# Standards-Compatible Data Storage in Laboratory-Based X-ray Instruments
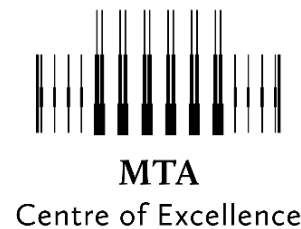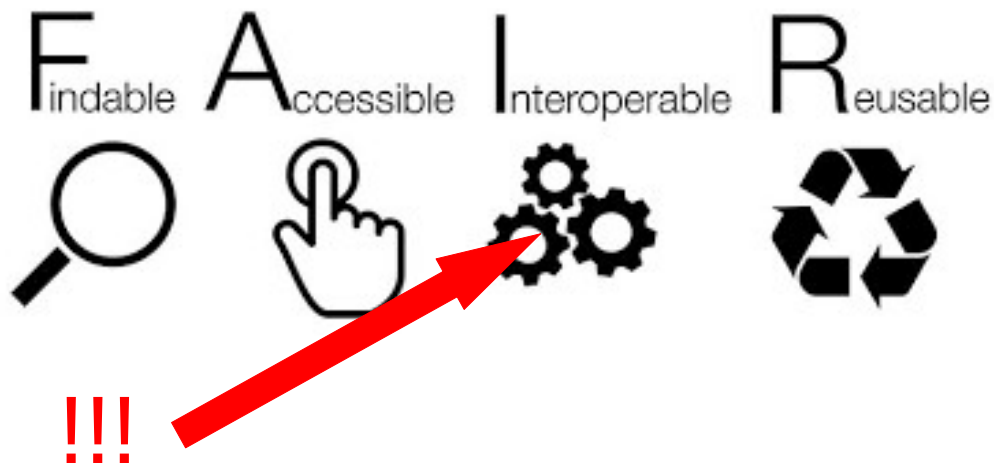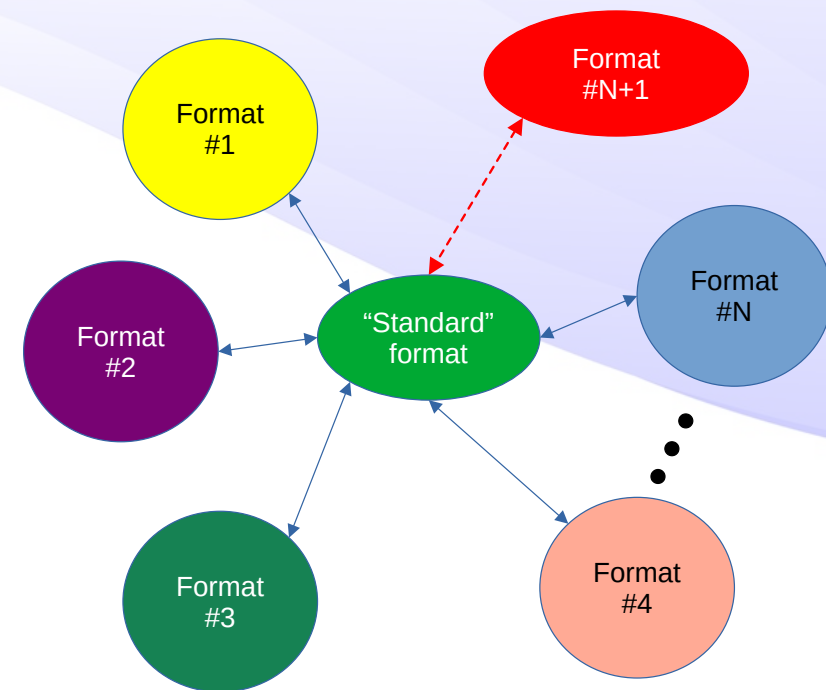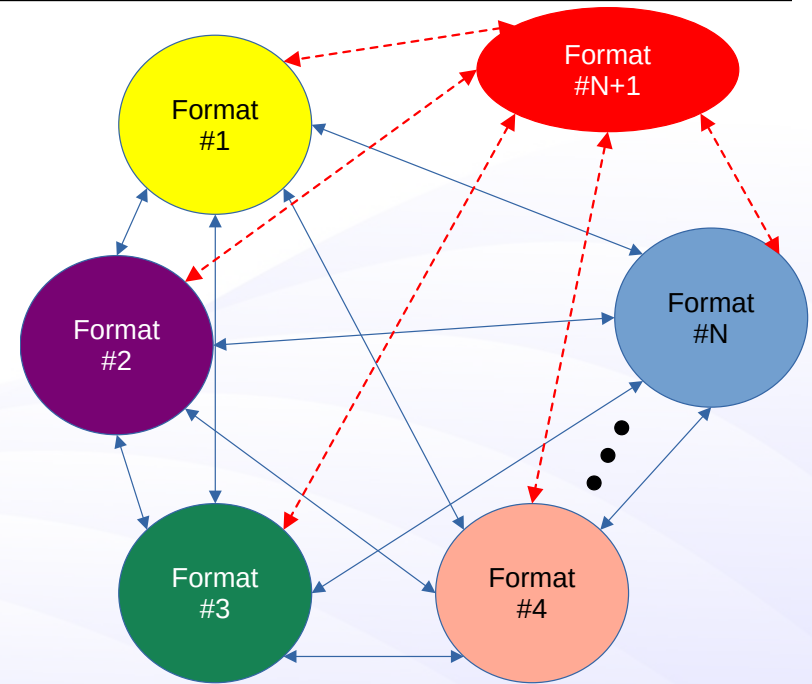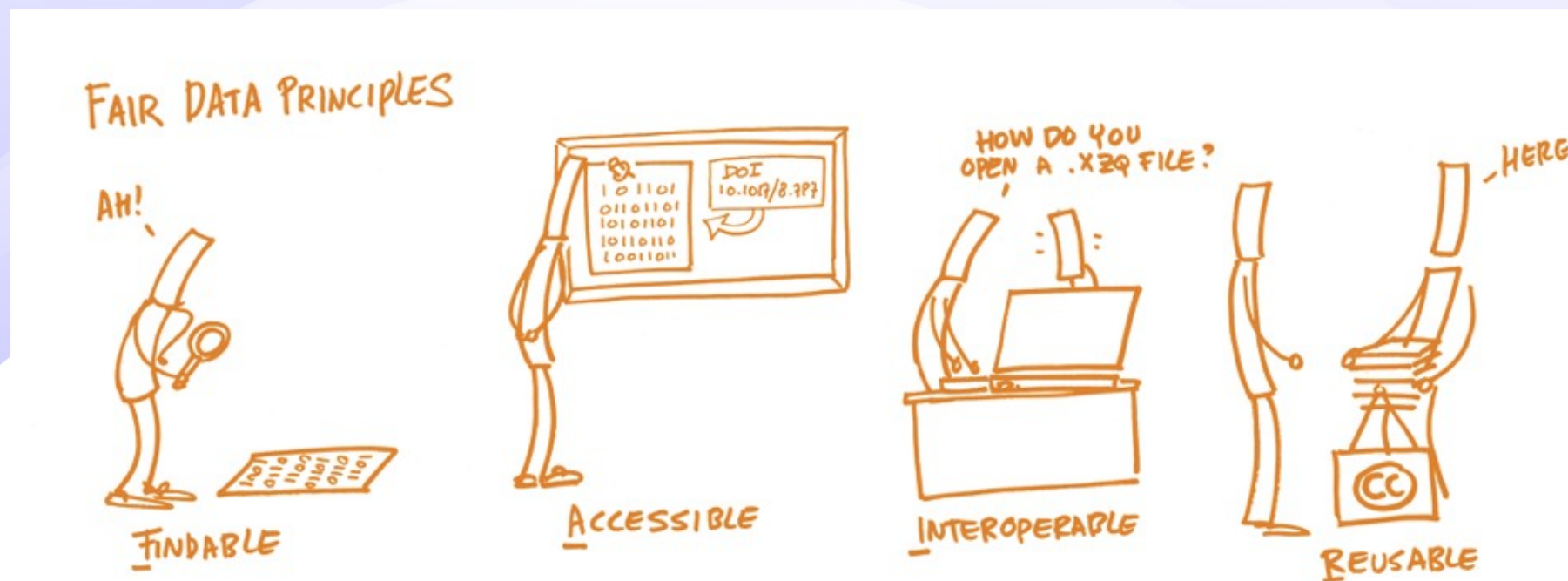
András Wacha

# Motivation: Standardized Data Storage Formats

- "It's good that we have so many standards to choose from"
- The evolution of storage formats together with hardware and software
- The question of interoperability
- Without standardization:
  - N formats, N×(N-1) converters
  - New format: 2N new converters needed
- With a standard format agreed upon:
  - N formats, 2N converters.
  - New format: 2 new converters needed
- FAIR principles of open data

Findable  Accessible  Interoperable  Reusable

!!!

# FAIR Principles in the Practice (https://www.openaire.eu/what-is-fair-data)

- **Findable:** "Discoverable with metadata, identifiable and locatable by means of a standard identification mechanism"
- **Accessible:** "Always available and obtainable; even if the data is restricted, the metadata is open"
- **Interoperable:** "Both syntactically parseable and semantically understandable, allowing data exchange and reuse between researchers, institutions, organizations or countries"
- **Reusable:** "Sufficiently described and shared with the least restrictive licences, allowing the widest reuse possible and the least cumbersome integration with other data sources"
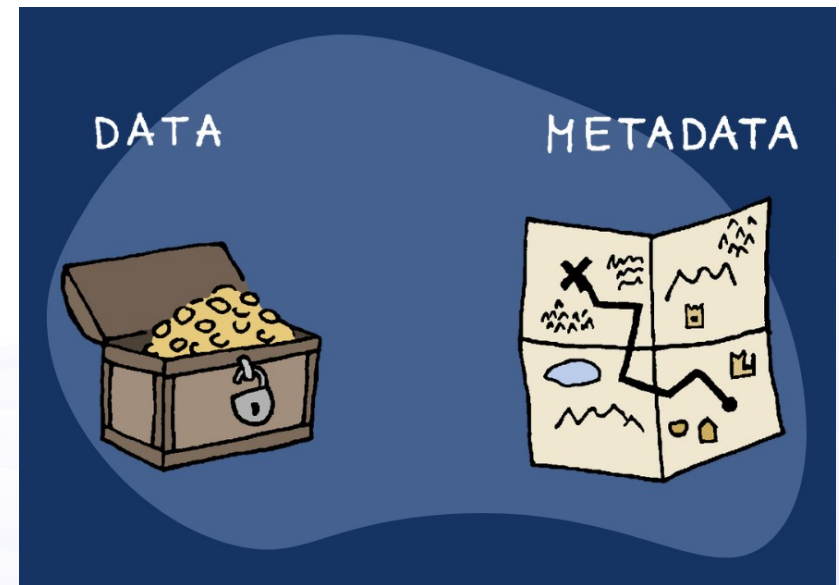
# Data and metadata, i.e., what to store

- Primary data: what we're most interested in
  - Typically non-scalar, possibly N-dim
  - Dependent variable as a function of the independent variable(s)
  - Fluorescence/absorption spectrum
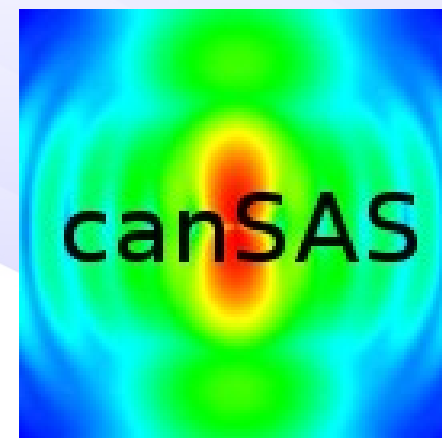  - Scattering pattern, scattering curve
  - Image
  - …
  - Raw vs. processed
- Metadata: describing the primary data
  - Information about:
    - The sample (name, composition, longer description…)
    - The instrument (name, state variables…)
    - Experimental conditions (temperature, *in situ* parameters)
    - Data reduction/evaluation procedures used
  - Ensuring correctness, reliability and reproducibility of the primary data

# FAIR Data in Photon (and Neutron) Science

- Photon and neutron open science cluster:
  PaNOSC (https://www.panosc.eu)
  - Part of the European Open Science Cloud (EOSC)
  - PaNOSC project: EC-financed, 2018-2022
  - Representing European photon and neutron research infrastructures
  - PaN-data Europe Deliverable D2.1: Common policy framework on scientific data
    - a generic data management policy
    - can be tailored by facilities to their own needs
  - Recommends the NeXus/HDF5 format for storing data and metadata
- CanSAS (https://www.cansas.org)
  - "collective action for nomadic small angle scatterers"
  - Providing the small-angle scattering user community with shared tools and information
  - First meeting: 1998
  - N-dimensional data: NeXus-based (NXcanSAS) 2017-06-06 (announced)
  - 1D data: XML-based format
    - cansas1d (v1.0: 2009-05-12, v1.1: 2013-03-29)
    - Recommendation from early 2017 on: store 1D also as NxcanSAS
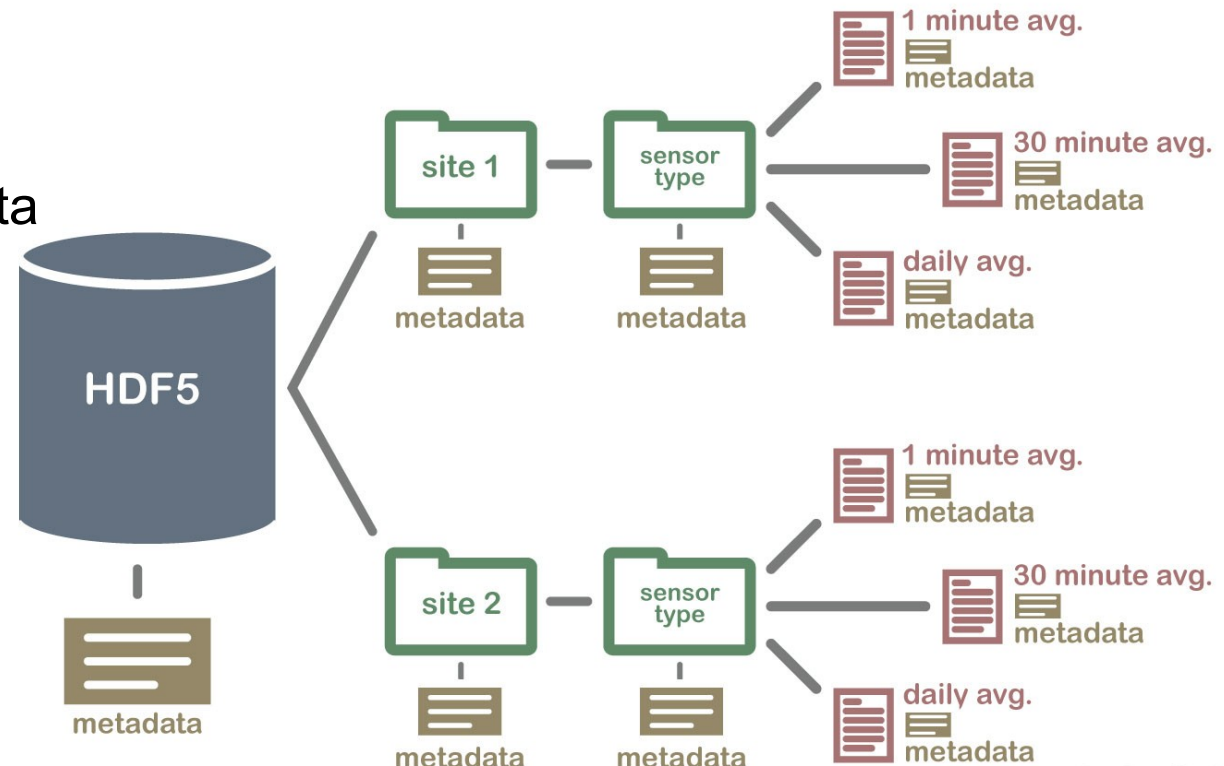
# How to choose a data format?

- Home-brewn format
  - Easy to implement if you are in a hurry, but usually not *future-proof*
  - Not well thought-out (not all relevant aspects are stored)
  - Difficult to extend (becomes "patchy", needs workarounds)
  - Difficult to maintain (protection against hardware and software obsolescence)
- Choose from an already existing standard
  - Takes efforts to implement
  - Be kind to your "nomadic users"
  - Good chances that many pieces of software already support it
- Other requirements
  - **Open:** specification, algorithms, libraries freely available
  - **Self-describing:** data are labeled, intuitively stored
  - **Compression:** store large datasets, preferably in a seamless way
  - **Fast read/write:** high throughput
  - **Fault tolerance / detection:** redundancy, checksums
  - **Straightforward API:** easy to access the data in many programming languages (scientists are not programmers)
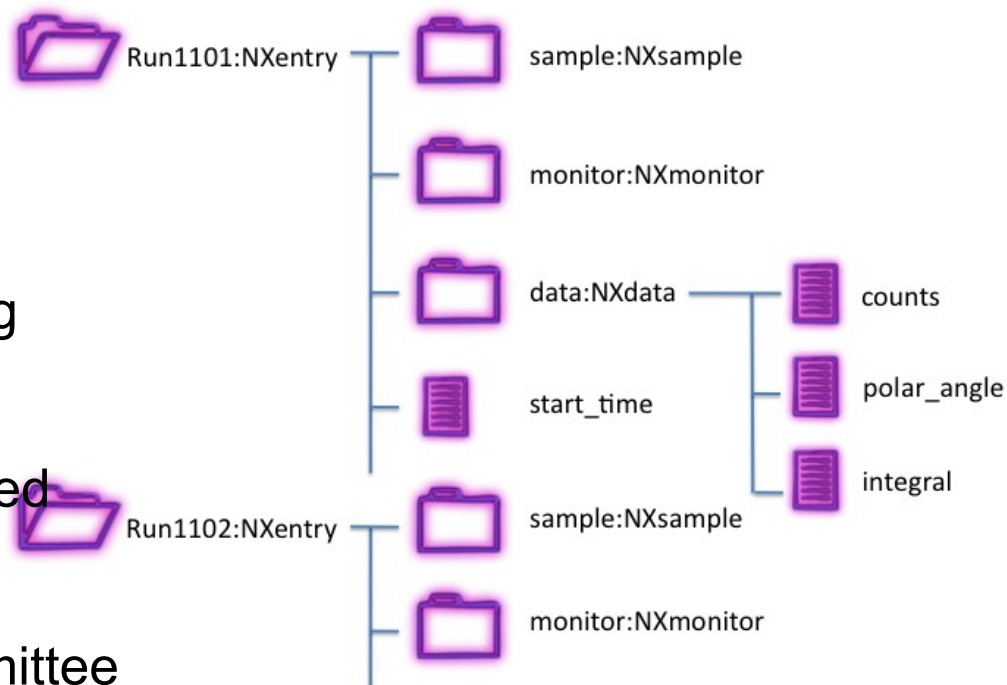
# Hierarchical Data Format       (https://www.hdfgroup.org/solutions/hdf5/ )

- Binary container format
- Cross-platform
- Many programming languages supported (C, C++, Java, Python…)
- Tree structure
  - Groups (~ folders)
    - containers
  - N-dimensional datasets (~ files)
    - N ≥ 0
    - Size, shape, data type
  - Symbolic or hard links
    - Pointers to the same data
  - External links
- Metadata (attributes) can be attached to groups and datasets
- Transparent IO filters
  - Compression
  - Shuffling
  - Fletcher32 checksum

# The NeXus Format

- *De facto* standard in photon, neutron and muon sciences
  - PaNOSC, canSAS…
- Common data exchange format
- Aims
  - Container for raw data, associated with a scientific instrument
  - Container for processed data
- Based on the HDF5 format
  - HDF5 is only a container
  - standardized structure
  - Standardized nomenclature
- Aims:
  - Domain-specific rules for organizing and arranging data
  - Quick default visualization
  - Standard definitions that can be used to validate files
- Governing body: NIAC
  - Nexus International Advisory Committee
  - Convenes every other year
- https://www.nexusformat.org
- J Appl Crystallogr (2015) 48(1) 301-305 (doi:10.1107/S1600576714027575)

# NeXus Base Concepts I: NeXus base classes

- The "type" of a HDF5 group, stored in the "NX_class" attribute of the group
- Prescription of the possible fields (datasets) and attributes
  - Field types are also specified
- Corresponds to real-world objects (sample, instrument, data reduction step…)
- https://manual.nexusformat.org/classes/base_classes/index.html
- Example 1: NXentry (describe the measurement)
  - @NX_class="NXentry" (attribute)
  - title (NX_CHAR): title of the entry
  - start_time (NX_DATE_TIME): starting time of the measurement
  - end_time (NX_DATE_TIME): ending time of the measurement
  - program_name (NX_CHAR): the program used to generate this file
  - Sample (NXsample): a group describing the sample
  - Instrument (NXinstrument): a group describing the instrument
  - ...
- Example 2: NXsample (describe the sample)
  - @NX_class="NXsample" (attribute)
  - name (NX_CHAR): name of the sample
  - chemical_formula (NX_CHAR): chemical formula
  - temperature (NX_FLOAT): temperature of the sample
  - ...

## NeXus Base Concepts II: Application definitions

- Base classes only define the nomenclature, not requirements
- Application definitions: domain-specific rules on obligatory and optional data
- Declared in the "definition" field of the NXentry
- Example: NXxas (X-ray absorption spectroscopy measurements)
- https://manual.nexusformat.org/classes/applications/NXxas.html
- Lower case names: required name
- Upper case names: arbitrary name
- Path specification in the file:
  - By name:
    - entry145/vonhamos/xraygenerator/name
  - Using NXclass attributes:
    - Nxentry/Nxinstrument/Nxsource/name
- Multiple methods (e.g. raw SAXS and processed SAXS): NXsubentry

---

**ENTRY**: (required) NXentry

  **@entry**: (required) NX_CHAR

    ▼ NeXus convention is to use "entry1", "entry2", … …

    NeXus convention is to use "entry1", "entry2", … for analysis software t

  **title**: (required) NX_CHAR ⏎

  **start_time**: (required) NX_DATE_TIME ⏎

  **definition**: (required) NX_CHAR ⏎

    ▶ Official NeXus NXDL schema to which this file conforms …

**INSTRUMENT**: (required) NXinstrument ⏎

  **SOURCE**: (required) NXsource ⏎

    **type**: (required) NX_CHAR ⏎

    **name**: (required) NX_CHAR ⏎

    **probe**: (required) NX_CHAR ⏎

      Obligatory value: x-ray

  **monochromator**: (required) NXmonochromator ⏎

    **energy**: (required) NX_FLOAT (Rank: 1, Dimensions: [nP]) ⏎

  **incoming_beam**: (required) NXdetector ⏎

    **data**: (required) NX_NUMBER (Rank: 1, Dimensions: [nP]) ⏎

  **absorbed_beam**: (required) NXdetector ⏎

    **data**: (required) NX_NUMBER (Rank: 1, Dimensions: [nP]) ⏎

    This data corresponds to the sample signal.

**SAMPLE**: (required) NXsample ⏎

  **name**: (required) NX_CHAR ⏎

  Descriptive name of sample

**MONITOR**: (required) NXmonitor ⏎

  **mode**: (required) NX_CHAR ⏎

    ▶ Count to a preset value based on either clock time (timer) …

# NeXus Base Concepts III: Contributed definitions

- Tentative, suggested extensions to the NeXus specification
- Proposed by the community
- Not yet standardized
- Both base classes and application definitions
- Example fields:
  - Optical spectroscopy
  - Multi-dimensional photoemission spectroscopy
  - Atom probe microscopy
  - Electron microscopy
  - …
- Curated, commented on and finally incorporated into the NeXus standard by the NIAC

# NeXus Base Concepts IV: Default Visualization

- Each measurement should have a default visualization...
- … which should be declared in the data file
- The NXdata class
  - "Encapsulates all the information required for a set of data to be plotted"
  - *Signals*: dependent variables (1- or more dimensions)
    - Default signal: "signal" attribute
  - *Axes*: independent variables (typically 1D, but can be more)
  - Names freely chosen (but *cf* application definitions)
  - https://manual.nexusformat.org/classes/base_classes/NXdata.html
  - Example:
    -

```
data: NXdata
   @signal = "data"
   @axes = ["x", "y"]
   data: float[10, 20]
   x: float[10]
   y: float[20]
```

```
data: NXdata
   @signal = "data"
   @axes = ["x", "y"]
   @x_indices = 0
   @y_indices = 1
   data: float[10, 20]
   x: float[10]
   y: float[20]
```

# More complex NXdata

- Multi-dimensional data, e.g., scans

```
data:NXdata
  @signal = "data"
  @auxiliary_signals = ["data2", "data3"]
  @axes = ["x", "y", "energy", "wavelength"]
  @x_indices = 0
  @y_indices = 1
  @energy_indices = 2
  @wavelength_indices = 2
  data: float[10, 20, 30]
  x: float[10]
  y: float[20]
  energy: float[30]
  wavelength: float[30]
  data2: float[10, 20, 30]
  data3: float[10, 20, 30]
```

```
data:NXdata
  @signal = "absorption"
  @axes = ["xpos", "ypos", "energy"]
  @xpos_indices = [0, 1, 2]
  @ypos_indices = [0,1,2]
  @energy_indices = 2
  absorption: float[10, 20, 30]
  xpos: float[10, 20, 30]
  ypos: float[10, 20, 30]
  energy: float[30]
  absorption_errors: float[10, 20, 30]
  xpos_errors: float[10, 20, 30]
  ypos_errors: float[10, 20, 30]
```

# NeXus utilities – Silx view

- Displays HDF5 tree and the default (customizable) plot
- https://www.silx.org

# Other NeXus utilities: programming & whatnot

- Reading and writing in Python
  - h5py (https://www.h5py.org): HDF5 wrappers for Python
  - PyTables (https://www.pytables.org): alternate HDF5 wrappers, sposnored by NUMFOCUS (Numpy, Scipy & Co.)
  - NeXpy (https://nexpy.github.io/nexpy/):  high-level Python interface (+ GUI)
- Validation + other utilities
  - Punx (https://github.com/prjemian/punx)
  - NeXus command-line utilities: nxbrowse, nxconvert, nxdir…
- Data analysis programs supporting NeXus files:
  - DAVE (https://www.ncnr.nist.gov/dave/): for inelastic neutron scattering
  - DAWN (https://www.dawnsci.org): generic visualization, domain-specific processing
  - Mantid (http://mantidproject.org): high-performance computing on neutron and muon data
  - PyMCA (https://pymca.sourceforge.net): X-ray fluorescence data analysis
  - … (see https://manual.nexusformat.org/utilities.html)

# How We Store NeXus Files in the CREDO System?

# SAXS data

- Modern SAXS measurements: 2D position sensitive detectors
- Isotropic (unoriented, powder, …) sample: dependence only on $q=|\mathbf{q}|$.
  - Azimuthal average → radial scattering curve ("intensity vs. $q$")
- Anisotropic scattering patterns:
  - Azimuthal averages in sectors
  - Radial average in an annulus: azimuthal scattering curve ("intensity vs. $\varphi$")
- Typical detector format: ~ 4-9 Mpixel
  - File size ~ 30-50 Mbyte
- Data reduction
  - External and internal background
  - Geometrical corrections
  - Normalization by exposure time and beam intensity
  - Correct for X-ray absorption by the sample



Scattered rays

Scatterer ("sample")

Incident rays

Unscattered + forward scattered radiation

Beam stop

Azimuthal averaging

Scattering pattern

Scattering curve

Intensity (1/cm × 1/sr)

q (1/nm)

# Storing NeXus files

- Application definitions:
  - Raw data: NXsas
  - Processed data: home-brewn (moving to NXcansas)
- https://nexdatas.github.io
  - Developed at DESY
  - A set of Tango servers:
    - NeXus file writer
    - Configuration server for NeXus file writer
    - Component selector
  - GUI
    - Component designer
  - Sardana/Taurus extensions
    - Macro GUI
    - Sardana Recorder which uses the Tango servers
- In-house developed Sardana macros
  - nxsbegin – ct – nxsend: most metadata written while the exposure is being made

**Environment cycle #1**

Set up

DAQ loop

Backgrounds & references

Sample #1

Sample #2

...

**Environment cycle #2**

Set up

DAQ loop

...

# Data storage strategy

- **File sequence: <prefix>_<file_sequence_index>.<extn>**
  - Examples: crd_52133.cbf (raw detector image), crd_52133.nxs (NeXus file, separate NXsubentry for raw and processed)
  - Prefix: independently counted images, based on the use (we use the same 2D detector for everything!):
    - scn: scan measurements
    - tra: transmission measurements
    - crd: "production", true measurements
    - tst: test shots
    - gsx: GISAXS
  - A SQL database is also written (raw; updated with processed)
    - Basis for **F**indability later on
- **End result: a single dataset for each sample**
  - ? time-evolution ?
  - Transitioning to NXcansas in the near future

# Data post-processing GUI

- Load a range of a file sequence

# Data post-processing GUI

- Analyze metadata

# Data post-processing GUI

- Average exposures corresponding to the same sample and same sample-to-detector distance
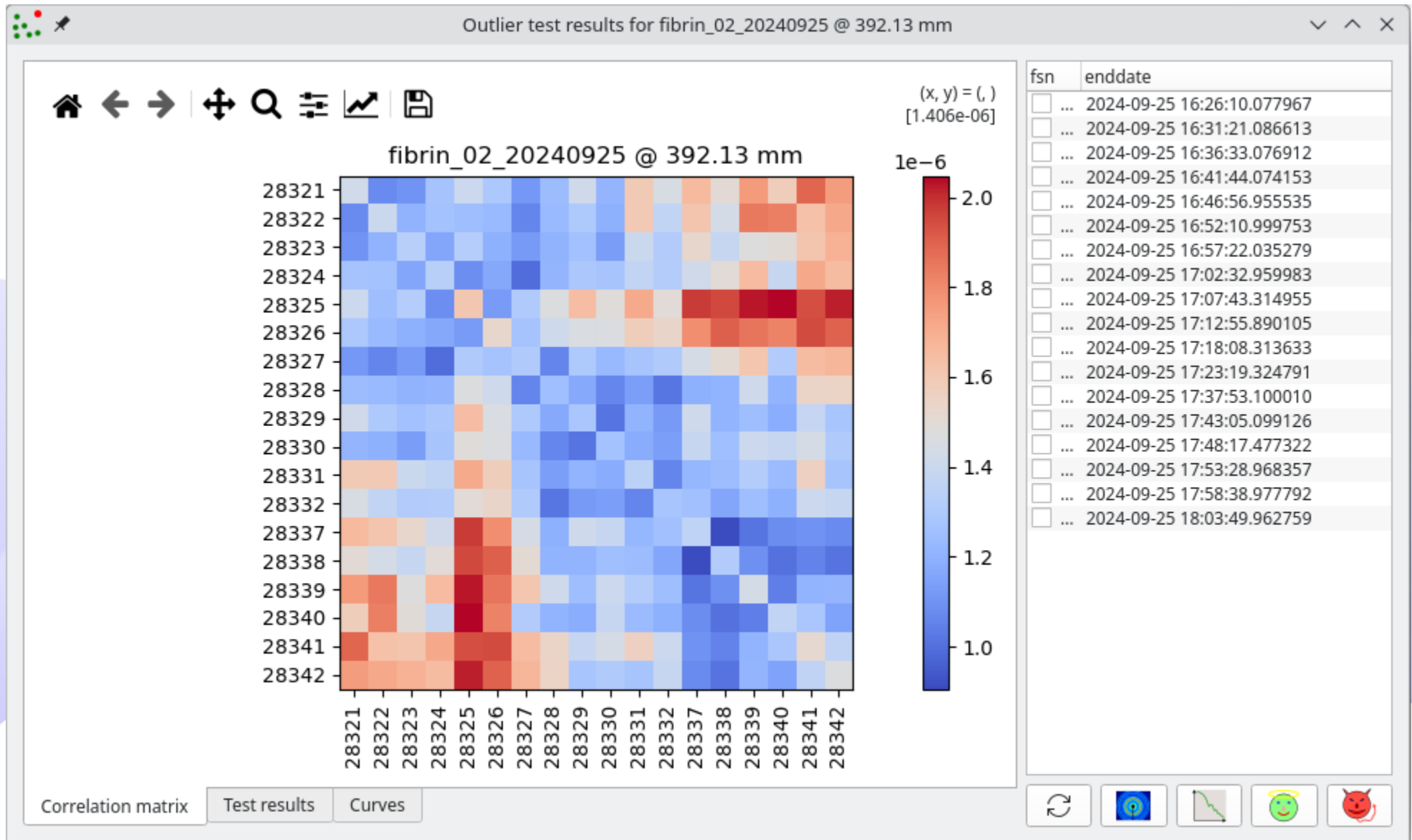
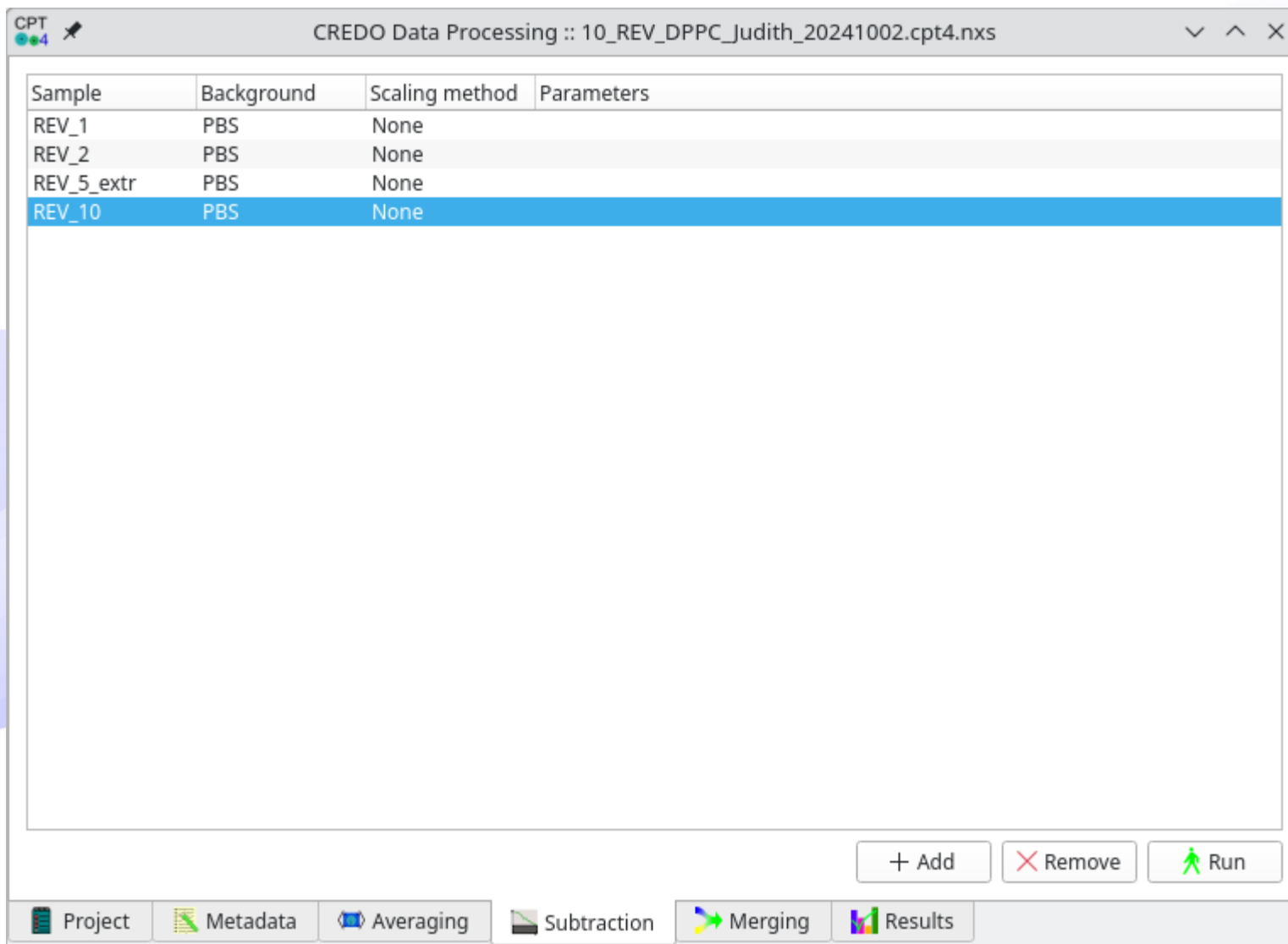# Data post-processing GUI

- Filter exposures with artefacts (Pilatus chip flashes), assess sample stability

# Data post-processing GUI

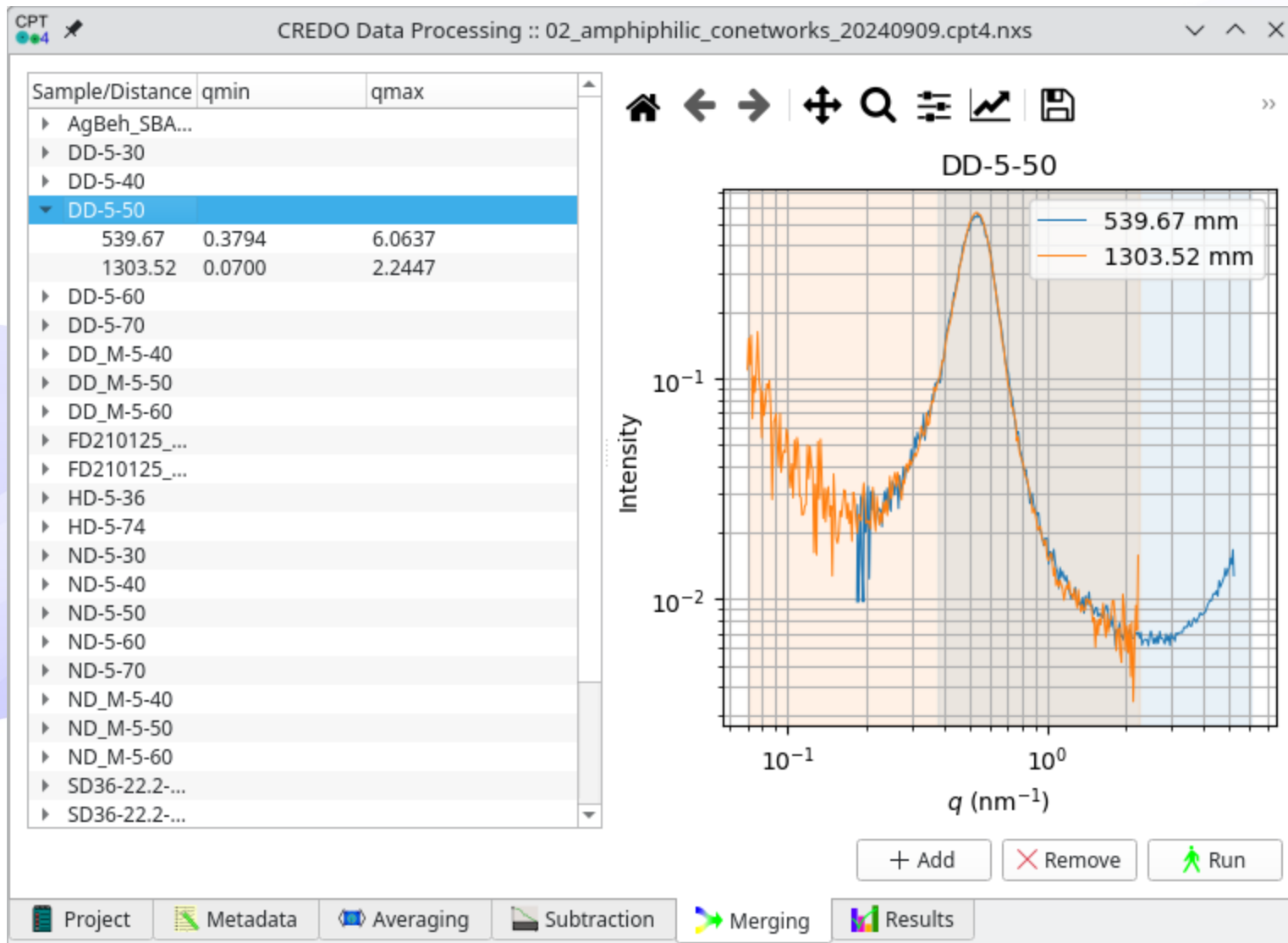- Filter exposures with artefacts (Pilatus chip flashes), assess sample stability

# Data post-processing GUI

- Subtract background (optional)

# Data post-processing GUI

- Merge curves from multiple sample-to-detector distances

# Data post-processing GUI

- Present results: draw curves, patterns, analyze anisotropy, export to various formats

# Conclusions

- **Storage format for in-house**
  - Home-brewn might be okay on the short term
  - Standardized format on the long term
- User facility → requirement of interoperability (and F+A+R, too)
- Completeness vs. simplicity
- Find the most agreed-upon data format for your domain
  - SAXS: NeXus, canSAS
  - XAS: NeXus?
  - HDF5 should be a good general choice
- Exporting processed data to other formats
  - "No matter how good a data storage format you adapt, most users will ask for ASCII text files"
  - "… or Excel workbooks."

# Thank you for your attention (again)!

- Research Group for Biological Nanochemistry, HUN-REN Research Centre for Natural Sciences (https://bionano.ttk.hu/biological-nanochemistry)
- CREDO SAXS laboratory (https://credo.ttk.hu)