



# Gold Standard Format Proposal

Herbert. J. Bernstein

Ronin Institute for Independent Scholarship

Notes for HDRMX Meeting at Diamond Light Source  
6-7 Nov 2019



rev 9 Dec 2019



# Changes in metadata

- Since we introduced the CBF format for the Dectris Pilatus detectors in 2007, there has been a recognition of the importance of controlling the metadata associated with images in order both to ensure that essential information is not lost and to minimize delays in handling the metadata
- When the Eiger detectors were introduced, the community agreed to adopt the NeXus/HDF5 format for efficiency in handling the much larger volume of data with fewer files to reduce filesystem and network burdens, but much of the metadata is carefully aligned between NeXus/HDF5 and CBF under an agreement between the NeXus International Advisory Committee (NIAC) and the IUCr Committee for the Maintenance of the CIF Standard (COMCIFS).
- With the cooperation of Dectris, the High Data Rate Macromolecular Crystallography (HDRMX) group and website were established to facilitate community discussion of the software, data and metadata.

# HDRMX Discussions on Metadata

- There are signs of divergence among beamlines in Eiger formats and it is time to add new metadata, for example to identify beamlines and facilities and to record metadata that will be helpful in PDB depositions.
- The primary objective is to ensure that sufficient metadata will be provided to allow processing at a facility other than the one at which the data was produced. In particular, detailed descriptions of axis chains to be used to process the data are needed, both for sample goniometers and detector positioners.
- There were two informal HDRMX dinner meetings in summer 2019 to discuss a new “gold standard” for the NeXus HDF5 Eiger format, one at the July 2019 ACA meeting in Covington, KY and one at the August 2019 ECM32 in Vienna, AT
- We hope to come to final agreement on the changes at this HDRMX meeting at Diamond Light Source, 6-7 Nov 2019

# Structure of the New Metadata

- In general the requested augmentation of metadata is divided into two groups
- First, metadata to be added via a templating mechanism in the Dectris software to be set-up before collection as static changes to the "master" files, and,
- Second, metadata to be added after collection, possibly via H5copy. For simplicity we refer to the former as static and the latter as dynamic.

## Static Metadata:

- Some tags for static (*i.e.* Dectris template) additions are already available. imgCIF defines AXIS tags needed for specification of arbitrary and very general axis chains. NeXus defines the equivalent information in the NXtransformations base class.
- Concern has been expressed about cluttering the templating mechanism with large numbers of tags used only in the most complex cases.

# Static Metadata

- To avoid such clutter the input to the template can be the path to either a CBF or a NeXus file with the appropriate axis information, along with the necessary software to automatically convert between CBF and NeXus axis conventions. One way or another all diffraction geometry and all detector geometry need to be described.
- Tags have been defined to carry metadata specifying the beamline and facility.
- Note that the detector distance, wavelength and beam center are already specified and very necessary.
- As integrating detectors or other detectors that do not count single photons come into use in this performance range, detector gain will need to be specified.
- Tags are needed for the HDF5 software version, to declare the use of non-standard local format conventions, to list the files comprising a dataset, and to give the format of each particular file.

# Static Metadata Example

- As a partial example consider a beamline called XXX (ID1) at site SYNC with an omega axis, and pin\_x, pin\_y and pin\_z translation axes stacked 5 millimetres apart, using hdf5\_1.8.14 and NXmx 1.4. Then a portion of the necessary information presented as a CIF file might be:

```
data_AMX_metadata
```

```
  loop_
```

```
    _axis.id _axis.type _axis.equipment _axis.depends_on
```

```
    _axis.vector[1] _axis.vector[2] _axis.vector[3]
```

```
    _axis.offset[1] _axis.offset[2] _axis.offset[3]
```

```
source . source . 0 0 1 ...
```

```
gravity . gravity . 0 -1 0 ...
```

```
pin_x translation goniometer . -1 0 0 0 0 0
```

```
omega rotation goniometer pin_x 1 0 0 -5 0 0
```

```
pin_y rotation goniometer omega 0 1 0 -10 0 0
```

```
pin_z rotation goniometer pin_y 0 0 -1 -15 0 0
```

# Static Metadata Example (cont.)

```
_array_intensities.gain      1.0 #counts/photon

_diffrn_source.source        SYNCHROTRON
_diffrn_source.type          'SYNC XXX (ID1)'
_diffrn_source.pdbx_synchrotron SYNC
_diffrn_source.pdbx_synchrotron_beamline 'XXX (ID1)'

_dataset_file_format.file_format 'hdf5_1.8.14 and NXmx 1.4'

_diffrn_radiation.beam_width 7 #micrometres
_diffrn_radiation.beam_height 5 #micrometres
_diffrn_radiation.beam_flux   400000000000 #ph/s in the beam
```

# Dynamic Metadata

Many tags for dynamic (non-Dectris-template) additions are already available. For example, the monochromator, the beam\_height, beam\_width, beam\_flux and sample sequence can all be placed by a beamline or user in a CIF or NeXus file for merging with H5copy into an existing master metadata file. **The existing imgcif and mmcif dictionaries provide possibilities**, and more can be added. The following have been discussed:

sample provenance, sample physical characteristics, sample imagery, protein sequence, detector and sample environments, incl. temperature, sample delivery method, serial crystallography parameters (incl. pump probes), spectroscopy, sample mount, detector ROI. beamline optics, source parameters, e.g. mode, current, collection strategy, scan type, scan mode, beam profile (Gaussian, tophat), monochromator bandpass, beam divergences, beam collimation.



# NXmx Gold Standard Proposal

These are the proposed changes from the NeXus 2016 NXmx 1.4. Under the Gold Standard items may be optional, recommended or required.

(entry): NXentry (required)

title: (optional) NX\_CHAR

start\_time: (recommended) NX\_DATE\_TIME

**was optional**

end\_time: (recommended) NX\_DATE\_TIME

**was optional**

definition: (required) NX\_CHAR

NeXus NXDL schema to which this file conforms

Obligatory value: NXmx

(data): (recommended) NXdata

**was required**

data[np, i, j, k]: (recommended) NX\_NUMBER

**was required**

(sample): (required) NXsample

**was under /entry/instrument**

name: (recommended) NX\_CHAR

**was optional**

depends\_on: (recommended) NX\_CHAR

**was optional**

The axis on which the sample position depends may be stored anywhere, but is normally stored in the transformations:NXtransformations group within the NXsample group.

temperature: (optional) NX\_CHAR {units=NX\_TEMPERATURE}

(transformations): (recommended) Nxtransformations

**was optional**

# NXmx Gold Standard Proposal

(instrument): (required) Nxinstrument

name: (required) NX\_CHAR

Name of instrument, i.e. beamline name

@short\_name: (required) NX\_CHAR

(attenuator): (optional) Nxattenuator

attenuator\_transmission: (optional) NX\_NUMBER {units=NX\_UNITLESS}

(detector\_group): (recommended) NXdetector\_group **was optional**

Optional logical grouping of detector elements.

group\_index[i]: (required) NX\_INT

An array of unique indices for detector elements or groupings of detector elements.

group\_names: (required) NX\_CHAR

group\_parent[group\_index]: (required) NX\_INT

An array of the hierarchical levels of the parents of detector elements or

A top-level element or grouping has parent level -1

# NXmx Gold Standard Proposal

(detector): (required)

Normally the detector group will have the name detector. However, in the case of multiple detector elements, each element needs a uniquely named NXdetector group.

depends_on: (recommended) NX_CHAR	<b>was required</b>
data[np, i, j, k]: (recommended) NX_NUMBER	<b>was required</b>
description: (recommended) NX_CHAR	
time_per_channel: (optional) NX_CHAR {units=NX_TIME}	
distance: (recommended) NX_FLOAT {units=NX_LENGTH}	<b>was optional</b>
dead_time: (optional) NX_FLOAT {units=NX_TIME}	
count_time: (recommended) NX_NUMBER {units=NX_TIME}	<b>was optional</b>
distance: (recommended) NX_FLOAT {units=NX_LENGTH}	<b>was optional</b>
beam_center_x: (recommended) NX_FLOAT {units=NX_LENGTH}	<b>was optional</b>
beam_center_y: (recommended) NX_FLOAT {units=NX_LENGTH}	<b>was optional</b>
angular_calibration_applied: (optional) NX_BOOLEAN	
angular_calibration[i, j, k]: (optional) NX_FLOAT	
flatfield_applied: (optional) NX_BOOLEAN	
flatfield[i, j, k]: (optional) NX_FLOAT	
flatfield_error[i, j, k]: (optional) NX_FLOAT	

# NXmx Gold Standard Proposal

pixel\_mask\_applied: (optional) NX\_BOOLEAN  
pixel\_mask[i, j]: (recommended) NX\_INT **was optional**  
countrate\_correction\_applied: (optional) NX\_BOOLEAN  
bit\_depth\_readout: (recommended) NX\_INT **was optional**  
detector\_readout\_time: (optional) NX\_FLOAT {units=NX\_TIME}  
frame\_time: (optional) NX\_FLOAT {units=NX\_TIME}  
gain\_setting: (optional) NX\_CHAR  
saturation\_value: (optional) NX\_INT  
sensor\_material: (recommended) NX\_CHAR **was optional**  
sensor\_thickness: (optional) NX\_FLOAT {units=NX\_LENGTH}  
threshold\_energy: (optional) NX\_FLOAT {units=NX\_ENERGY}  
type: (optional) NX\_CHAR  
**(transformations): (recommended) Nxtransformations **was optional****  
**(collection): (optional) Nxcollection**  
**(detector\_module): (required) NXdetector\_module**  
    data\_origin: (required) NX\_INT  
    data\_size: (required) NX\_INT  
    data\_stride: (optional) NX\_INT

# NXmx Gold Standard Proposal

module\_offset: (optional) NX\_NUMBER {units=NX\_LENGTH}  
    @transformation\_type: (required) NX\_CHAR  
    Obligatory value: translation  
    @vector: (required) NX\_CHAR  
    @offset: (required) NX\_CHAR  
    @depends\_on: (required) NX\_CHAR

fast\_pixel\_direction: (required) NX\_NUMBER {units=NX\_LENGTH}  
    @transformation\_type: (required) NX\_CHAR  
    Obligatory value: translation  
    @vector: (required) NX\_CHAR  
    @offset: (required) NX\_CHAR  
    @depends\_on: (required) NX\_CHAR

slow\_pixel\_direction: (required) NX\_NUMBER {units=NX\_LENGTH}  
    @transformation\_type: (required) NX\_CHAR  
    Obligatory value: translation  
    @vector: (required) NX\_CHAR  
    @offset: (required) NX\_CHAR  
    @depends\_on: (required) NX\_CHAR

# NXmx Gold Standard Proposal

(beam): (required) NXbeam **was under sample, now under instrument**

incident\_wavelength: (required) NX\_FLOAT {units=NX\_WAVELENGTH}

In the case of a monochromatic beam this is the scalar wavelength.

In the case of a polychromatic beam this is an array of the wavelengths.

incident\_wavelength\_weight: (optional) NX\_FLOAT

In the case of a polychromatic beam this is an array of the relative weights of the corresponding wavelengths in incident\_wavelength.

incident\_wavelength\_spread: (optional) NX\_FLOAT {units=NX\_WAVELENGTH}

The wavelength spread FWHM for the corresponding wavelength(s) in incident\_wavelength.

flux: (optional) NX\_FLOAT {units=NX\_FLUX}

flux incident on beam plane area in photons per second per unit area

total\_flux: (required) NX\_FLOAT {units=NX\_FREQUENCY}

**was optional**

flux incident on beam plane in photons per second

incident\_beam\_size[2]: (recommended) NX\_FLOAT {units=NX\_LENGTH}

profile: (recommended) NX\_CHAR

Gaussian | Airy | top-hat | rectangular

incident\_polarisation\_stokes[np, 4]: (recommended) NX\_CHAR

incident\_wavelength\_spectrum: (optional) NXdata

# NXmx Gold Standard Proposal

(source): (required) NXsource

The neutron or x-ray storage ring/facility.

distance: (optional) NX\_FLOAT {units=NX\_LENGTH}

Effective distance from sample

name: (required) NX\_CHAR

@short\_name: (optional) NX\_CHAR

**was required**

type: (optional) NX\_CHAR

**was required**

type of radiation source (pick one from the enumerated list and spell exactly)

probe: (optional) NX\_CHAR

**was required**

type of radiation probe (pick one from the enumerated list and spell exactly)

power: (optional) NX\_FLOAT {units=NX\_POWER}

**was required**

Source power

emittance\_x: (optional) NX\_FLOAT {units=NX\_EMITTANCE}

**was required**

Source emittance (nm-rad) in X (horizontal) direction.

emittance\_y: (optional) NX\_FLOAT {units=NX\_EMITTANCE}

**was required**

Source emittance (nm-rad) in Y (horizontal) direction.

sigma\_x: (optional) NX\_FLOAT {units=NX\_LENGTH}

**was required**

particle beam size in x

# NXmx Gold Standard Proposal

sigma_y: (optional) NX_FLOAT {units=NX_LENGTH} particle beam size in y	<b>was required</b>
flux: (optional) NX_FLOAT {units=NX_FLUX} Source intensity/area (example: s-1 cm-2)	<b>was required</b>
energy: (optional) NX_FLOAT {units=NX_ENERGY} Source energy. For storage rings, this would be the particle beam energy. For X-ray tubes, this would be the excitation voltage.	<b>was required</b>
current: (optional) NX_FLOAT {units=NX_CURRENT}	<b>was required</b>
voltage: (optional) NX_FLOAT {units=NX_VOLTAGE}	<b>was required</b>
frequency: (optional) NX_FLOAT {units=NX_FREQUENCY}	<b>was required</b>
period: (optional) NX_FLOAT {units=NX_PERIOD}	<b>was required</b>
target_material: (optional) NX_CHAR Pulsed source target material	<b>was required</b>
number_of_bunches: (optional) NX_INT	<b>was required</b>
bunch_length: (optional) NX_FLOAT {units=NX_TIME}	<b>was required</b>
bunch_distance: (optional) NX_FLOAT {units=NX_TIME}	<b>was required</b>
pulse_width: (optional) NX_FLOAT {units=NX_TIME}	<b>was required</b>
mode: (optional) NX_CHAR source operating mode	<b>was required</b>



# NXmx Gold Standard Proposal

top_up: (optional) NX_BOOLEAN	<b>was required</b>
last_fill: (optional) NX_NUMBER {units=NX_CURRENT}	<b>was required</b>
@time: (required) NX_DATE_TIME	<b>was required</b>
date and time of the most recent injection.	
notes: (optional) NXnote	
<b>was required</b>	
bunch_pattern: (optional) NXdata	<b>was required</b>
title: (required) NX_CHAR	<b>was required</b>
name of the bunch pattern	
pulse_shape: (optional) NXdata	<b>was required</b>
source pulse shape	
geometry: (optional) NXgeometry	<b>was required</b>
“Engineering” location of source	
distribution: (optional) NXdata	<b>was required</b>
The wavelength or energy distribution of the source	

# Validate Images

- Especially with new metadata being added, a fast data-driven tool for NeXus/HDF5 image validation is needed.
- The best available tool is `cnxvalidate` by Mark Koennecke,
  - <https://github.com/nexusformat/cnxvalidate> (current)
  - <https://github.com/HDRMX/cnxvalidate> (proposed)
- which is data driven working against
  - <https://github.com/nexusformat/definitions> (current)
  - <https://github.com/HDRMX/definitions> (proposed)

Typical call and output are

```
nxvalidate -a NXmx -l ~/definitions -e thau2_25dps_tr0p05_1_master.h5  
message="Missing required global file_name attribute"  
... sev=error dataPath=/ dataFile=thau2_25dps_tr0p05_1_master.h5
```

# Validation Plans

- After the DLS HDRMX meeting in November 2019, the agreed changes will be integrated with the development version of cnxvalidate and submitted to Dectris and NIAC for review.
- If all goes well users should start seeing validated gold standard images in use in early 2020.

# Useful Links

<http://hdrmx.medsbio.org>

<http://github.com/HDRMX>

<http://github.com/HDRMX/cnxvalidate>

<http://github.com/HDRMX/definitions>

<http://hdrmx.medsbio.org/manual/build/html>

<https://zenodo.org/record/3385862> -- Small example Eiger 2X 16M data set from Diamond Light Source I04, Graeme Winter

<https://zenodo.org/record/3484187> -- Small example Eiger 2X 16M data set from Diamond Light Source I04, Graeme Winter, revised for clean cnxvalidate error report, Graeme Winter, Aaron Brewster, Herbert J. Bernstein

<https://zenodo.org/record/3352358> -- 68 image lysozyme dataset recorded on the Jungfrau 16M detector at SwissFEL and formatted as a NeXus file, Aaron Brewster, Meitian Wang

<https://zenodo.org/record/3526738> -- 68 image lysozyme dataset recorded on the Jungfrau 16M detector at SwissFEL and formatted as a NeXus file, revised for clean cnxvalidate error report, Aaron Brewster, Meitian Wang, Herbert J. Bernstein

# Acknowledgments

- Work done in part at NSLS-II, Brookhaven National Laboratory
- Supported in part by Dectris, Ltd.
- US Department of Energy Offices of Biological and Environmental Research and of Basic Energy Sciences (grant No. DE-AC02-98CH10886; grant No. E-SC0012704);
- U.S. National Institutes of Health (grant No. P41RR012408; grant No. P41GM103473; grant No. P41GM111244; grant No. R01GM117126; grant No. P30GM133893, grant No. R21GM129570)
- Our thanks to Frances C. Bernstein for helpful consultations