Metadata for raw data from X-ray diffraction and other structural techniques

A Satellite Workshop to the 29th European Crystallographic Meeting

Practical session: 2. Towards a distributed / virtual repository of primary datasets

DDDWG Work Plan

Metadata

DDDWG Triennial Report 2011-2014

- In the Madrid Congress the Commissions were charged at the DDD inaugural meeting to define the metadata that should accompany their raw data ...
- While each IUCr Commission needs to specify 'technical' metadata *i.e.* those specific to their experimental raw data there is also a need to review 'generic' metadata *e.g.* who 'owns' a data set, details of research grants, embargo periods etc. A higher-level classification of the domain of study may be needed. *E.g.* a synchrotron facility might need to define different data storage policies for, say, X-ray diffraction images *versus* X-ray tomography images. Such policies could be automatically implemented if data sets had characteristics identifying what sort of scientific study they represent.

DDDWG Work Plan

Metadata

DDDWG Meeting Montreal 2014

- A centralised crystallographic repository of raw dataset metadata should be scoped and piloted.
- With such a repository in place, we should revisit the proposal that authors shall provide a permanent and prominent link from an article to the associated raw datasets.

Raw crystallographic datasets *are* now beginning to be deposited at a number of locations.

• Store.Synchrotron (MyTARDIS)

	16 14262 51.3 GB 6th February 2013 Put
Toggle Full Description Mirrored from http://rawdata.chem.uu.nl/	16 Datasets
Description Metadata Sharing Transfer Datasets	Just start typing to filter datasets based on descriptions
Institution Utrecht University	4DD0
Creative Commons Attribution 3.0 Australia (CC BY 3.0).	
Administrators Steve Androulakis Download All	
	4DD1 unwarp
	4DD2
	4DD3

Raw crystallographic datasets *are* now beginning to be deposited at a number of locations.

- Store.Synchrotron (MyTARDIS)
- Zenodo

zenodo

Search Communities

Browse - Upload Get started -

16 June 2011

Dataset Open access

Simultaneous X-ray diffraction from multiple single crystals of macromolecules

Paithankar, Karthik ; Sørensen, Henning ; Wright, Jonathan ; Schmidt, Soren ; Poulsen, Henning ; Garman, Elspeth

(show affiliations)

X-ray diffraction datasets from the publication: K. S. Paithankar, H. O. Sørensen, J. P. Wright, S. Schmidt, H. F. Poulsen and E. F. Garman*

Acta Cryst. (2011). D67, 608-618 doi:10.1107/S0907444911015617

The potential in macromolecular crystallography for using multiple crystals to collect X-ray diffraction data simultaneously from assemblies of up to seven crystals is explored. The basic features of the algorithms used to extract data and their practical implementation are described. The procedure could be useful both in relation to diffraction data obtained from intergrown crystals and to alleviate the problem of rapid diffraction decay arising from the effects of radiation damage.

Files			*
Name	Date	Size	
HEWL_7_crystals.tar.lzma	13 Aug 2014	918.4 MB	a Download
HEWL_4_crystals.tar.lzma	13 Aug 2014	913.7 MB	📥 Download
HEWL_3_crystals.tar.bz2	13 Aug 2014	1.7 GB	a Download
HEWL_2_crystals.tar.lzma	13 Aug 2014	1.9 GB	a Download
README.txt	13 Aug 2014	906 Bytes	a Download
insulin_4_crystals.tar.lzma	13 Aug 2014	964.4 MB	a Download
insulin_3_crystals.tar.lzma	13 Aug 2014	1.8 GB	📥 Download
insulin_2_crystals.tar.lzma	13 Aug 2014	1.8 GB	📥 Download

Research. Shared.

🔊 Sign In 🛛 🐼 Sign U

Publication date: 16 June 2011

DOI

DOI 10.5281/zenodo.11277

Keyword(s):

macromolecular crystallography methods development radiation damage

Collections:

Communities > Macromolecular Crystallography Datasets

Open Access

License (for files):

Creative Commons Attribution Share-Alike

Uploaded on:

13 August 2014



Select citation style...

Raw crystallographic datasets *are* now beginning to be deposited at a number of locations.

- Store.Synchrotron (MyTARDIS)
- Zenodo
- University of Manchester eScholar



The University of Manchester Library, The University of Manchester, Oxford Road, Manchester, M13 9PP, UK. | Contact details | Feedback The University of Manchester, Royal Charter Number: RC000797

Raw crystallographic datasets *are* now beginning to be deposited at a number of locations.

- Store.Synchrotron (MyTARDIS)
- Zenodo
- University of Manchester eScholar
- eCrystals / Atlas data store
- Protein Data Bank
- Wladek Minor Laboratory, U. Virginia
- Experimental facilities

Raw crystallographic datasets *are* now beginning to be deposited at a number of locations.

~ ~ C m	Scripts.iuc	1.org/cgr-bin/sendsup?aj3204	
JAC Journ	al of Appl	ied Crystallography	IUC
A			search IUCr Journ
home archiv	<i>r</i> e editors	for authors for readers submit subscribe open access	
JAC RESEARCH PA J. Appl. Cryst. (20	PERS 13). 46 , 108-119	3	
doi:10.1107/5002	1889812044172		OPEN 3 ACCESS
	Experience v a series of ly	rith exchange and archiving of raw data: comparison of data from two diffractometers and four software packag sozyme crystals	es on
	S. W. M. Tanle	ry, A. M. M. Schreurs, J. R. Helliwell and L. M. J. Kroon-Batenburg	
	A systematic a diffraction data associated me	nalysis of diffraction data of 11 different lysozyme crystals (used for cisplatin and carboplatin binding studies), obtained with fou processing software packages and from two diffraction diffractometers, serves as a pilot study for archiving raw diffraction data adata. The availability of the raw diffraction images allows for independent assessment of software packages.	r and
	Keywords: da	ta exchange; data archiving; metadata.	
	Read article	Similar articles	
			_
	Supporting	information	
	LINK	Link http://dx.doi.org/10.15127/1.219230 Raw data: PDB code 4dd0 HEWL_cisplatin_aqueous_glycerol	
	LINK	Link http://dx.doi.org/10.15127/1.215883 Raw data: PDB code 4dd1; HEWL_cisplatin_aqueous_paratone	
	LINK	Link http://dx.doi.org/10.15127/1.219240 Raw data: PDB code 4dd2; HEWL_carboplatin_aqueous_glycerol	
	LINK	Link http://dx.doi.org/10.15127/1.219241 Raw data: PDB code 4dd3; HEWL_carboplatin_aqueous_paratone	
	LINK	Link http://dx.doi.org/10.15127/1.219233 Raw data: PDB code 4dd4; HEWL_cisplatin_DMSO_glycerol	
	LINK	Link http://dx.doi.org/10.15127/1.219236 Raw data: PDB code 4dd6; HEWL_cisplatin_DMSO_paratone	http://o
	LINK	Link http://dx.doi.org/10.15127/1.219242 Raw data: PDB code 4dd7; HEWL_carboplatin_DMSO_glycerol	data: F
	LINK	Link http://dx.doi.org/10.15127/1.219257 Raw data: PDB code 4dd9; HEWL_carboplatin_DMSO_paratone	
	LINK	Link http://dx.doi.org/10.15127/1.219259 Raw data: PDB code 4dda; HEWL_NAG	
	LINK	Link http://dx.doi.org/10.15127/1.219238 Raw data: PDB code 4ddb; HEWL_cisplatin_DMSO_glycerol_pH6.5	
	LINK	Link http://dx.doi.org/10.15127/1.219260 Raw data: PDB code 4ddc; HEWL_cisplatin_NAG_7.5%_DMSO	
	LINK	Link http://rawdata.chem.uu.nl/≠0001 Raw data: archive at Utrecht University containing images measured at Manchester University	
	LINK	Link http://vera183.its.monash.edu.au/protein_cisplatin_carboplatin Raw data: mirror of the raw data from Tardis at Monash University	

Tanley, S. W. M., Schreurs, A. M. M., Helliwell, J. R. & Kroon-Batenburg, L. M. J. (2013). Experience with exchange and archiving of raw data: comparison of data from two diffractometers and four software packages on a series of lysozyme crystals. *J. Appl. Cryst.* **46**, 108–119.

Link http://dx.doi.org/10.15127/1.219230 Raw data: PDB code 4dd0 HEWL_cisplatin_aqueous_glycerol

Link http://dx.doi.org/10.15127/1.219230 Raw data: PDB code 4dd0 HEWL_cisplatin_aqueous_glycerol



LINK

LINK

Link http://rawdata.chem.uu.nl/#0001 Raw data: archive at Utrecht University containing images measured at Manchester University

⊢ → C 🖌 🙆 https://www.esch C 🖌 🗋 scripts.iuc → C f [] rawdata.chem.uu.nl/c001/index.html Q ☆ = ← 1.Experience with exchange and archiving of raw data: comparison of data from two diffractometers and four software packages on a series MANCHESTER **Journal of App** The Un of lysozyme crystals The University of Manchester 0 3 Simon W. M. Tanley, Antoine M. M. Schreurs, John R. Helliwell and Loes M. J. Kroon-Batenburg Search resources Academic suppor Journal of Applied Crystallography, 2013, Volume 46, pages 108-119 RESEARCH PAPERS he University HEW reprint (PDF file, 1.8 Mb) J. Appl. Cryst. (2013). 46, 108-11 doi:10.1107/50021889812044172 Manchester Library Tanley, Search resources (Resea Snapshot Nr of Scans Images Expanded PDB Sample Image Tarfile(s) Size (Mb) Diffractometer Acces Size (Mb) Manchester eScholar cisplatin Experience a series of cisplatin Search X.tar.gz unpacks cisplatir cisplatir cisplatir in original format pna S. W. M. Tar into subdirectory X ► Browse by A systematic Metrics diffraction d 4DD0 4DD0_01_0001.osc ∆hstra 1465 associated n ▶ Help 360 4DD0.tar.gz 6191 Rigaku R_AXIS IV Abstract About Keywords: d Abstract archiving Diffractio represen with asso the bindir Contact us Read article 4DD2 4DD2_01_0001.osc 360 4DD2.tar.gz 2657 6191 diffractor FEEDBACK AND diffraction Pt compo with EVA B factors ENOUIRIES + Supporting 4DD3 4DD3 01 0001.osc 360 4DD3.tar.gz 2293 6191 noticea By makin diffraction Bibliog 4DD9 4DD9 01 0001.osc 360 4DD9.tar.gz 2716 6191 LINK Digita Man 4DDA 4DDA 01 0001.osc 180 4DDA.tar.gz 1036 3096 LINK 619 4DDB 4DDB 01 0001.osc 360 4DDB.tar.gz 2249 LINK LINK 614 4DD1.tar.gz -480 630 Exter 4DD1 4DD1_01_0001.sfrm 554 1236 4DD1unwarp.tar.gz LINK Bruker PLATINUM¹³⁵ 4DD4.tar.gz 591 792 LINK 4DD4 4DD4 01 0001.sfrm 77 Institut 4DD4unwarp.tar.gz 624 1564 Univer LINK Acade 1440 4DD6.tar.gz 1067 1464 4DD6 4DD6_01_0001.sfrm 1117 2897 4DD6unwarp.tar.gz LINK Record 1862 4DD7.tar.gz 1464 1913 4DD7 4DD7_01_0001.sfrm Mano LINK 4DD7unwarp.tar.gz 1566 3746 1480 proteumnac.p4p 1440 4DDC.tar.gz 1149 LINK 4DDC 4DDC 01 0001.sfrm 4DDCunwarp.tar.gz 1220 2897 | Disclaimer | Privacy | Copyright notice LINK The University of Ma Last modified: 6/19/2015, 1:29:22 PM

Link http://vera183.its.monash.edu.au/protein_cisplatin_carboplatin Raw data: mirror of the raw data from Tardis at Monash University



LINK

Requirements

- Identification
- Provenance
- Disambiguation
- Categorization
- Context
- Relationship
- Size
- Licence

Communicating with repositories

OAI-PMH

Protocol for Metadata Harvesting (OAI-PMH)

The OAI-Protocol for Metadata Harvesting (OAI-PMH) defines a mechanism for harvesting records containing metadata from repositories. ... The metadata that is harvested may be in any format that is agreed by a community ... although unqualified Dublin Core is specified to provide a basic level of interoperability. Thus, metadata from many sources can be gathered together in one database, and services can be provided based on this centrally harvested, or "aggregated" data. The link between this metadata and the related content is not defined by the OAI protocol. It is important to realise that OAI-PMH does not provide a search across this data, it simply makes it possible to bring the data together in one place. ... To provide services, the harvesting approach must be combined with other mechanisms.

Although **the OAI-PMH is technically very simple,** building coherent services that meet user requirements remains complex. The OAI-PMH protocol could become part of the infrastructure of the Web, as taken-for-granted as the HTTP protocol now is, if a combination of its relative simplicity and proven success by early implementers in a service context leads to widespread uptake by research organisations, publishers, and "memory organisations".

http://www.oaforum.org/tutorial/english/page1.htm

Communicating with repositories

Dublin Core

Dublin Core Metadata Element Set Version 1.1

- 1. Title
- 2. Creator
- 3. Subject
- 4. Description
- 5. Publisher
- 6. Contributor
- 7. Date
- 8. Туре
- 9. Format
- 10. Identifier
- 11. Source
- 12. Language
- 13. Relation
- 14. Coverage
- 15. Rights

- Identification Provenance Disambiguation Categorization Context Relationship
 - Size
 - Licence

Communicating with repositories

OAI-PMH

- *So*, Dublin Core **could** serve as a metadata descriptor language in repositories containing crystallographic data
- *But*...

The OAI-Protocol for Metadata Harvesting (OAI-PMH) defines a mechanism for harvesting records containing metadata from repositories. ... The metadata that is harvested may be in any format that is agreed by a community ... although unqualified Dublin Core is specified to provide a basic level of interoperability. Thus, metadata from many sources can be gathered together in one database, and services can be provided based on this centrally harvested, or "aggregated" data. The link between this metadata and the related content is not defined by the OAI protocol. It is important to realise that **OAI-PMH does not provide a search across this data**, it simply makes it possible to bring the data together in one place. ... **To provide services, the harvesting approach must be combined with other mechanisms**.

OAI-PMH has three potentially useful features

1. Sets

• A **set** is an optional construct for grouping items for the purpose of *selective harvesting*. Repositories *may* organize items into sets. Set organization *may* be flat, *i.e.* a simple list, or hierarchical. Multiple hierarchies with distinct, independent top-level nodes are allowed.

2. Friends

 The friends container . . . is used by repositories that want to point harvesters to other repositories, by listing their base URLs. Usage of the friends container is *recommended*; it may support harvesters in discovering the network-location of repositories.

3. Metadata negotiation

• **ListMetadataFormats** [is a verb] used to retrieve the metadata formats available from a repository. An optional argument restricts the request to the formats available for a specific item. If this argument is omitted, then the response includes all metadata formats supported by this repository. Note that the fact that a metadata format is supported by a repository does not mean that it can be disseminated from all items in the repository.

Recommendations

1. Sets

- Define one or more set specifications that characterise crystallographic data, and educate repository managers to use these consistently, *e.g.* subject:crystallography-images *or* discipline:crystallography:category:experimental:type:diffraction-images
- For: establishes and supports a natural classification/organizational hierarchy
- Against: requires construction of the classification scheme; must ensure impossibility of conflict with other *ad hoc* schemes

Recommendations

2. Friends

- Crystallographic repositories should include a friends container listing other known crystallographic repositories
- For: allows incremental growth of network from local knowledge
- Against: possible duplication or multiple misallocation of related repositories; may benefit from a central registry of known 'friends' in this discipline

Recommendations

3. Metadata negotiation

- The crystallographic community (this WG?) should define an appropriate 'middle-layer' metadata scheme for optimising handling of high-level crystallographic metadata by repositories
- For: The existence of such a scheme in a repository guarantees that you are retrieving the desired sort of data; freedom to build complex systems
- Against: need to devise, test and implement such a scheme; may be needlessly duplicating other, more general efforts [*e.g.* Jisc Research at Risk consultation identified need for a discipline neutral metadata scheme –
 http://researchdata.jiscinvolve.org/wp/2015/07/03/research-data-metadata/
 (thanks to Chris Gibson, U. Manchester)

Outcome

What should be the outcome of a mechanism for locating distributed data sets?

A central database of deposited data sets

- Identifier (DOI)
- Provenance
- Link to publication
- Link to structure deposition (CSD, PDB, COD, ICSD)
- Nature (raw, processed, derived; X-ray, electron, neutron; diffraction, microscopy, NMR)
- Quality (resolution, completeness, level of interest)
- Size
- Licence / access

Stakeholders

Who benefits? Who pays? Who develops? Who maintains?

IUCrData?