# Integration of CBF, NeXus and HDF5

HERBERT J BERNSTEIN,[a]* GRAEME WINTER,[b] TOBIAS S. RICHTER,[b]

MARK BASHAM,[b] JAMES PARKHURST,[b] NEXUS INTERNATIONAL ADVISORY

COMMITTEE[c] AND COMMITTEE ON THE MAINTENANCE OF THE CIF STANDARD[d]

[a]*Department of Mathematics and Computer Science, Dowling College, Oakdale, NY 11769 USA*, [b]*http://wiki.nexusformat.org/NIAC*, and

[c]*http://www.iucr.org/resources/cif/comcifs*. *E-mail: yayahjb@gmail.com*

## Abstract

This is a report on progress in the integration of CBF, NeXus and HDF5 in support of the high data rate/high data volume (BIG DATA) demands of the new generation of X-ray pixel array detectors. Fast Dectris Pilatus detectors and XFEL detectors such as the Cornell-SLAC pixel array detector (CSPAD) are already straining file systems, and the new generation of even faster detectors, such as the Dectris Eiger, are bringing this issue to a head. In addition, the modular nature of these detectors provides the opportunity to construct more complex detector arrays (e.g. the Dectris Pilatus detector at I23 at DLS), which in turn requires a more complete description of the detector geometry. Taken together these give rise to a need to combine the best of CBF/imgCIF (the Crystallographic Binary File, which has a complete description of the experiment), NeXus (a common data framework for neutron, X-ray and muon science, which gracefully handles large data sets) and HDF5 (Hierarchical Data Format,

version 5, the high-performance data format used by NeXus) for the management of such data at synchrotrons. A proof-of-concept API based on CBFlib and the HDF5 API is now in use and will continue to be developed. Among the issues considered are: The roles of CBF, NeXus and HDF5, the NXmx (NeXus macromolecular crystallography) application definition, the specific needs of the DECTRIS EIGER and the CSPAD XFEL detector, recent progress on multiple module, multiple data array support, and the handling of structure factors.

## 1. Introduction

This is a report on progress in the integration of CBF (Bernstein & Hammersley, 2005), NeXus (Filges, 2001) and HDF5 (Filges, 2001) in support of the high data rate/high data volume (BIG DATA) demands of the new generation of X-ray pixel array detectors. Since 2010 (Bernstein, 2010) (Bernstein, 2012) (Bernstein *et al.*, 2013) (Bernstein, 2013) (Bernstein *et al.*, 2014*a*) (Bernstein *et al.*, 2014*b*), we have been exploring the feasibility of combining the capabilities of CBF, NeXus and HDF5. Fast Dectris Pilatus detectors and XFEL detectors such as the Cornell-SLAC pixel array detector (CSPAD) are already straining file systems, and the new generation of even faster detectors, such as the Dectris Eiger are bringing this issue to a head.

CCD X-ray detectors provide images at moderate data rates of one every few to several seconds (see *e.g.* `http://www.adsc-xray.com/Q4techspecs.html`). Current higher performance X-ray detectors, such as the DECTRIS Pilatus, are capable of collecting six-megapixel images at 10 to 25 frames per second (Trueb *et al.*, 2012), while the newest Pilatus3 6M detectors can operate at 100 frames per second (see `https://www.dectris.com/pilatus3\_specifications.html`). The coming next generation of high performance X-ray detectors for MX such as the DECTRIS Eiger will be capable of collecting 16+ megapixel images at more than 125 frames per second

(Willmott, 2011) (Johnson *et al.*, 2012). See Table 1.

In addition, the modular nature of these detectors provides the opportunity to construct more complex detector arrays (e.g. the Dectris Pilatus detector at I23 at DLS, see Fig. 1), which in turn requires a more complete description of the detector geometry. Taken together these give rise to a need to combine the best of CBF/imgCIF (the Crystallographic Binary File, which has a complete description of the experiment), NeXus (a common data framework for neutron, X-ray and muon science, which gracefully handles large data sets) and HDF5 (Hierarchical Data Format, version 5, the high-performance data format used by NeXus) for the management of such data at synchrotrons. Resolving all the issues involved will take a collaborative effort extending over many years. This paper reviews the nature of the problem and the progress made to date. A proof-of-concept API based on CBFlib and the HDF5 API that is being developed in a collaboration among Dowling College, Brookhaven National Laboratory and Diamond Light Source is now in use and will continue to be developed. Releases of CBFlib since CBFlib 0.9.2.12 can store arbitrary CBF files in HDF5 and recover them, can support use of all CBFlib compressions in HDF5 files, and can convert sets of miniCBF files to a single NeXus file. The latest release, CBFlib 0.9.5, is operational for HDF5 handling of single detector module, monochromatic MX data compatibly with imgCIF. Among the issues considered here are: The roles of CBF, NeXus and HDF5, the NXmx (NeXus mcoromolecular crystallography) application definition, the specific needs of the DECTRIS EIGER and the CSPAD XFEL detector, recent progress on multiple module, multiple data array support, and the handling of structure factors.

## 2. HDF5 and NeXus

The Hierarchical Data Format Version 5 (HDF5) is a self-describing file format with a robust, well documented API routinely handling multi-gigabyte files of data. It has a diverse user community covering a wide range of disciplines and is fully supported (Dougherty *et al.*, 2009). HDF5 is particularly well suited to the management of very large volumes of complex scientific data, has been adopted as the primary data format in a wide range of disciplines (`http://www.hdfgroup.org/HDF5/users5.html`) and provides the inner workings of important frameworks, such as NetCDF (Rew *et al.*, 2004) and NeXus. To avoid confusion we use the term format to describe the logical organization of data on a storage medium. An ontology is a dictionary of terms that may include descriptions of the relationships between terms. An ontology can be realized in one or more formats. We are therefore dealing with the HDF5 format, the NeXus ontology, a CBF format, and an imgCIF ontology. The HDF5 format, XML format and NeXus ontology together form the NeXus data transfer framework. The CBF format, CIF format and imgCIF ontology form the imgCIF data transfer framework. HDF5 is tree-oriented, which is a very powerful and useful characteristic allowing file-system-like nesting of groups of data within groups of data, in order for information to be easily, reliably and efficiently searched. However, tables are more useful for loading information into a relational database management system (Codd 1970). NeXus (Filges, 2001)(Koennecke, 2006) is a tree-oriented ontology for use wth HDF5 (and XML and HDF4) of importance in managing neutron and X-ray data. NeXus adds rules for storing data in files and a dictionary of documented names to HDF-5 in order to make HDF-5 applicable to the problem domain of synchrotron, neutron and muon scattering. NeXus is a convenient thin layer over HDF5 that is widely used at many physics research centers, including at synchrotrons. Together NeXus and HDF5 provide a portable, extensible and efficient framework for the storage and

management of data.

## 3. Why a change is needed

Today for MX alone Diamond Light Source employs three Pilatus 6M fast and two Pilatus 3 6M, giving a combined data rate of over 1 GB/sec and over 200 files/sec. These new detectors are creating the need to manage hundreds of thousands of images being received at rates from sixty megapixels to 2.5 gigapixels per second and beyond. For the Advanced Beamlines for Biological Investigations with X-rays (ABBIX) that are being built for NSLS-II (Hendrickson, 2012), just two of the beam lines, the Frontier Macromolecular Crystallography (FMX) beamline and the Automated Macromolecular Crystallography (AMX) beamline (Schneider *et al.*, 2012), are expected to produce an aggregate of approximately 10 terabytes per operational half day, 660 terabytes per week or 38 petabytes per year. The anticipated beamline flux is $10^{13}$ photons per second for FMX and $2 \times 10^{13}$ photons per second for AMX, approximately 50 times the NSLS X25 and X29 fluxes. One subtle effect of these high fluxes is that there will be more photons per pixel in images, making them more difficult to compress.

The normal practice in designing crystallographic diffractometers and beam lines has been to store one image frame per file. As data rates rise, the number of files being opened per unit time rises. File systems that were designed to accept hundreds of files being opened each second begin to choke when thousands of files are being opened in the same timeframe. The rates implied for a single beam line using an EIGER X 16M are shown in Fig 2. Because many beam lines may share the same common file system, and *in situ* processing may require several files opens for each image in rapid succession, so having several beamlines, all using such detectors, mandates strongly for grouping some large number of images per file to reduce the file-opening burden

on the file system. While a CBF can hold a large number of images per file, the use of HDF5, a practice that is common for high energy physics, gives a facility more control over the combined burden.

### 3.1. Jan 2013 DECTRIS Eiger Workshop and Followup

The attendees at the January 2013 DECTRIS Workshop agreed on the use of an HDF5-based NeXus framework for the DECTRIS Eiger pixel array detector. The workshop charged Herbert J. Bernstein with following up on mapping additional terms to the new format. Tobias Richter, Jonathan Sloan and Herbert J. Bernstein worked on a CBF-NeXus concordance and supporting software based on CBFlib and HDF5 with the cooperation of Bob Sweet, Graeme Winter and Mark Koennecke. Discussions with NIAC were held and then discussions with COMCIFS were held prior to ECM 28 in August 2013. There was general agreement that it was a good idea to have CIF and NeXus interoperate. COMCIFS and NIAC have agreed to start on a single crystal monochromatic macromolecular crystallography experiment NeXus application definition. An application definition in NeXus is a specification of the required metadata and data for that application. The NXmx application definition is in use and has proven workable. It is now being expanded to handle cases in which multiple data arrays are required, either because, as with the EIGER, the data volumes are too large to allow a single data array to be reliably stored under HDF5, or because, as with the Pilatus 12M-DLS or the CSPAD, data arrays with different index axes are required. An application definition for reflection data is also being tested.

## 4. Reflection Data API

*** This section is under development ***

## 5. Compression

The discussion of file opening, above, argues for use of an HDF5-based format with multiple images stored in single arrays. This, however, does not address the high data volumes themselves. To conserve disk space, network bandwidth and energy, compression of images as early as possible is needed. Table 2 show the range of current choices of compressions when applied to typical MX data. High speed, high compression ratio compression is a critical issue for the next generation of detectors. Some compressions raise license issues. Some popular compressions are slow or inefficient or both. Some compressions can be handled in processing programs such as XDS if license and language issues can be addressed. Low pixel density fine-slicing with clean backgrounds makes some compressions more effective. CBFlib provides useful compressions. A plugin has been written to allow HDF5 to read and write CBFlib compressions.

For the DECTRIS Pilatus 300K image shown in Fig 3, the compressions are shown in Table 3. All the HDF5 presentations of the data see a modest increase in size due to the overhead of the more complex format. For larger images this would not be as significant a percentage. This particular data, having a noisy background and significant spots, does not compress well. For many experiments using fast detectors, it is now feasible to take very large numbers of fine-sliced images that have very few photons per image, resulting in images that consist primarily of pixels containing zero with a small number of pixels with very few counts. Fortunately, such images often can be faithfully compressed by factors of 10 to 60.

## 6. Mapping from NeXus to CBF

All NeXus base classes now have proposed slots in CIF categories. Handling of the DECTRIS Eiger HDF5 format is under development with a preliminary concordance having been specified. For this project, organizing data and metadata according to the

conventions of the IUCr Crystallographic Information File (Hall *et al.*, 1991) using imgCIF (Bernstein & Hammersley, 2005) and its open source supporting software CBFlib (Ellis & Bernstein, 2005) provides a database-friendly tabular structure. The imgCIF ontology provides the metadata needed for the analysis of diffraction images and is supported by all the major detector manufacturers. This aspect is particularly important for instruments with complex geometries, e.g., the Pilatus 12M-DLS for the long wavelength beamline I23 at Diamond Light Source. The necessary metadata for use with CBF for this detector has been has been generated. In order to map that metadata to NeXus, the canSAS approach to associating axes with array indices will be used.

The embedding of CIF tables in HDF5 files was demonstrated at the "HDF5 as hyperspectral data analysis format" workshop in January 2010 (Götz *et al.*, 2010). The workshop recommendation was, in part, to "adopt as much as possible from imgCIF and sasCIF". Tables are easily embedded into trees but going in the other direction is more difficult. There is serious effort required to make general trees into tables suitable for use in a relational database management system, involving a process known as "normalization" (Codd, 1972). One of the tasks of this project is to extend the imgCIF ontology to ensure workable database access to metadata in the HDF5 tree that has not already been normalized into CIF categories. For example, Digital Object Identifiers (DOIs) and SHA2 or SHA3 checksums from multiple experiments will need to be brought forward into a common table for post-experiment forensic validation.

***To be continued after a clean presentation of the comparison mappings within the limitations of journal pages if worked out. At present it appears much of it will need to be in supplementary material.***

# Appendix A
# NXmx Application Definition

*** The NXMx details will either appear as an appendix or in supplementary data***

Acknowledgements

We gratefully acknowledge the work by Jonathan Sloan, formerly at Diamond Light Source, contributing to the code in CBFlib for mapping between CBF and NeXus.

Our thanks to James Hester, Chair of COMCIFS and Mark Koennecke, Chair of NIAC, for supporting and encouraging this effort and to all participants in COMCIFS and NIAC for helpful comments. Our thanks to Nick Sauter and Aaron Brewster for extending this work to CSPAD. Our thanks for years of supporting efforts at the BNL PXRR Group: Robert M. Sweet, Dieter Schneider, Howard Robinson, John Skinner, Matt Cowan, Leonid Flaks, Richard Buono; at DLS: Alun Ashton, Bill Pulford; and at the Dowling College ARCiB Lab Group: Mojgan Asadi, Kostandina Bardhi, Keti Bardhi Frey Lewis, Limone Rosa, Our thanks to DECTRIS, BIOIHDF and the HDF Group. Our thanks to Frances C. Bernstein. This work was funded in part by NIGMS, DOE, NSF, PaNdata ODI (EU 7th Framework Programme) and DECTRIS.

## References

Bernstein, H. J. (2010). In *HDF5 as hyperspectral data analysis format Workshop,11 – 13 January 2010, ESRF, Grenoble, FR*.

Bernstein, H. J. (2012). In *DIALS-East, Partnering Data Collection and Reduction in the Beamline Environment workshop, July 27, 2012, Harvard Medical School, Boston, MA*.

Bernstein, H. J. (2013). In *DECTRIS Eiger Workshop, Baden, Switzerland, 24-25 January 2013*.

Bernstein, H. J. & Hammersley, A. P. (2005). In *International Tables For Crystallography*, edited by S. R. Hall & B. McMahon, vol. G: Definition and Exchange of Crystallographic Data, chap. 2.3, pp. 37 – 43. International Union of Crystallography, Springer, Dordrecht, NL.

Bernstein, H. J., Sloan, J. M., Winter, G., Richter, T. S., NIAC & COMCIFS (2013). In *American Crystallographic Association, 2013 Annual Meeting, Honolulu, HI, July 19 – 24, 2013*.

Bernstein, H. J., Sloan, J. M., Winter, G., Richter, T. S., NIAC & COMCIFS (2014*a*). *Computational Crystallography Newsletter*, **5**, 12 – 18.

Bernstein, H. J., Sloan, J. M., Winter, G., Richter, T. S., NIAC & COMCIFS (2014*b*). *Computational Crystallography Newsletter*, **5**, 12 – 18. `http://www.phenix-online.org/newsletter/CCN_2014_01.pdf{\#}page=12`.

Codd, E. F. (1972). *Courant Computer Science Symposium 6*, chap. Further Normalization of the Data Base Relational Model, pp. 33 – 64. Prentice-Hall, 1972.

Dougherty, M. T., Folk, M. J., Bernstein, H. J., Bernstein, F. C., Eliceiri, K. W., Benger, W., Zadok, E. & Best, C. (2009). *Communications of the ACM*, **52**(10), 42 – 47.

Ellis, P. J. & Bernstein, H. J. (2005). *Definition and Exchange of Crystallographic Data, International Tables For Crystallography*, chap. CBFlib: an ANSI C library for manipulating image data, pp. 544 – 556. International Union of Crystallography, Springer, Dordrecht, NL.

Filges, U. (2001). In *VITESS Workshop Berlin, 25 – 27 June 2001*.

Götz, A., Solé, V., Madonna, C. & Maydew, A. F., (2010). Elisa vedac workshop report, workshop title: Hdf5 as hyperspectral data exchange and analysis format, grenoble, january 11th to january 13th, 2010. `http://vedac.esrf.eu/public-discussions/hdf5-workshop/workshop-report`.

Hall, S. R., Allen, F. H. & Brown, I. D. (1991). *Acta Crystallographica Section A: Foundations of Crystallography*, **47**(6), 655 – 685.

Hendrickson, W. A. (2012). *NSLS-II - Status of the Life Sciences Program*. Tech. rep. Brookhaven National Laboratory, X6A Science Advisory Committee. `protein.nsls.bnl.gov/mediawiki/images/e/e3/Hendrickson_2012.pdf`.

Johnson, I., Bergamaschi, A., Buitenhuis, J., Dinapoli, R., Greiffenberg, D., Henrich, B., Ikonen, T., Meier, G., Menzel, A., Mozzanica, A. *et al.* (2012). *Journal of Synchrotron Radiation*, **19**(6), 0–0.

Koennecke, M. (2006). *Physica B: Condensed Matter*, **Vol.385386**. Proceedings of the Eighth International Conference on Neutron Scattering.

Rew, R., Ucar, B. & Hartnett, E. (2004). In *20th Int. Conf. on Interactive Information and Processing Systems*.

Schneider, D. K., Sweet, B. M. & Skinner, J. (2012). Projection of MX ABBIX needs at AMX and FMX for data acquisition, data processing, software, data archiving and networking. Private Communication.

Trueb, P., Sobott, B., Schnyder, R., Loeliger, T., Schneebeli, M., Kobas, M., Rassool, R., Peake, D. & Broennimann, C. (2012). *Journal of Synchrotron Radiation*, **19**(3), 347 – 351.

Willmott, P. (2011). *An Introduction to Synchrotron Radiation: Techniques and Applications*. Chichester, UK: John Wiley and Sons. Page 6.

Table 1. *Examples of detector speeds.*

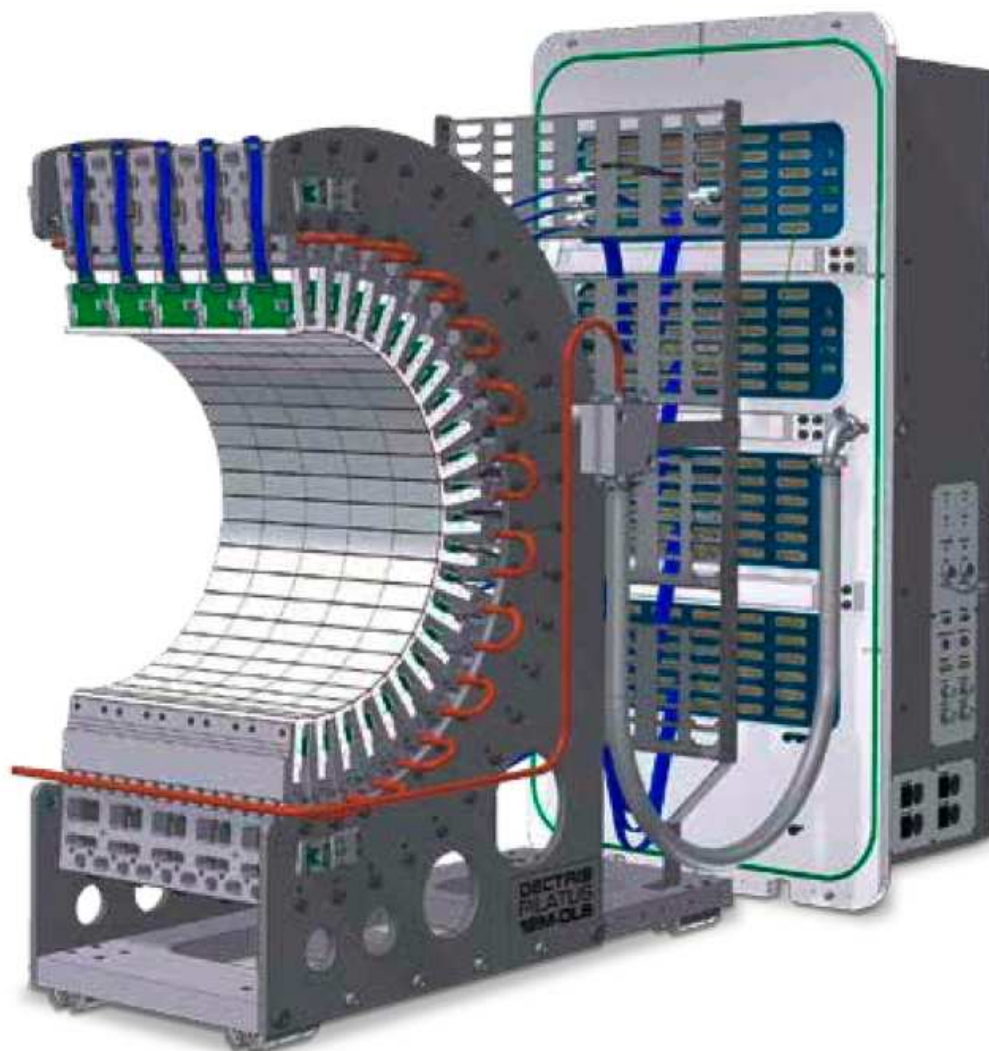| Detector | Raw Image size (MB) | Frame Rate (Hz) | Compressed Rate (GB/s / Gb/s) |
|---|---|---|---|
| ADSC Q315 (2x2 binned) | 18 | 0.37 | .0016 / .013 |
| Pilatus 2 6M | 24 | 10 | .06 / .48 |
| Pilatus 2 Fast 6M | 24 | 25 | .15 / 1.2 |
| Pilatus 3 6M | 24 | 100 | .6 / 4.8 |
| Eiger X 16M | 72 | 133 | 2.4 / 19 |

Fig. 1. The DECTRIS I23-DLS detector. This is an approximately cylindrical detector requiring a different set of module axes to describe each module.
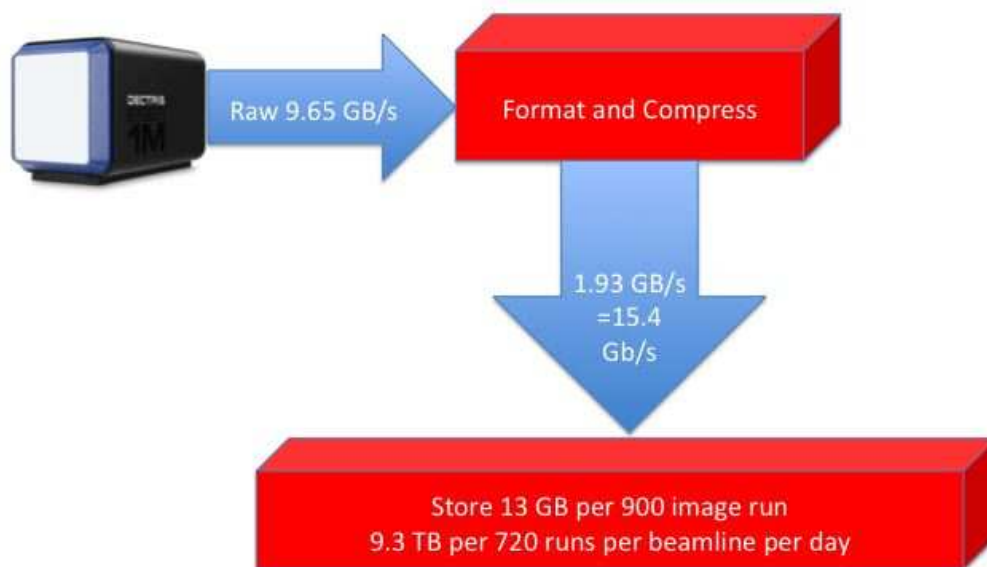
Fig. 2. Expected data flows at an NSLS-II beam line with an EIGER X 16M detector.

Table 2. *Current compression choices for MX data*

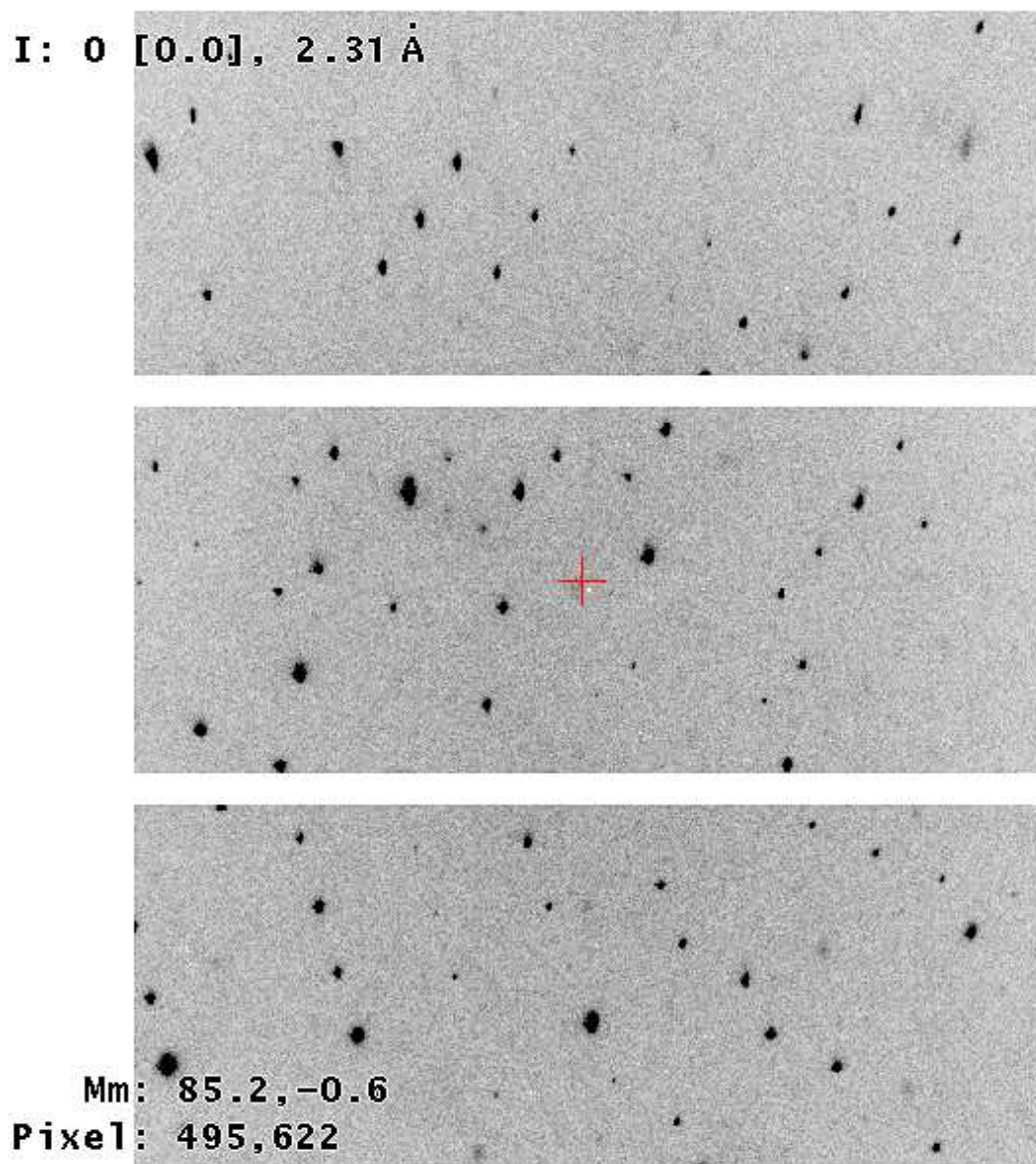| Compression Method | Compression Ratio | Relative Time |
|---|---|---|
| external bzip2 compression | 20.4:1 | 5.6 |
| HDF5 CBFlib canonical compression | 15.7:1 | 3.9 |
| HDF5 CBFlib nibble offset compression 2 Fast 6M | 11.5:1 | 2.9 |
| HDF5 CBFlib packed V2 compression 3 6M | 11.0:1 | 2.8 |
| HDF5 zip compression 16M | 9.7:1 | 2.4 |
| external LX4 compression (C1 one pass) | 8.7:1 | 2.2 |
| HDF5 CBFlib packed compression | 8.6:1 | 2.2 |
| external LX4 compression (C0 2 passes) | 5.2:1 | 1.3 |
| HDF5 CBFlib byte offset compression | 4.0:1 | 1.0 |

Fig. 3. DECTRIS Pilatus 300K image

Table 3. *Compressions of a DECTRIS Pilatus 300K image Fig 3*

| Compression | CBF size (MB) | HDF5 size (MB) |
|---|---|---|
| raw binary | 1.212 | 1.296 |
| byte offset | 0.309 | 0.393 |
| HDF5 zlib | n/a | 0.370 |
| nibble offset | 0.207 | 0.290 |
| packed | 0.184 | 0.267 |
| canonical | 0.178 | 0.262 |
| external bzip2 | 0.164 | 0.169 |

## Synopsis

This is a report on progress in the integration of CBF, NeXus and HDF5 in support of the high data rate/high data volume (BIG DATA) demands of the new generation of X-ray pixel array detectors.