Workshop on Raw Diffraction Data Reuse:
"The Good, the Bad and the Challenging"

# Handling big data at the European XFEL

Fabio Dall'Antonia, European XFEL Data Analysis Group
August 22, 2023

# Outline

- Introduction to the European XFEL

- SFX at the European XFEL

- The data challenge: storage space and reuse options

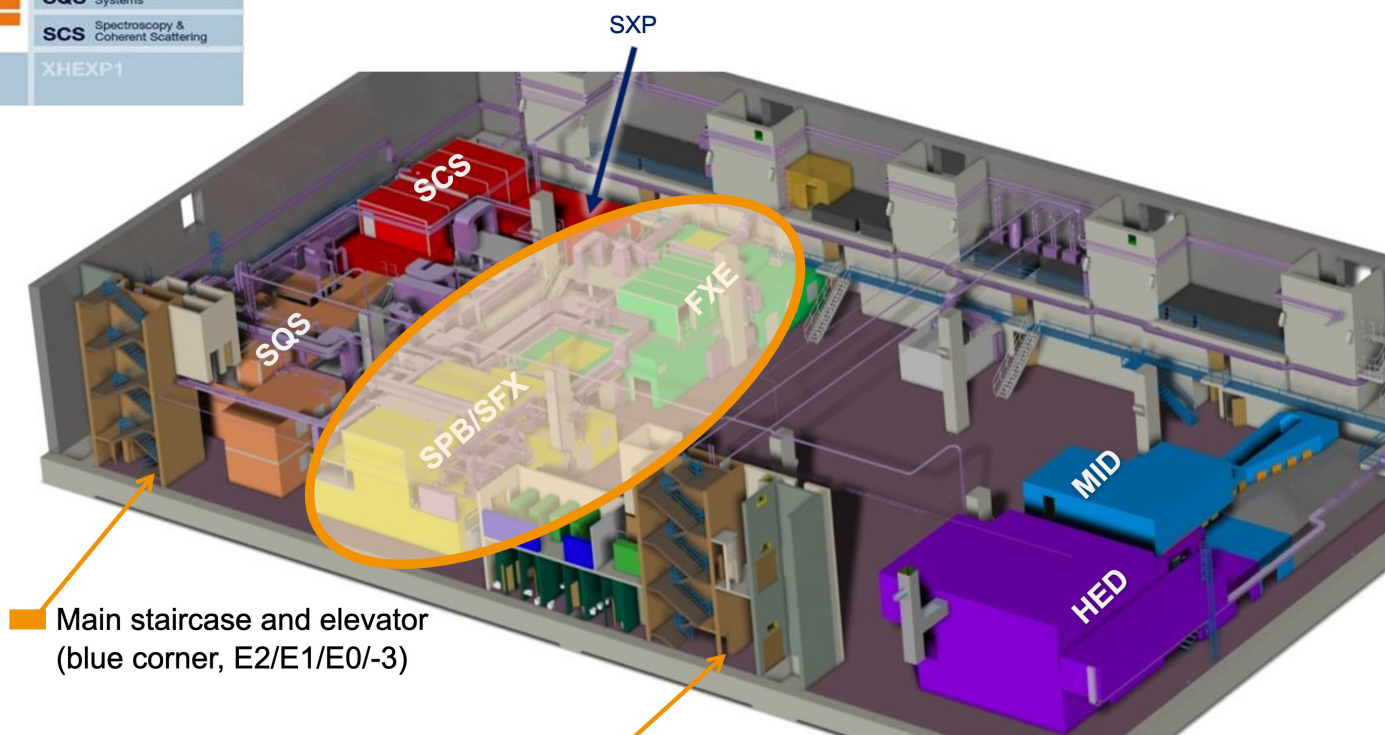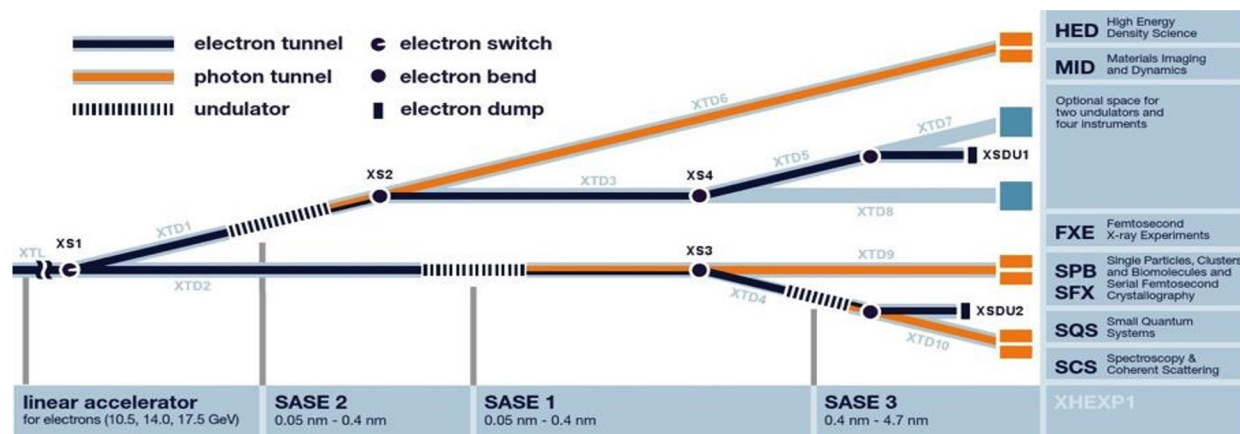- The way forward

**European XFEL**

# European XFEL at a glance

- Non-profit company, eleven shareholder countries

- Photon science user facility (plus own science), operation since September 2017
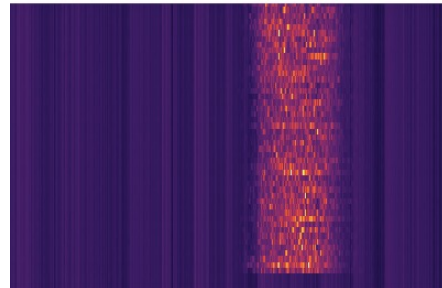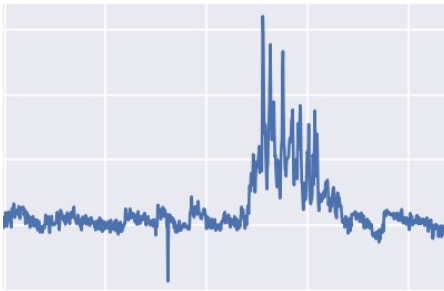
- currently > 500 employees





- linear electron accelerator, 3.4 km tunnel

- Seven scientific instruments

**European XFEL**

# Beamlines and instruments



SXP

Main staircase and elevator
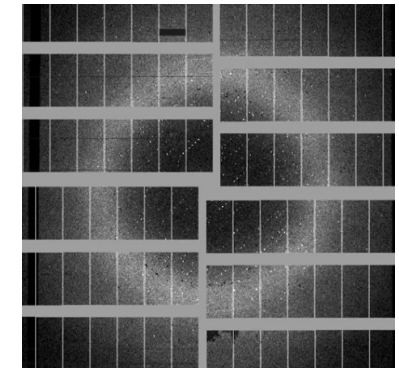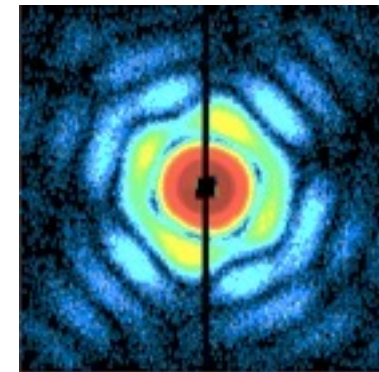(blue corner, E2/E1/E0/-3)

**European XFEL**
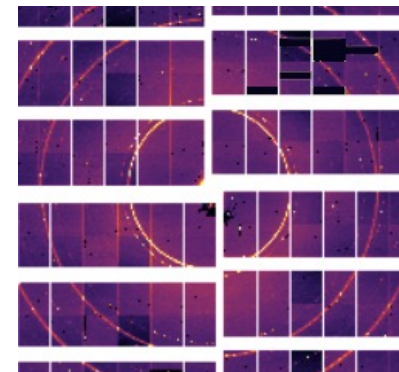
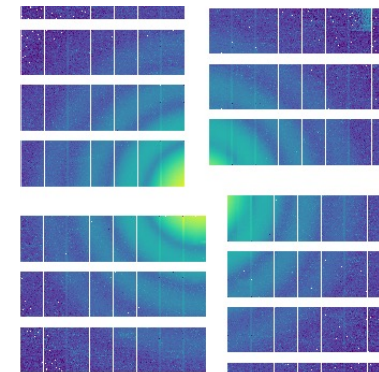# Major types of data

- Diffraction / scattering data: pixel / area detectors
  - integrating, mainly custom-built, multi-gain, MHz (burst mode)

- Spectroscopic data: 1D detectors, partly pulse-resolved (e. g. Gotthard)



[1]

- Digitizers, diagnostics devices X-ray-gas monitors etc., control data (motors, valves etc.)...

[1] taken from Kirkwood et al. (2022) Nat. Sci. Data 9

**European XFEL**

# Flow of experiment data

Real-time processing pipeline

Live monitoring

Control system



Detector (often adaptive gain)

DAQ

Informs operator, decision on acquisition

File-based data analysis

**European XFEL**

# SFX at the European XFEL

**European XFEL**

# The typical SFX setup at EuXFEL (SPB/SFX)
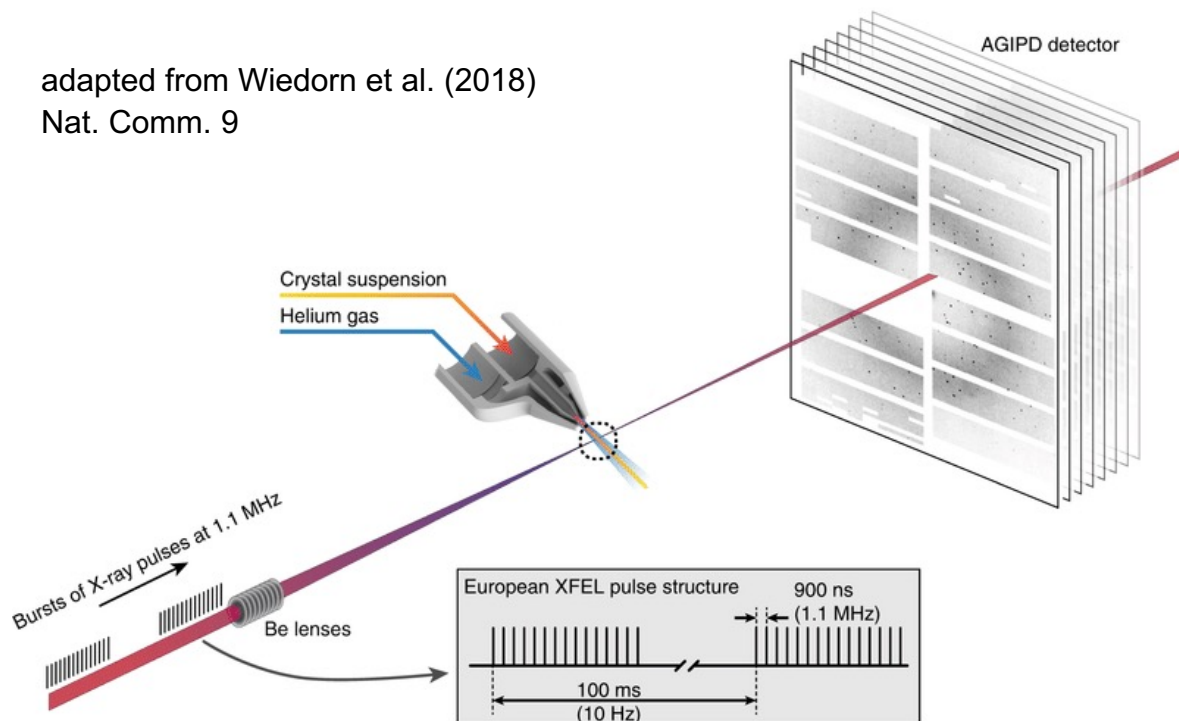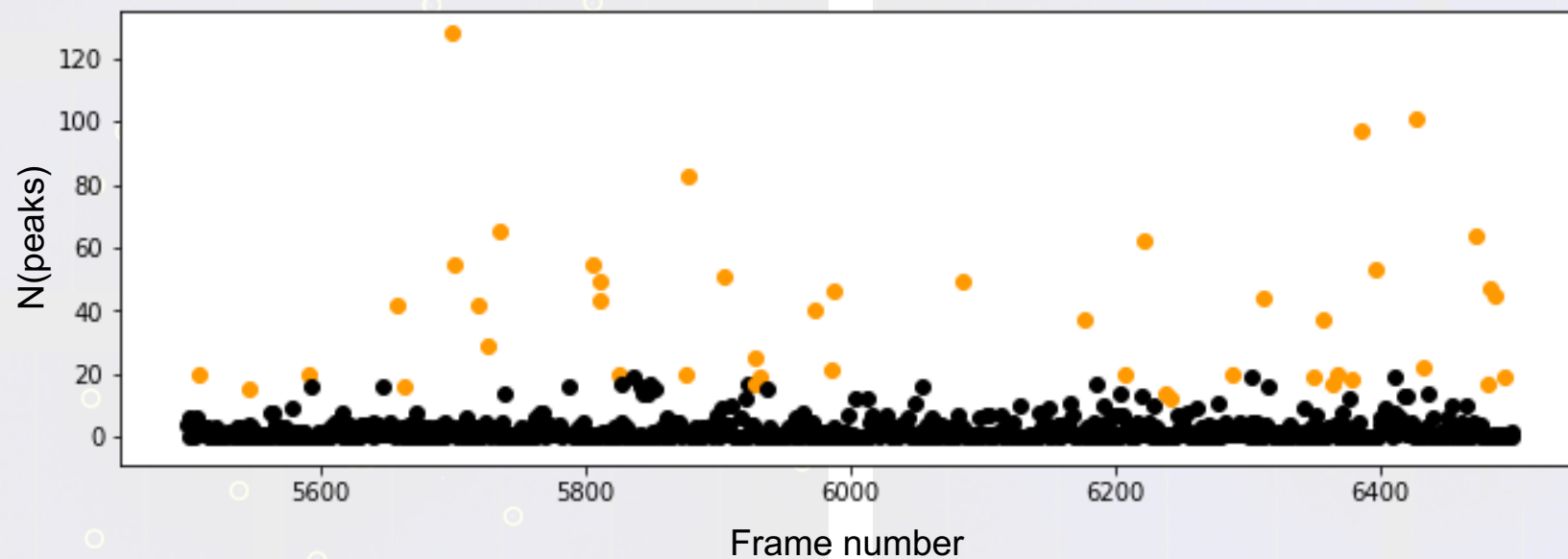
adapted from Wiedorn et al. (2018) Nat. Comm. 9



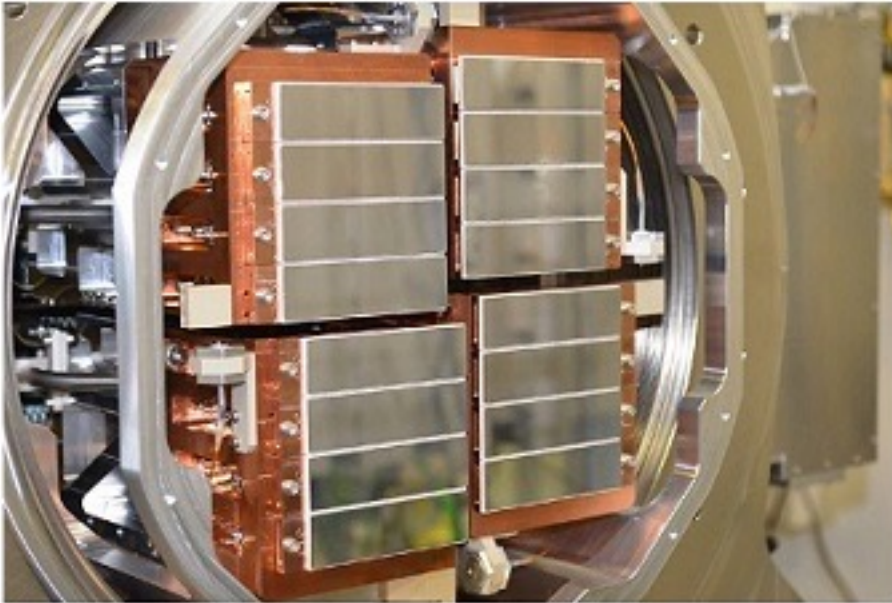▶ Diffraction before destruction

▶ Each image stems from a differently oriented crystal, indexing cannot make use of multi-frame information

▶ Requirement of high crystal isomorphism, each image must represent (closely) the same crystal "template"

▶ Still images, no oscillation range, integration must use a 3D model for 2D pixel data

▶ Many empty frames due to missing shots

**European XFEL**

# Crystal hits and misses



▶ Parameters, e. g. N(peaks) threshold have to be chosen carefully
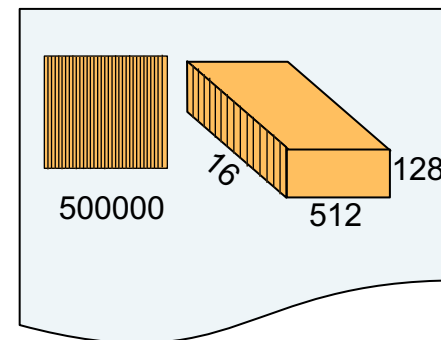
**European XFEL**

*SPB/SFX, AGIPD-1M*

# Detectors: multi-module topology



▶ AGIPD-1M detector, 16 modules, separate read-out

▶ four quadrants of four modules, motor-moveable

**European XFEL**

# Data format / shape



(N ×) 16 HDF5 files
in a run folder
(several TB in total)

single HDF5 file:
virtual data layout
(~ 50 MB)

Geometry description
(text file or HDF5, ... )

# Extra-Xwiz is a pipeline for data processing with CrystFEL

# The data challenge: storage space and re-use options

# File systems for data storage

## European XFEL Storage Overview

### Schenefeld

**Online GPFS**

- Cache

- Extremly high performance
- Data available immediately
- Optimised for concurrency
- High redundancy
- Dedicated storage for each SASE
- Very high cost per PB
- Capacity for a few days

**Offline GPFS**

- Performance

- High performance
- Large scale data analysis
- High redundancy
- High cost per PB
- Shared within XFEL
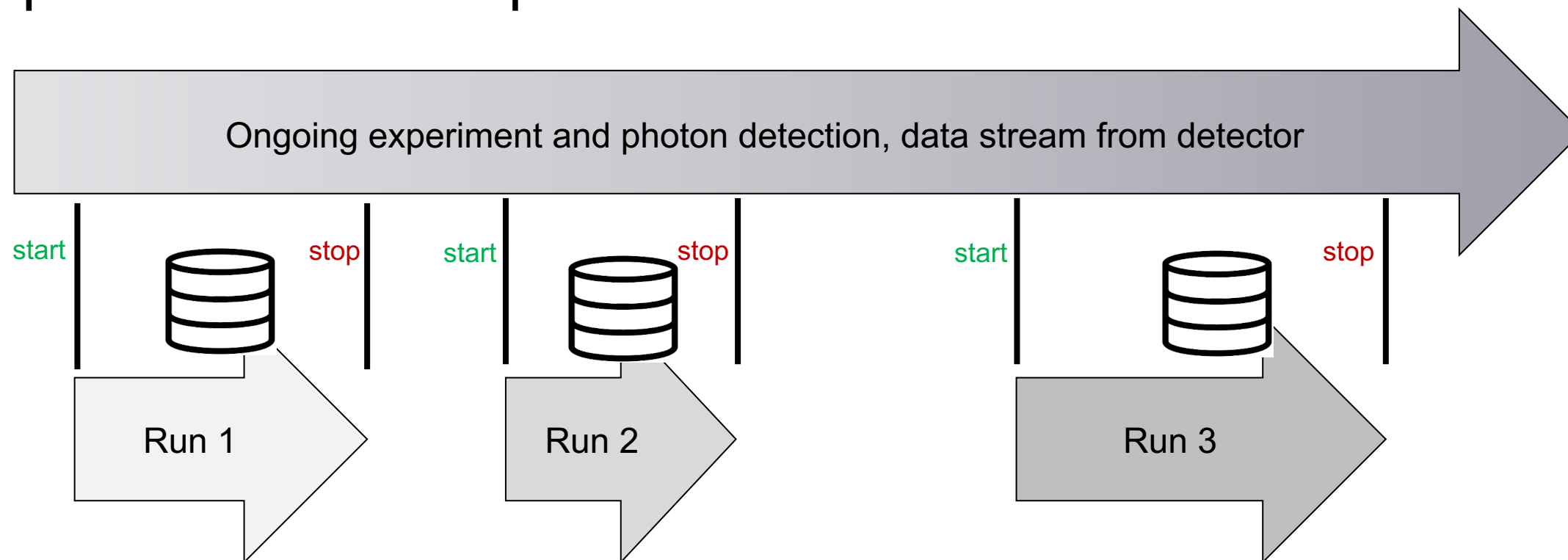- Large capacity

### DESY Data Center

**dCache**

- Capacity

- Lower performance
- Lower cost
- Scalability
- Shared within XFEL
- High capacity

**Tape Archive**

- Safety

- Very slow
- Even lower cost
- Very high capacity
- Safety (second copy)
- Shared within DESY campus
- Long term

**European XFEL**

# Experiment data acquisition to runs

# Volume of stored data



- accelerator can produce up to 27000 pulses per second (4.4 MHz max intra-train)

- detectors sync to pulse-train structure, typically 3500 to 8000 frames / second

- maximum of 52 TB/h (for the most frequent MHz detector), up to 4 PB raw data per beam-time, assuming highest efficiency



*Krzysztof Wrona & Janusz Malka*

- Experimental efficiency in terms of data taking time fraction has increased within the last year

**European XFEL**

# The EuXFEL scientific data policy in a nutshell

- Follows PaN Data Europe WG recommendations (2011)

- Users accept it to get beam-time awarded

- Raw data are embargoed for 3 years
  - then become automatically open (not stated what this exactly means)

- Access to results-data is restricted

- Derived data are not retained long-term
  - best-effort curation by EuXFEL, but not including e.g. SW for reproducibility

- Metadata: "good practice is encouraged" (no mention of concrete arrangements)

**European XFEL**

# Access of EuXFEL data via myMDC

# Open SFX data from EuXFEL

# The PaNOSC project

- EuXFEL was one of six contributing PaN facilities, 2018-2022

- Outcomes include:

  - FAIR recommendations for PaN data
  - Data analysis software and remote data analysis services
  - Simulation software
  - Solutions for unified authentication and other federated services
  - Training resources (learning platform, summer schools)

**GO FAIR**   FAIR Principles   Implementation Networks   News   Events   Resources   About GO FAIR

**What is the difference between "FAIR data" and "Open data" if there is one?**

- Open data and FAIR data are two related but different things

- Openness lies in the A of FAIR, access levels (authentication, authorization) need to be adjusted case by case

**European XFEL**

# Finding open EuXFEL data: the PaN search portal

Example data from real experiments (XMPL)



https://data.panosc.eu

▶ Metadata model of common items for search and filter terms

- Agreed-on subset, intersection of different catalogue systems
- PaNET ontology for experimental techniques

European XFEL

# The way forward

# Data reduction - the aim

Max. 6 months after
beam-time

Defined restricted volume of data
to be retained long-term on disk
(50 TB or 10% of collected raw data)

## Optional content of the RED data volume

▶ Selected runs of raw data

▶ All runs of raw data with event-selected detector frames

▶ Only processed data, likely selected and/or reduced, including corrected data in European XFEL format

▶ A mixture of selected raw data and processed data

RED data becomes OPEN data automatically after end of the embargo period, allowing for earlier OPEN data subsets upon publication

RED

OPEN    DOI 1    DOI 2    DOI 3

3 yrs.

**European XFEL**

# Data reduction - the methods

| Selective reduction | Transformative reduction | Numerical reduction |
|---|---|---|
| By experiment setup: runs | Temporal: averaging over trains, pulses | Rounding to less digits, lower bit-size, or integer |
| Temporal: trains, pulse pattern | Spatial: area/azimuthal integration | Lossless compression |
| By event: "hit" frames | | Lossy compression |
| Spatial: ROIs, modules | | |

- Most of these methods are performed in users' data analysis "anyway"
- In the context of the RED process, these methods will be implemented to facility procedures and deployed as largely automated services

**European XFEL**

# Data reduction – conceptual examples



Spatial selection (modules, ROIs)

Event selection

**European XFEL**

# Data lifecycle and data management plans



▶ DMPs formalize the collaborative planning among stakeholders before, during and after the experiment, given:

- Data usage requirements (storage, CPUs, software)
- Boundary conditions between stakeholders

European XFEL

Matthews, B. et al. (2012). Model of the data continuum in Photon and Neutron Facilities. PaN-data ODI Deliverable 6.1. https://doi.org/10.5281/zenodo.3897910

# Features of data management plans



▶ DMPs will be integrated to the existing facility workflow

▶ DMPs will be automatically pre-filled as much as possible with user office data and other proposal information

▶ DMPs are dynamic documents adapted at each stage from proposal acceptance to the open access phase

▶ DMPs will be editable by each of the partners, with automatic notification of changes

**European XFEL**

# Open data access: integration to myMDC



**10.22003/XFEL.EU-DATA-700000-00** ⧉

**Example Data**

The European XFEL (EuXFEL) example data proposal contains experimental data original beam-times, currently covering the techniques of serial coherent diffraction imaging (single particle imaging, SPI), X-ray powder scattering (SAXS) and X-ray photon correlation spectroscopy (XPCS).

**Released**

**Facility**

**Type**

**Services**

VISA

PaNdata Software Catalogue

**This proposal data is open**

**Would you like to get access to this proposal datasets?**

Please contact us through the open.data@xfel.eu email address

Thank you for visiting!

**European XFEL**

**Proposal Runs**

ℹ Automatically assess new runs (after being closed by DAQ) as: [To be evaluated manually ▾]

ℹ Automatically start run calibration after migration: [No ▾]

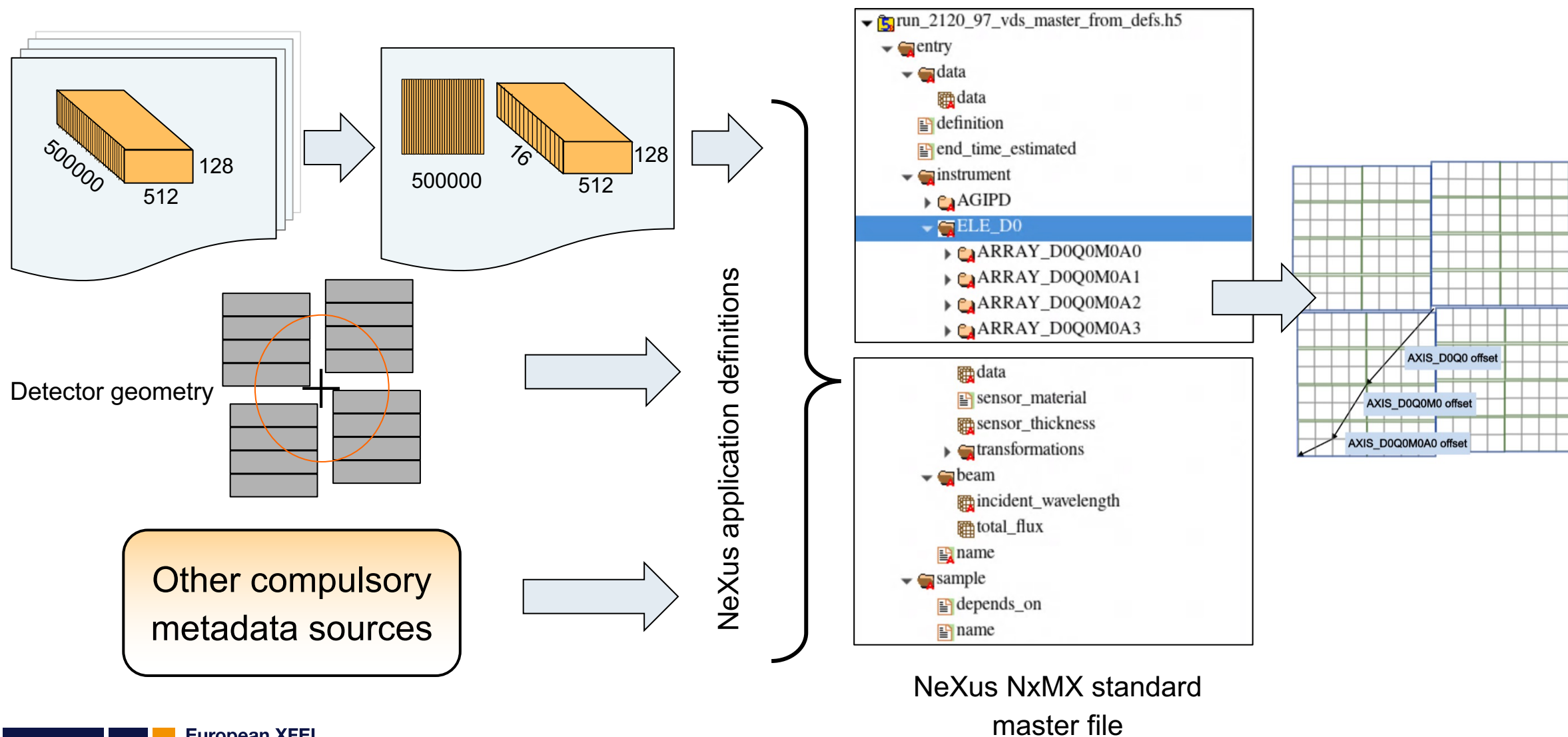| Run Number (alias) | Run type | Sample Name | Techniques | Start date | Run status | Data Assessment | Calibration | Run Comment | Edit |
|---|---|---|---|---|---|---|---|---|---|
| 0034 (SPI on sucrose solution, AGIPD detector at SPB instrument) | Single Particle Diffraction | Sucrose Solution 3% v/v | coherent diffraction imaging | 2021-06-01 02:25:08 +0200 | Closed | Good | ☑ ▾ | 👁 ⋮ ☑ | ☑ |
| 0033 (SAXS on vycor sample, AGIPD detector at MID instrument) | scattering | Vycor | small angle x-ray scattering | 2021-04-10 14:48:20 +0200 | Closed | Good | ☑ ▾ | 👁 ⋮ ☑ | ☑ |
| 0031 (SFX on Hen egg-white lysozyme, AGIPD detector) | Diffraction data | Lysozyme | serial femtosecond crystallography | 2021-04-15 10:48:26 +0200 | Closed | Good | ☑ ▾ | 👁 ⋮ ☑ | ☑ |
| 0030 (SFX on Hen egg-white lysozyme, AGIPD detector) | Diffraction | Lysozyme | serial femtosecond crystallography | 2020-03-09 01:20:02 +0100 | Closed | Good | ☑ ▾ | 👁 ⋮ ☑ | ☑ |
| 0029 (SFX on Hen egg-white lysozyme, AGIPD detector) | Diffraction | Lysozyme | serial femtosecond crystallography | 2020-03-09 01:07:51 +0100 | Closed | Good | ☑ ▾ | 👁 ⋮ ☑ | ☑ |
| 0027 (SAXS on 50 nm silica, AGIPD detector at MID instrument) | scattering | Silica 50nm | small angle x-ray scattering | 2019-09-21 01:12:49 +0200 | Closed | Good | ☑ ▾ | 👁 ⋮ ☑ | ☑ |
| 0026 (Time-resolved SAXS on Ni75-11 MLs, DSSC detector at SCS) | SAXS 500kHz // no pump laser | Ni75-11 MLs-b | small angle x-ray scattering | 2019-08-23 07:08:02 +0200 | Closed | Good | ☑ ▾ | 👁 ⋮ ☑ | ☑ |

Direct access to open data, requires enhanced Globus service

# Data interoperability: NeXus support



NeXus NxMX standard
master file

# Reusing open EuXFEL data: analysis in the cloud



Cloud cluster,
e.g. OpenStack

Web
application

VM

jupyter

https://visa.xfel.eu

- Enables data analysis from anywhere with a web browser, requires no local SW requirements or installations
- Data analysis is performed where the facility data is: no need to download tons of data
- Employs cluster resources with simple and transparent configuration of virtual machines
- Brings Jupyter notebooks and GUI programs with remote desktop under the same roof

**European XFEL**

# Summary and Outlook

- European XFEL data stems from a unique time structure with very high FEL pulse and detector repetition rates and accounts for extreme storage volumes

- Data reduction methods are known and need to be well implemented to automated facility services and anchored in an adequate scientific data policy.

- SFX data offers considerable reduction potential due to low hits rates in typical experimental setups

- Data analysis (incl. re-use) services for users exist – now they have to be expanded for guests interested in open, i. e. post-embargo, data

- To adopt FAIR principles better, developments for catalogue integration to search services, remote/cloud-based data analysis and standards for interoperability are ongoing

**European XFEL**

# Thanks and credits to all colleagues from European XFEL, collaborators of the PaNOSC/ExPaNDS projects, and users – for their work, help and ideas in fostering FAIR science adoption at our facility



## Data Analysis Group

Ivette Jazmin Bermudes Macias
Cammille Carinan
Matheus do Carmo Teodoro
Danilo Enoque Ferreira de Lima
Hadi Firoozi
Luca Gelisio
Thomas Kluyver
Amna Majid

Thomas Michelat
Robert Rosca
Philipp Schmidt
Egor Sobolev
Yue Sun
Oleksii Turkot
James Wrigley

## Data Department
## IT & Data Management group

Steve Aplin
Luis Maia
Janusz Malka
Krysztof Wrona



## PaNOSC & ExPaNDS



Andy Götz
Juncheng E
Michael Schuh
Tim Wetzel

**Thank you for your attention!**



Always FAIRplay!