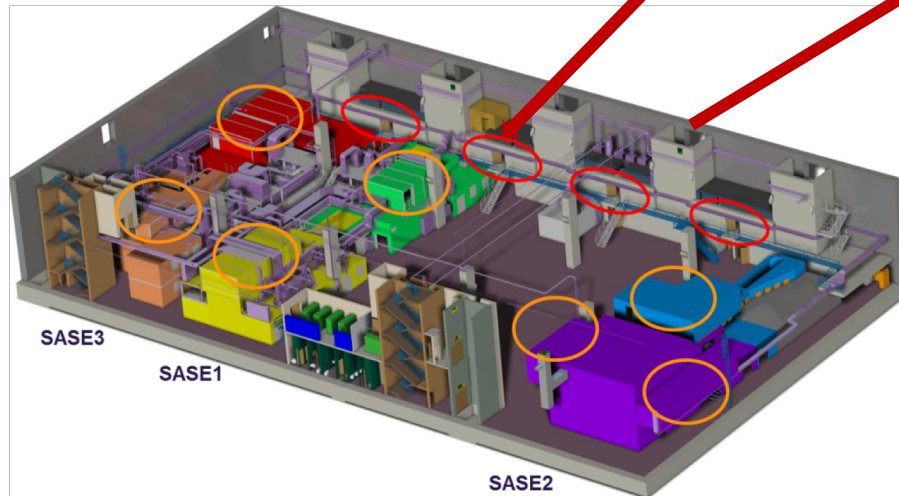
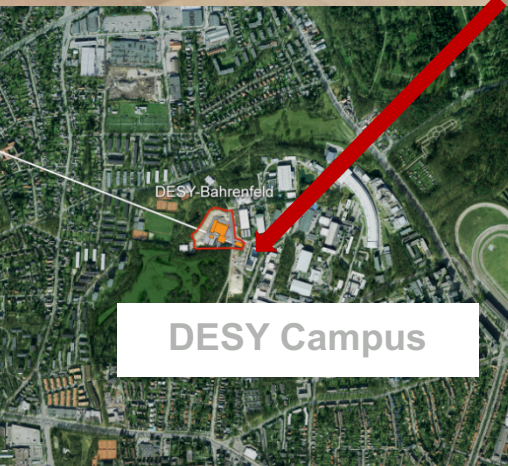
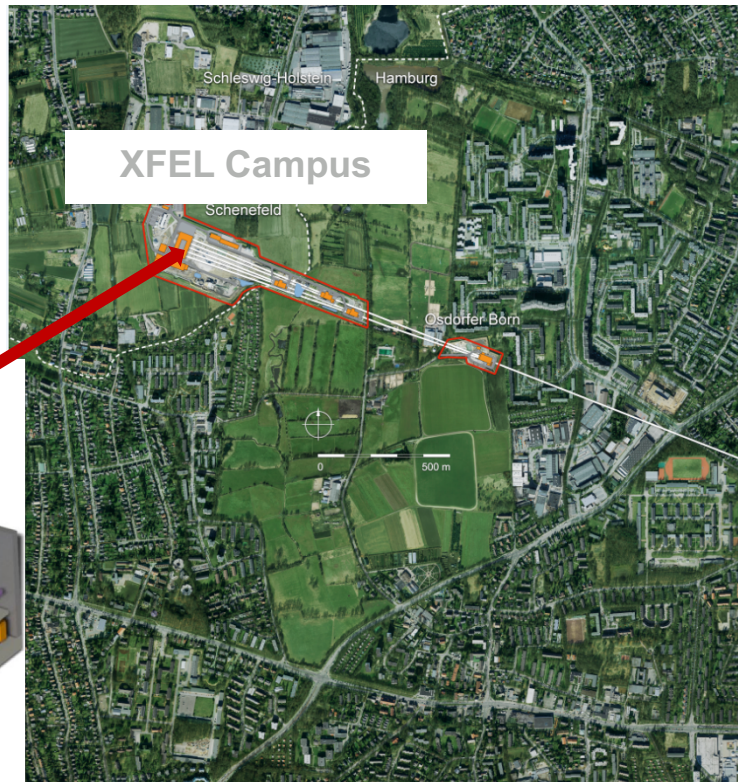
- GPFS
 - policy scans, XATTR, ACLs, multi-protocol access
 - GNR, end to end checksumming
 - fast – although it's a POSIX filesystem – see CORAL
 - startbeamtime/stopbeamtime preparation – automatic setup
- DOOR – beamtime metadata
- Gamma Portal – manage your beamtime data (visibility, ACL, ...)
- external access – FTP
- dCache & Tape copy
- HiDRA – online data analysis/live viewer, cover 'first mile'
- dashboard – overview per beamline
- home build InfiniBand monitoring
 - clean and running network – first on the list ;-)
- many scripts 'glueing' – mostly python



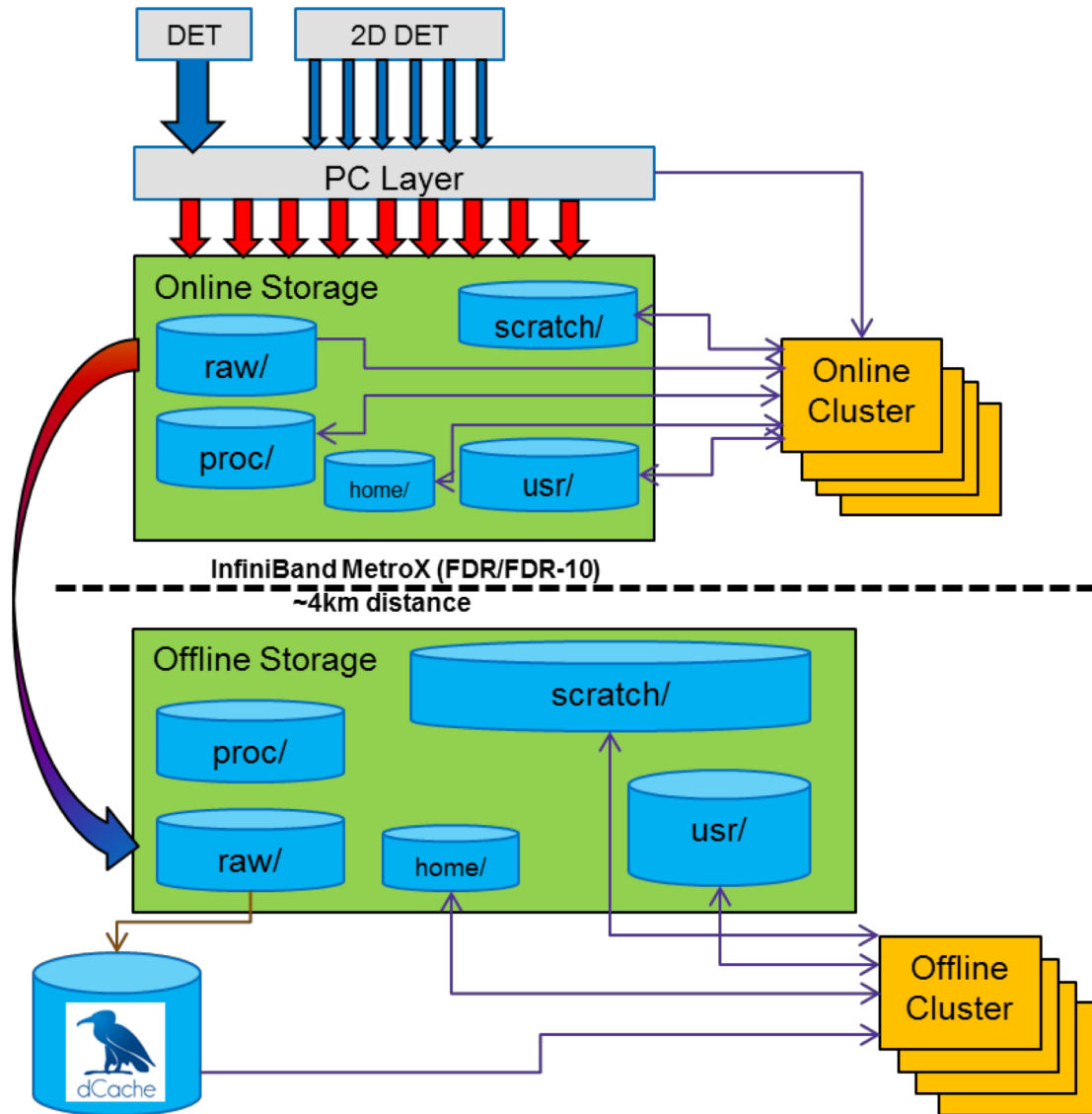
EuXFEL – a new research facility/instrument / September 2017



IT infrastructure at XFEL and DESY



DAQ & Offline architecture



- PC Layer
 - Pixel reordering
 - Data reduction, FPGA based compression
 - Veto
 - File creation on the online storage (HDF5)

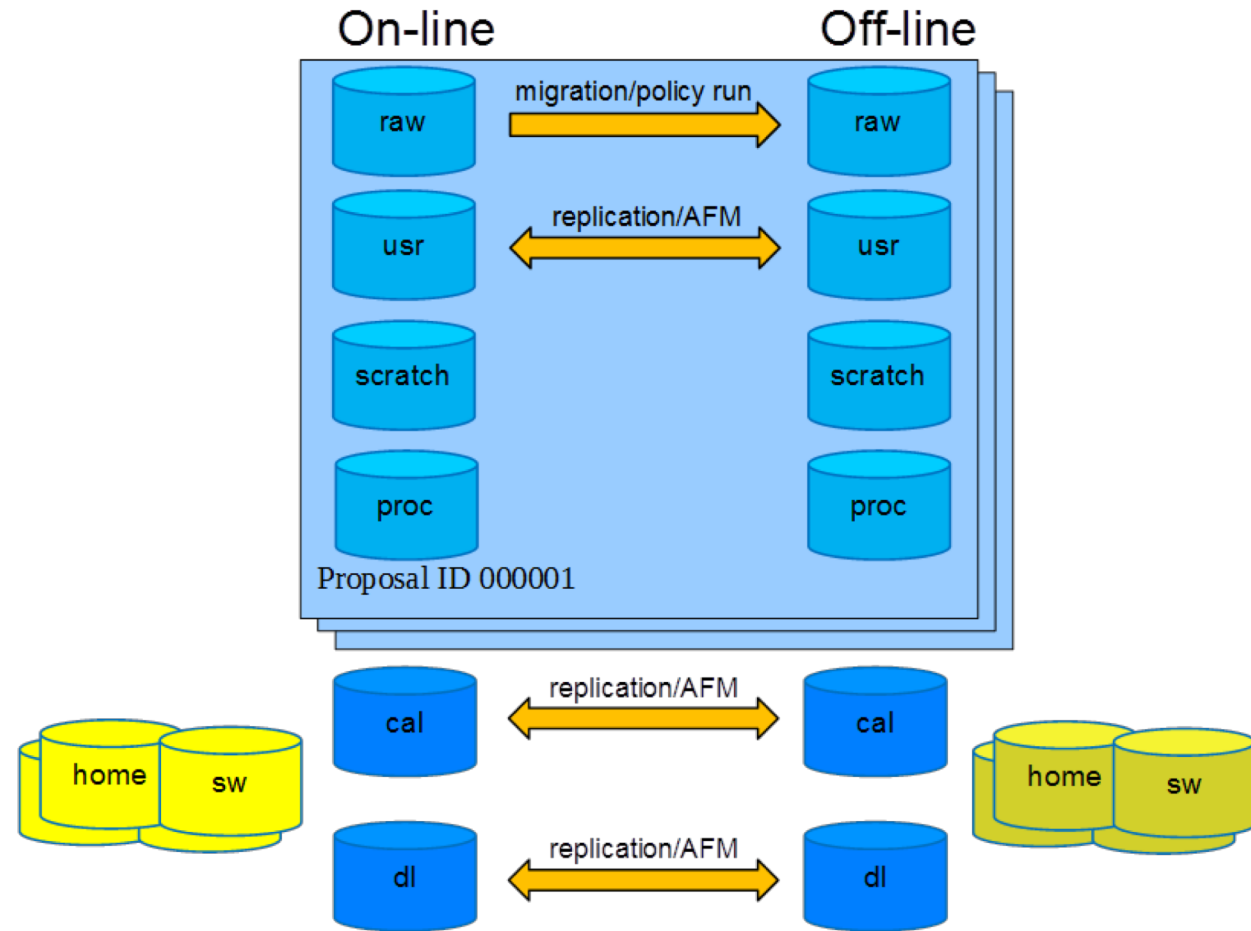
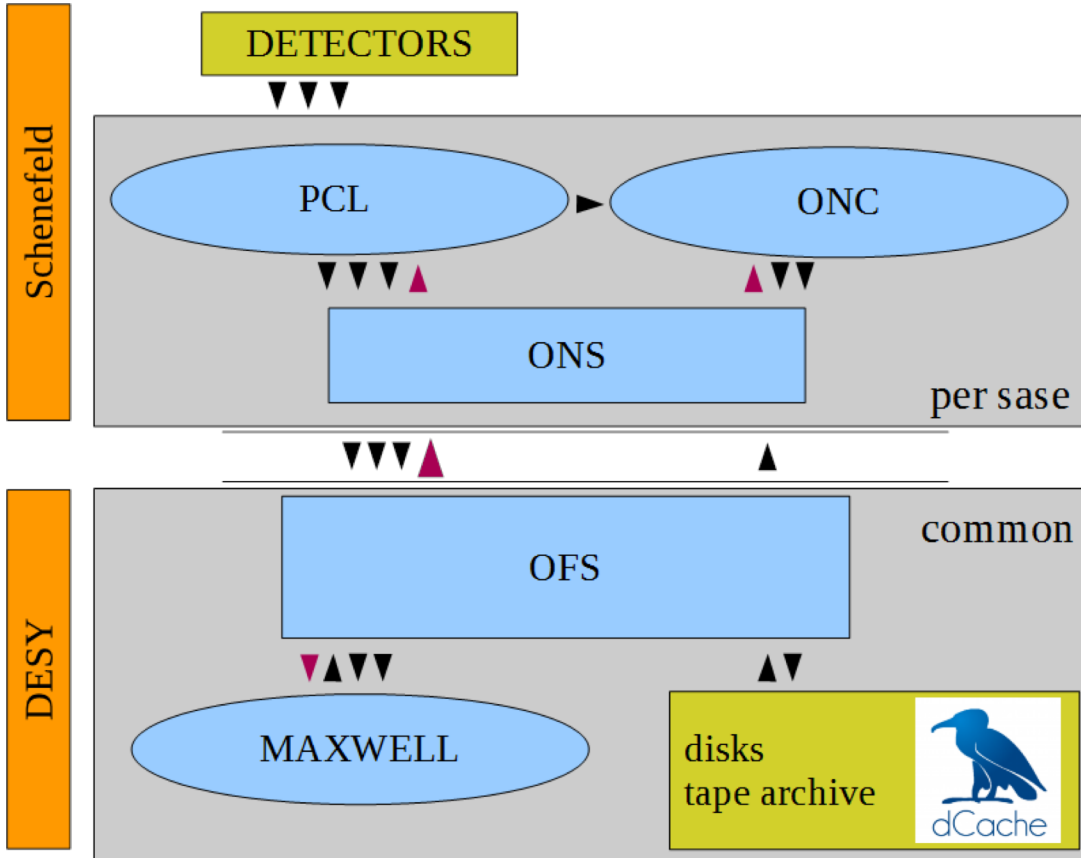
- Online cluster
 - CPU and GPU nodes
 - Online data analysis (fast feedback)
 - Calibration and data correction

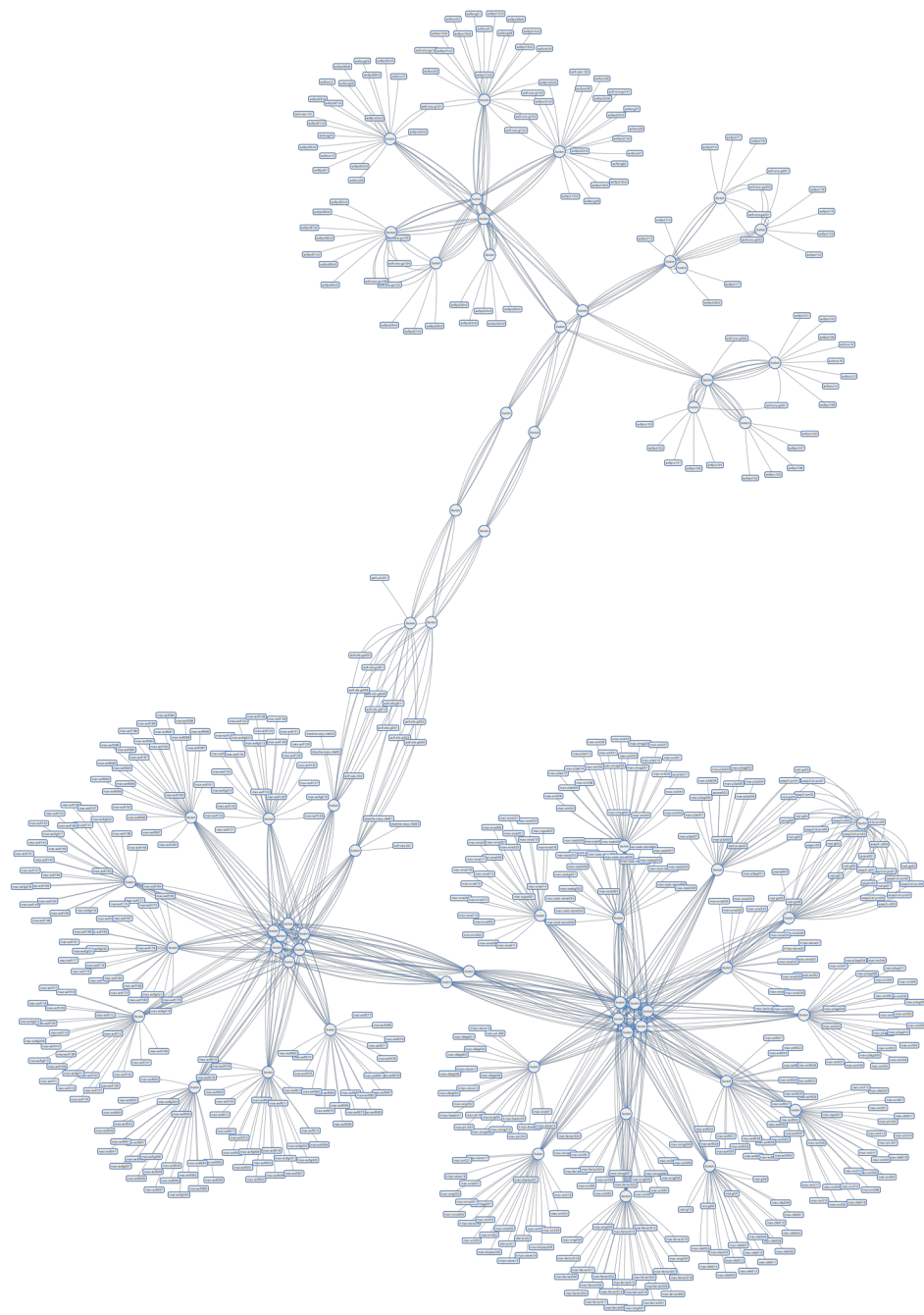
- Long-Haul Infiniband Metro-X
 - Supports 6 long haul ports (FDR-10 40Gb/s)
 - Scalable by multiple switches

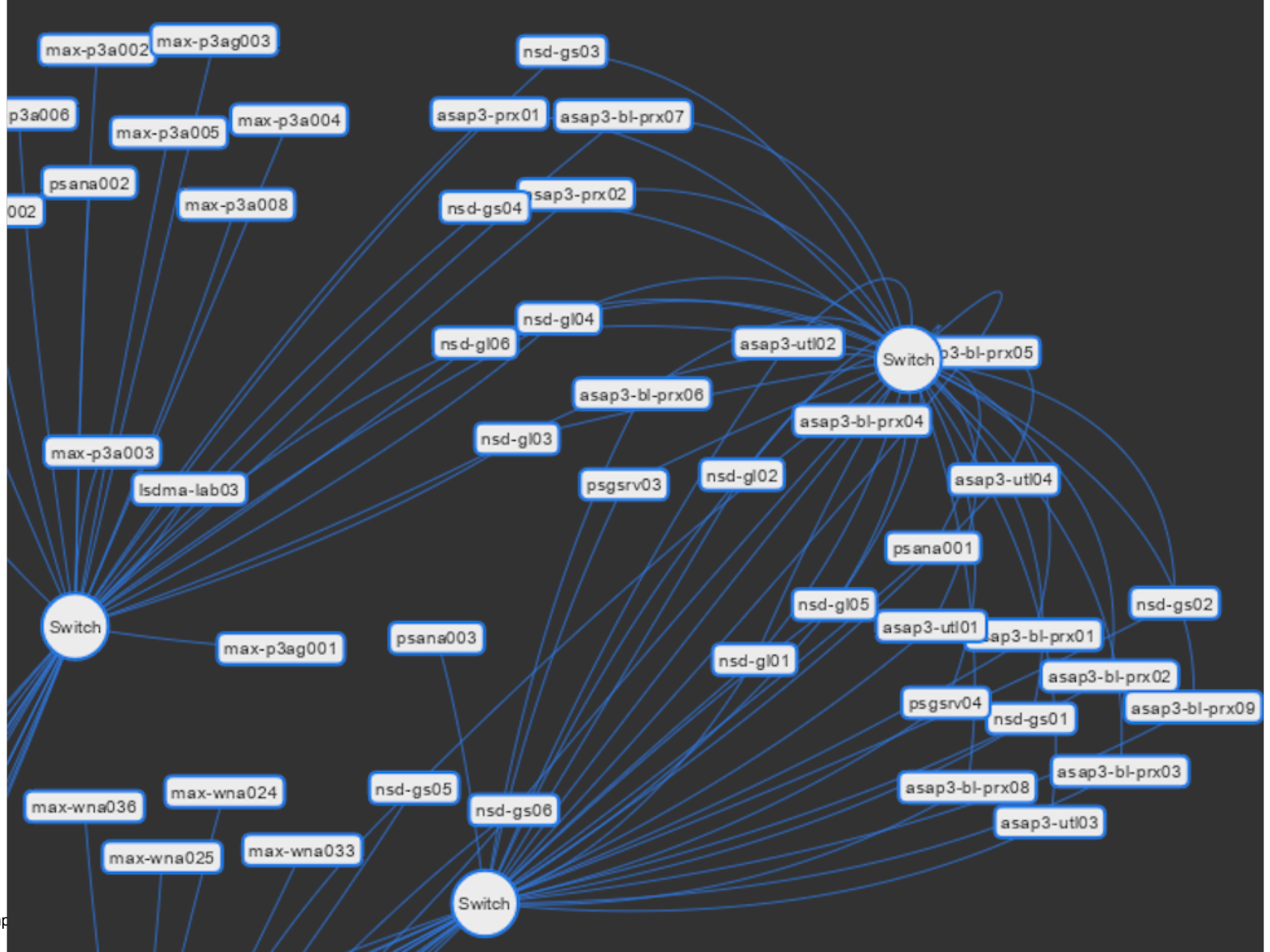
- Transfer online to offline storage
 - Custom scripts with policy runs
 - GPFS AFM

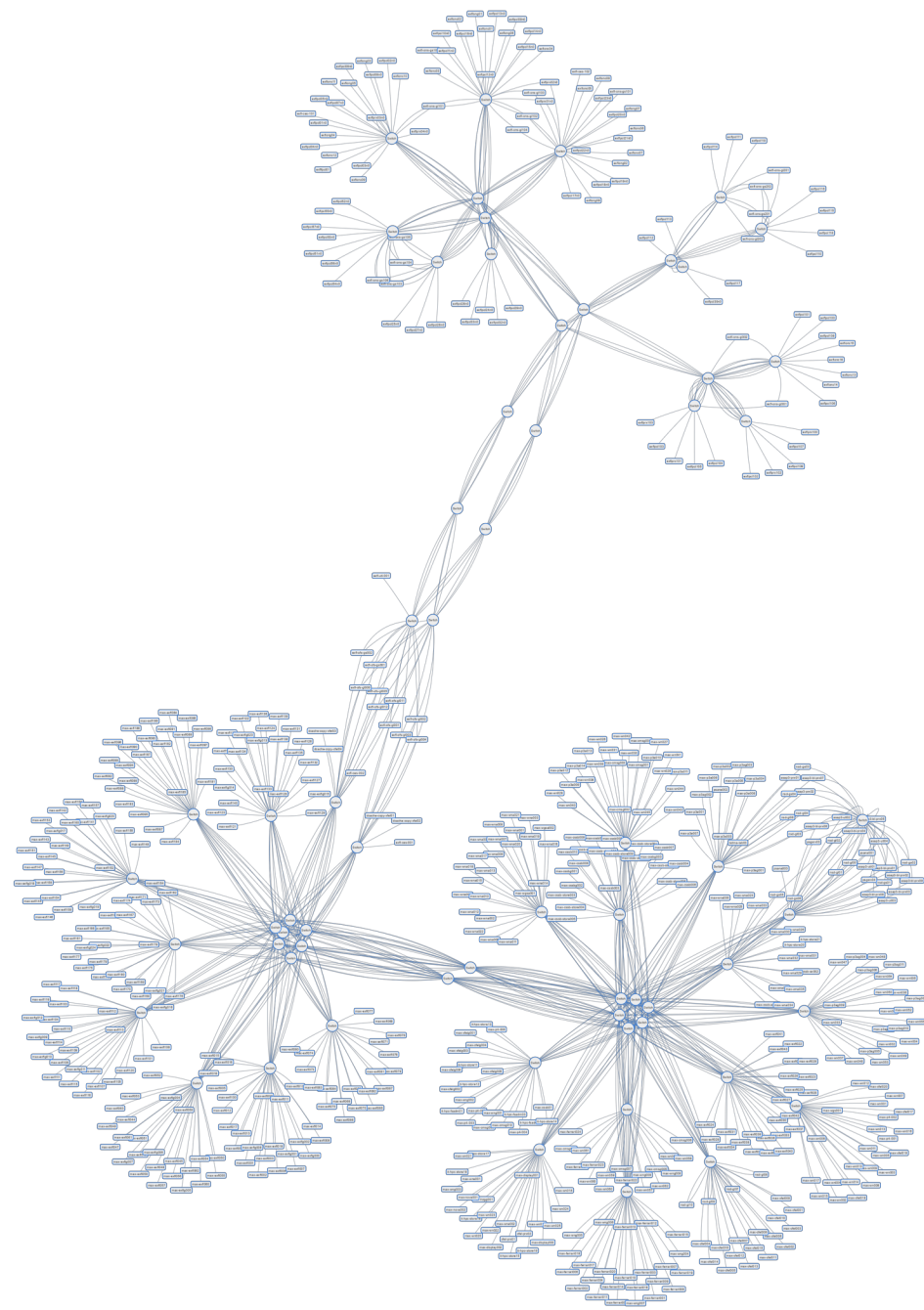
- Offline storage
 - Shared across experiment stations
 - Data migrated to offline storage after quality checks
 - Copy data to dCache, ACLs
 - Raw data access only from dCache
 - Calibrated data stored on GPFS
 - User analysis based on calibrated data

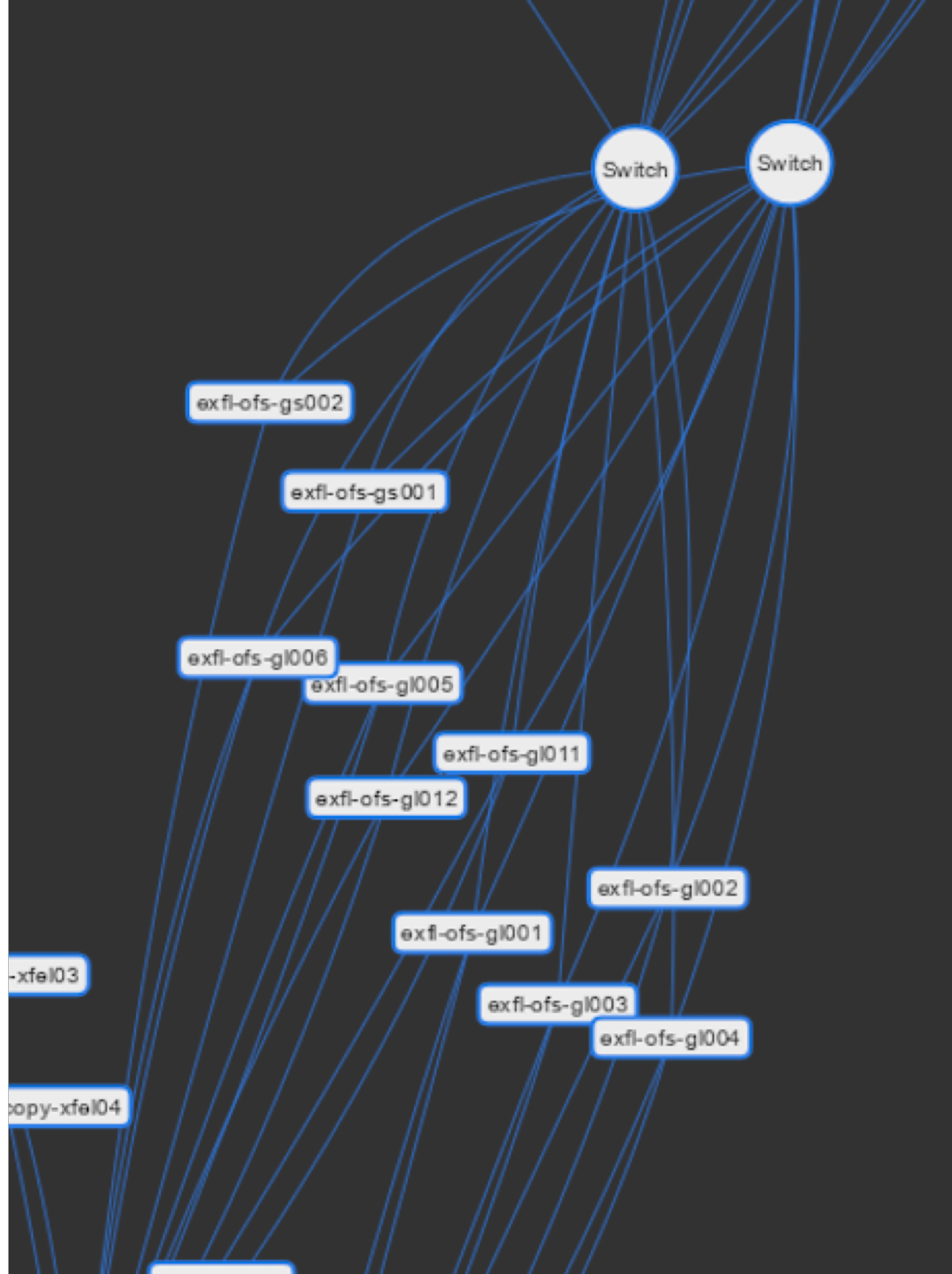
File systems and data placement





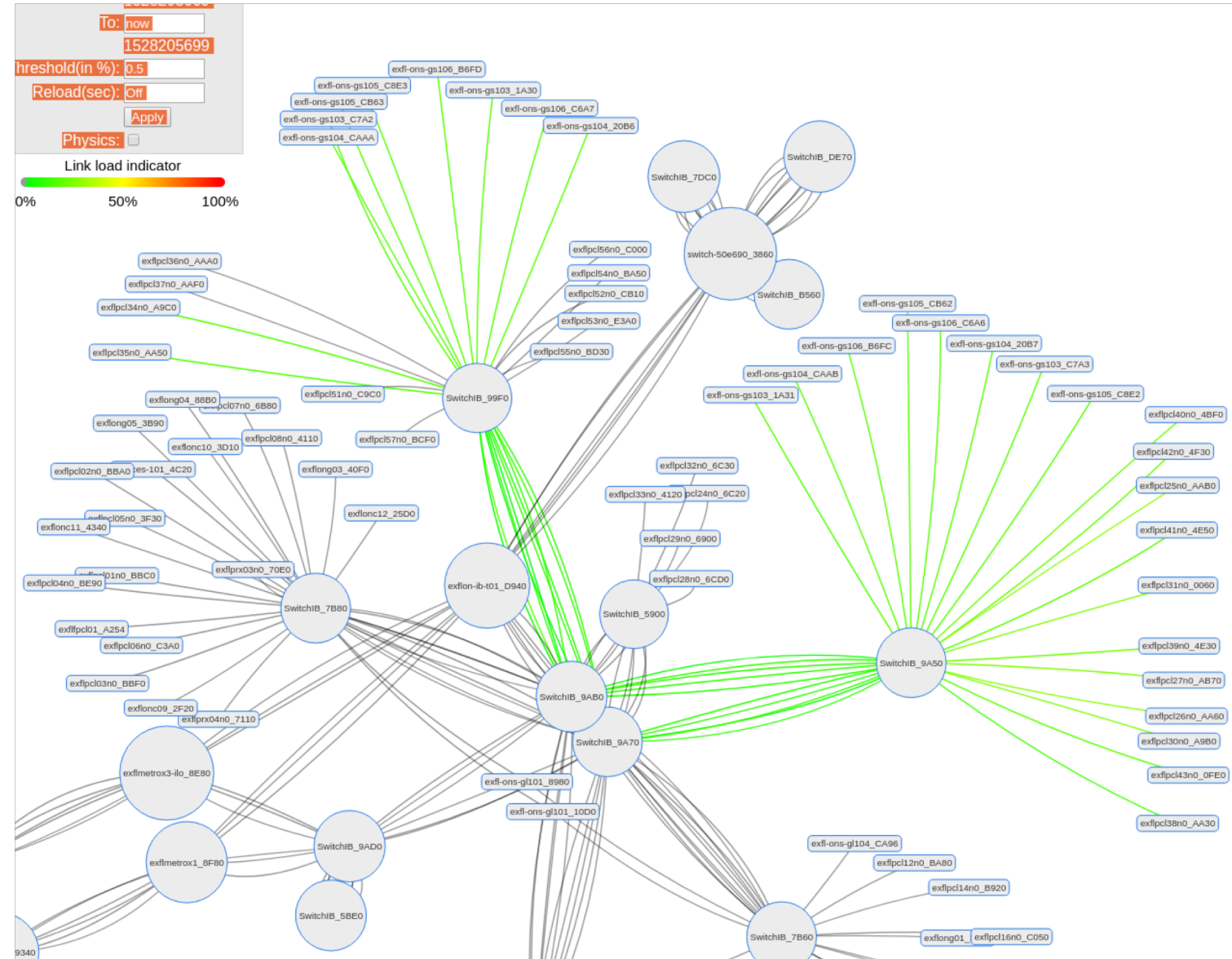




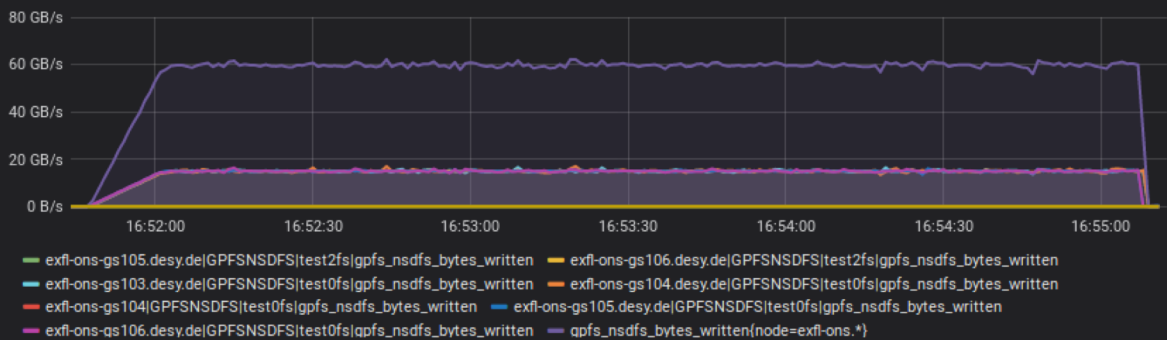


All Flash for DAQ writes (online storage)

- GS4S
 - 4 x 24 x 4TB (SAS connected)
 - EDR (16 ports)
 - measured 30GBs per system
 - agnostic against 'crazy IO'
 - limited by Linux SCSI stack and PCIgen3
- Petra 3 & FLASH (ASAP³)
- EuXFEL (3x)
 - tests – EuXFEL – 15 nodes each writing 200 x 4GB file every second
 - jitter < 20%



Online storage - gpfs_nsdfs_bytes_written - testts

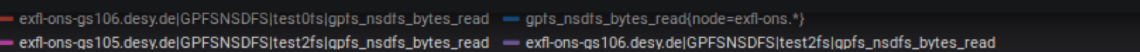


From:

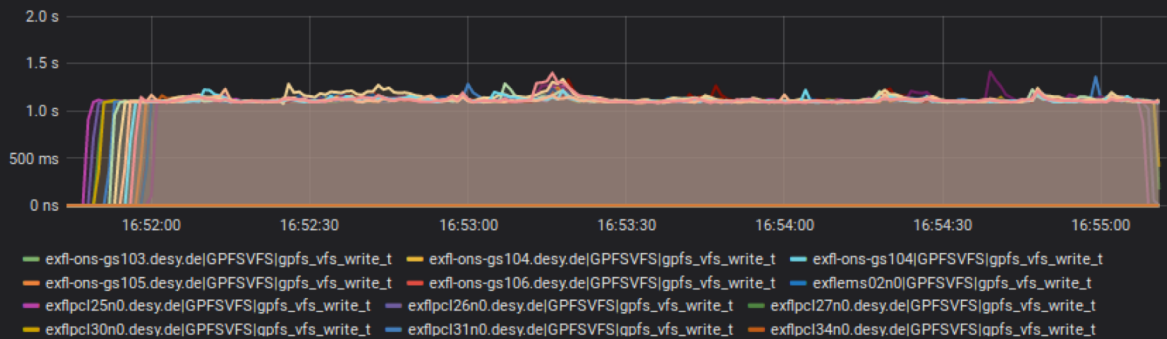
To:

Refreshing every:

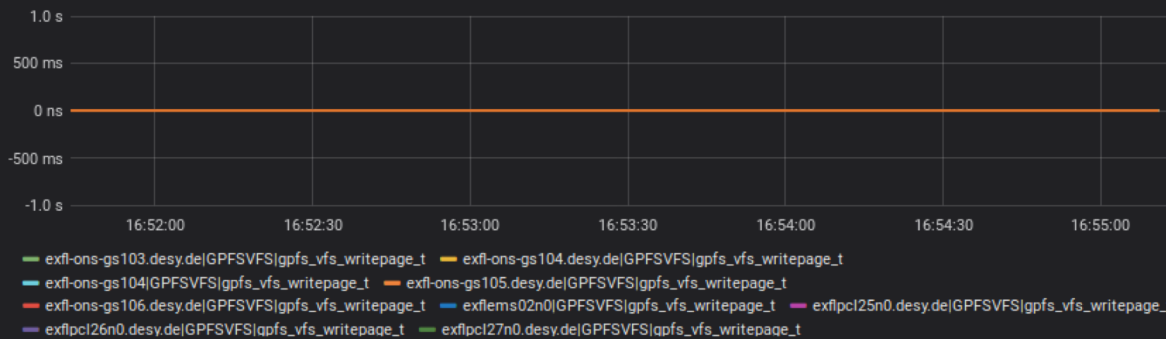
Last 2 days	Yesterday	Today	Last 5 minutes
Last 7 days	Day before yesterday	Today so far	Last 15 minutes
Last 30 days	This day last week	This week	<u>Last 30 minutes</u>
Last 90 days	Previous week	This week so far	Last 1 hour
Last 6 months	Previous month	This month	Last 3 hours
Last 1 year	Previous year	This month so far	Last 6 hours
Last 2 years		This year	Last 12 hours
Last 5 years		This year so far	Last 24 hours



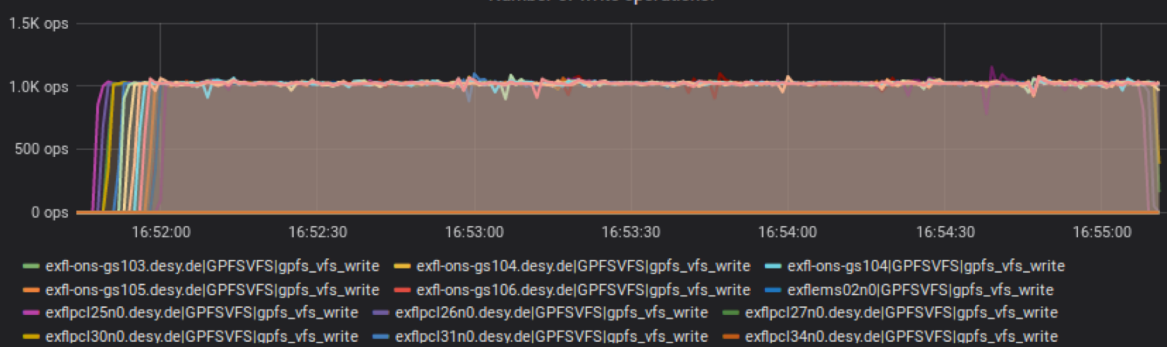
Amount of time in seconds spent in write operations.



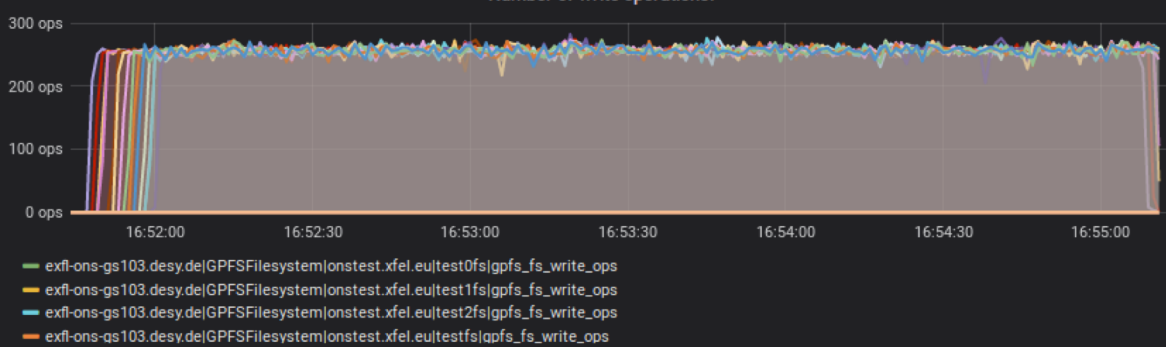
Amount of time in seconds spent in writepage operations.



Number of write operations.



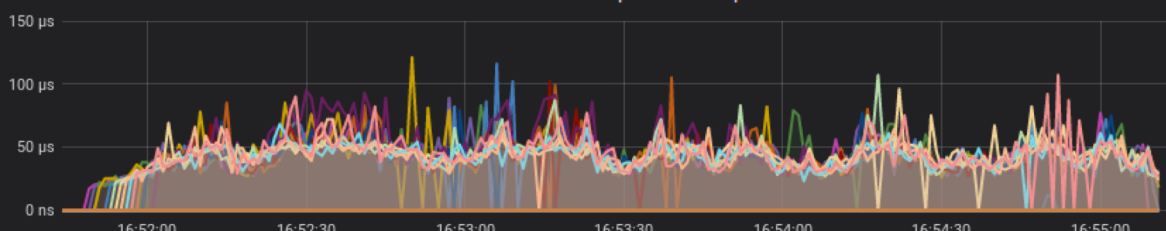
Number of write operations.



Amount of time in seconds spent in create operations.



Amount of time in seconds spent in close operations.



on the list

- GPFS events (cluster wide inotify)
 - events managed through Kafka
 - first tests completed
- GNR on network (MeStore)
 - promising lab results
 - smoother scaling at very high utilization
 - potential for – faster/cheaper burst buffer, capacity (disk) configurations
- logbook (digital)
 - everybody has a ‘not well beloved’ system, looking for better alternative
 - everybody acknowledge importance
- easier and scalable ‘online data analysis’ services and configurations

observations – end of 2017

primarily from IT perspective

- doubling detector rates faster than Moore's law (8-20 months)
- more than 1GBs from beamline to storage (faster + more detectors)
 - NFS & SMB ruled out
- online data analysis (aka. fast-feedback-loop)
- primary problem at the origin – getting data out of detector server memory fast enough (continuous mode)
- more Petra III Extensions in operations & FLASH in the game
- first experiments with 'Online Analysis' – hand crafted setup (Maxwell nodes access beamline FS)
 - could not simply be established as regular service – IO resource overrun, authentication/authorization, ...
 - docker + volume driver + ...
- load increasing / beamline & core FS
- connection/cooperation to other labs - 'capable of development' – too weak
 - standardization, common practices and developments, ...
 - CHEP and HEPiX like conference/workshop missing
- synergies with XFEL systems (long-range InfiniBand, all-flash systems, ...)

HiDRA etc.

summary

- stable and performant operation of GPFS
- profit from CORAL developments
 - less locks
 - buffer & rpc code
 - NUMA awareness (network & gpfs & block device threads)
- massive scale of new X GBs detectors – everybody wants that
 - several 10 of these – a little problem ;-)
 - data reduction (or ‘find the good ones’) becomes most important
 - Petra 4 on the horizon – heavy planning activities – next big jump in data rate (another x1k)
- PS is different than HEP – in many ways - computing too ;-)