

Back-end and Data Storage – DESY

New Concepts in Ultra Fast Data Acquisition – PSI April 10th 2018

Steve Aplin CFEL DESY





DESY : A History of Big Data

Serving Multiple Clients

HEP + Synchrotrons + FEL's : Accelerators

10 Hz ↔ 100 MHz

Kilobytes *+* Petabytes

Kb/s \leftrightarrow Tb/s

Serving Multiple Clients



Serving Multiple Clients



Serving Multiple Clients



Storage consumption in size (per Beamline)

4 Beamlines generate 80 % of the data on GPFS

Serving Multiple Clients



Serving Multiple Clients



Large Scale Strategic Infrastructure Over Island Solutions







Running a Service, Not Simply the Hardware



PETRA III Data Flow – ASAP³

Running a Service, Not Simply the Hardware



PETRA III Data Flow – ASAP³ Beam-line File System

- "Wild West"
- Host based authentication
- Access through NFSv3, SMB, or ZeroMQ
- Optimised for performance
 - NFSv3: ~ 600 MB/s
 - SMB: ~ 300 600 MB/s
- Tiered Storage
 - Tier 0: SSD Burst Buffer 2 TB
 - Tier 1: 200 TB Capacity

Running a Service, Not Simply the Hardware



PETRA III Data Flow – ASAP³ Core File System

- "Clean World"
- Full user authentication
- NFSv4 ACLs
- Access through NFSv4, SMB, or native GPFS
- GPFS Policy Runs copy data
 - Beamline Core Filesystem
 - Single UID/GID
 - ACL inheritance activated
 - Raw data classified as immutable

\rightarrow

• Single File-set per beamtime

User Adoption

Storage consumption in number of files per beam-line



Storage consumption in size per beam-line



Discrete drops in data, represent removal of previously archived data from disk based storage



| UP | | |
|--------------------|--------------|--|
| scan\$var.serialno | NXentry |] |
| OUP | | |
| P03 | NXinstrument | |
| ROUP | | |
| Pilatus1M | NXdetector | |
| FIELD | | |
| x_pixel_size | NX_FLOAT64 | 217 |
| FIELD | | |
| y_pixel_size | NX_FLOAT64 | 217 |
| FIELD | | |
| layout | NX_CHAR | area |
| FIELD | | |
| description | NX_CHAR | Pilatus 1M |
| FIELD | | STRATEGY |
| data | NX_UINT32 | POSTRUN /data/p03/2013B/xxyyzz |
| | | ATTRIBUTE |
| FileDir | NX_CHAR | \$datasources.P1_fileDir FINAL |
| | | ATTRIBUTE |
| FilePostfix | NX_CHAR | \$datasources.P1_filePostfix FINAL |
| | | |
| - • • • | | |
| cianal | | ATTRIBUTE |
| signal | | 1 |
| GROUP | | |
| extra_info | NXcollection | |
| | | CTRATECH |
| delay time | NX FLOAT64 | statasources P1 delavTime |
| delay_time | NA_I LOAI 04 | |
| - • • • | | |
| FIELD | | STRATEGY |
| nb_exposures | NX_FLOAT64 | \$datasources.P1_nbExposures FINAL |
| OUP | | |
| data | NXdata | |
| | | - |
| data | | NXentry/NXinstrument/Pilatus1M:NXdetector/data |
| 00.00 | | , |

How to Store the Data

- NeXus: a common data format for neutron, x-ray, and muon science
- NeXus file is a collection of components
 - Primitive devices: counter, MCAs, 2D detectors
 - Composite devices: monochromator, diffractometer
 - Structures: beam-line description
 - Can be inserted and removed without side effects
 - Have the strategy attribute: INIT, STEP, FINAL, POST
 - Refer to data sources: counters, MCAs, etc.
- Data Sources
 - Tango Server attributes (read from hardware)
 - JSON strings from the control client
 - Python script output
 - DB queries
- Integration tasks
 - Create a library of template components
 - Instantiate components by supplying local data source names
 - Integrate this scheme into the experiment control clients



How to Share the Data

HASYNAB wse archive > Users(ACLs) List Users Home Browse archive Period: from 21-APR-11 to 27-APR-11 Staging status Beamtime: 10005888 Migration status Access Lastname Firstname User Role 🔫 Create Date Change Date ccoun Permissio Enju 0 elima Lima leader download 27-FEB-2013 15:13 27-FEB-2013 15:13 Pernot Petra participant 27-FEB-2013 15:13 27-FEB-2013 15:13 pernot 0 download Wiegart Lutz participant 🗘 download 27-FEB-2013 15:13 27-FEB-2013 15:13 wiegart 1 - 3 Delete Cancel Save ADD DOOR USER TO EXPERIMENT FOR DATA ACCESS Door account Add User ACL MANAGEMENT Update Acls Show Acls

- Gamma Portal provides users with an interface to manage and access data
- Uses DOOR credentials
- Browse data: file discovery, searchable catalog
- Data access management for PI's and Data-Managers



How we keep all the data



- DESY Data Policy
 - The PI has the full responsibility of the data
 - The PI may grant data access to other persons
 - DESY offers to store the data over a complete data life cycle (10 year) at the expense of the PIs institution.



- PaNdata Data Policy
 - The facility acts as the custodian of the data
 - Data are open access after an embargo period of 3 years, can be extended

PETRA III/IV Future Developments

Paradigm Change for User Experiments:

Increasingly complicated experiments:

- Users do not have resources or might not be experienced enough to process data by themselves
- → Data analysis becomes integral part of experiment



→ DESY needs extra resources for data analysis as service to users!



Christian G. Schroer | Workshop "Digitale Agenda" | February 23, 2018 | page 22

Data Science and Scientific Computing requirements

→ Growing demand for measuring local properties of heterogeneous samples and following their evolution in-situ/operando → space and time resolved

Data science:

> growth of data volume (next 10 years):

100 x - 1000 x increase in brilliance (PETRA IV)5 x more beamlines taking large data sets10 x automation of data acquisition

 \rightarrow 10⁴ - 10⁵ x increased data rate

Scientific computing:

> today: data evaluation (pure number crunching time) is factor 10 - 100 too slow

> algorithms typically scale with (data size)^{α} with $\alpha \ge 1$, e. g. tomography: ~ size log(size)

→ new perspective on "raw data" (do not keep raw data: reconstructed data is "raw")

- → clever algorithmic developments needed (> 10⁵ x increased computational power/efficiency)
- → not solvable by brute force



Tried and Trusted Methods ... ?

How we keep all the data



Tried and Trusted Methods ... ?

How we keep all the data



"Do they look at the data as physicists ... or biologists?" Anonymous Detector Expert

"Do they look at the data as physicists ... or biologists?" Anonymous Detector Expert

A large part of the experimental infrastructure comes from experimental particle physics, e.g. the accelerators, the detector technology, people ...

It's easy to forget the science does not ...

Where's My Data?!?!

User Support is Part of the Service



Where's My Data?!?!

User Support is Part of the Service

| Online Storage Online Scratch FS Scratch | Fast 2D pixel detector data | P2P Online analysis | | | |
|---|--|---|--|--|--|
| 2. aplin@max-exfl0 | 01: /gpfs/exfel/exp/SP | B/201701/p002012/raw/r0101 (ssh) | | | |
| → [max-exfl001:aplin] ls -lh head total 235G | | | | | |
| -rr 1 xdata xdata 3.7G | Sep 18 01:12 | RAW-R0101-AGIPD00-S00000.h5 | | | |
| -rr 1 xdata xdata 3.7G | Sep 18 01:12 | RAW-R0101-AGIPD00-S00001.h5 | | | |
| -rr 1 xdata xdata 3.7G | Sep 18 01:13 | RAW-R0101-AGIPD00-S00002.h5 | | | |
| -rr 1 xdata xdata 3.7G | Sep 18 01:13 | RAW-R0101-AGIPD00-S00003.h5 | | | |
| -rr 1 xdata xdata 3.7G | Sep 18 01:12 | RAW-R0101-AGIPD01-S00000.h5 | | | |
| -rr 1 xdata xdata 3.7G | Sep 18 01:12 | RAW-R0101-AGIPD01-S00001.h5 | | | |
| -rr 1 xdata xdata 3.7G | Sep 18 01:13 | RAW-R0101-AGIPD01-S00002.h5 | | | |
| -rr 1 xdata xdata 3.7G | Sep 18 01:13 | RAW-R0101-AGIPD01-S00003.h5 | | | |
| -rr 1 xdata xdata 3.7G | Sep 18 01:12 | RAW-R0101-AGIPD02-S00000.h5 | | | |
| /gpfs/exfel/exp/SPB/201701/p002012/raw/r0101 | | | | | |
| → [max-exfl001:aplin] | | | | | |
| | Online monitoring Rapid feedback Rapid feedback Rapid feedback | Rapid feedback Scratch space | | | |
| | Not shown is technical infrastructure such as switches. | Raw data Calibrated data Calibration data | | | |

Alignment datasets are shipped with the data products and tools for

coordinate system conversion are provided by the facility.

PC-layer Trainbuilder-format Data Cal. constants Data Cal. constants

Tragedy of the Commons

we are all in this together

The **tragedy of the commons** is an economic theory of a situation within a shared-resource system where individual users acting independently according to their own self-interest behave contrary to the common good of all users by depleting or spoiling that resource through their collective action.





Massive increases in data rates at the experimental measurement stations has created the pressure to consolidate technical solutions into strategic infrastructure

At the same time the amount of data users are forced to handle has created the need for well defined data handling strategies and policies

What do we Need to spend our money on? What do we Want to spend our money on?

It's the End-to-End solution which counts at the end of the day.