The vital role of Crystallographic Information Files in chemical and biological crystallography to underpin the databases' validation reports

Brian McMahon

International Union of Crystallography, 5 Abbey Square, Chester CH1 2HU, UK

Email: bm@iucr.org



The clue is in the name...

66

There is a traditional hierarchy of components of understanding: **data**, **information**, **knowledge**, **wisdom** (the DIKW model).

Crystallographic **information** is the component that bridges the gap between the raw experimental **data** and the global **knowledge** bases represented by the Protein Data Bank, Cambridge Structural Database, International Centre for Diffraction Data, Inorganic Crystal Structure Database, Crystallography Open Database etc.

Summary: CIF and the *checkCIF* paradigm

- CIF for small molecules led to the *checkCIF* validation service
- checkCIF has a threefold approach:
 - Completeness
 - Correctness
 - Context
- CIF embraces all validation requirements (publication/database)
 - Handles metadata on same footing as data
 - Suitable for data query
 - Extensible

Types of crystallographic data described by CIF

'Raw' data

Numerical data collected directly from an experimental apparatus



Annotation



'Processed' data

Reduced, calibrated, processed numerical observations

h,k,l, Fc-squared, Fo-squared, sigma(Fo-squared) and status flag data_6 shelx title ' 01SRC413 in P2(1)/n' shelx_refln_list_code

shelx F calc maximum exptl crystal F 000 183.83 1144.00 refins_d_resolution_high 0.7705

symmetry_equiv_pos_as_xyz 'x, y, z' '-x+1/2, y+1/2, -z+1/2' '-x, -y, -z' 'x-1/2, -y-1/2, z-1/2'

cell length a 11.8293 10.3312 21.6318 cell length b cell length c cell_angle_alpha 90.000 cell_angle_beta 100.203 cell angle gamma 90.000

shelx_F_squared_multiplier 1.000

loop _refln_index_h _refln_index_k refln_index_1 refin_F_squared_calc refin_F_squared_meas refln F squared sigma refin_observed_status 772.37 856.47 1445.15 1446.80 1130.79 1097.08 1347.13 1490.27 3273.01 3545.64 12 48.20 40.50 14 79.87 63.02 2093.70 1975.83 33795.10 34884.29 2298.16 2035.72 9.73 36.06 449.80 506.89 1.81 43.36 28.81 64.18 48.51 10 1412.22 1628.54

242.68

14.96

16.87

16.46

2443.71

23397.80

20572.37

8854.88

0.00

11

12 13

14

15

2

2

'C	Derive	d' d	ata

Numerical description of the parameters of a calculated structure model



'Reference' data

Bibeo (

28.20 o

39.55 0

30.62 0

55.41 0

4.56 0

7.91 0

154.91

1287.71 c

38.24

5.59

11.92 c

5.59 0

6.79

6.02 0

9.70 0

3.84 0

4.56 0

7.91 0

5.59 0

61.27 o

546.30 c

520.01 o 169.57 o

45.96 0

7.91

279.96

10.52

15.76

7.91

3.95

2679.14

23770.90

19502.51

8282.53

Wyck	off Por	sitions	of Group P4/nnc (No. 126) [origin choice
Multiplicity	Wyckoff	Site	Coordinates
16	ĸ	1	$\begin{array}{l} (x,y,2) & (x+1/2,y+1/2,2) \; (y+1/2,x,2) \; (y,x+1/2,x) \\ (x+1/2,y,z+1/2) \; (x,y+1/2,z+1/2) \; (y,x,z+1/2) \; (y+1/2,x+1/2,z+1/2,y+1/2,z+1/2,z+1/2,y+1/2,z+1/2,z+1/2,y+1/2,z+1/2,$
8	J	2	(x,3/4,1/4) (-x+1/2,3/4,1/4) (3/4,x,1/4) (3/4,-x+1/2,1/4) (-x,1/4,3/4) (x+1/2,1/4,3/4) (1/4,-x,3/4) (1/4,x+1/2,3/4)
8	1	2	(x,1/4,1/4) (-x+1/2,1/4,1/4) (1/4,x,1/4) (1/4,x+1/2,1/4) (-x,3/4,3/4) (x+1/2,3/4,3/4) (3/4,-x,3/4) (3/4,x+1/2,3/4)
8	h	2	(x,x,1/4) (-x+1/2,-x+1/2,1/4) (-x+1/2,x,1/4) (x,-x+1/2,1/4) (-x,-x,3/4) (x+1/2,x+1/2,3/4) (x+1/2,-x,3/4) (-x,x+1/2,3/4)
8	9	2	(1/4,3/4,z) (3/4,1/4,z) (1/4,3/4,-z+1/2) (3/4,1/4,-z+1/2) (3/4,1/4,-z) (1/4,3/4,-z) (3/4,1/4,z+1/2) (1/4,3/4,z+1/2)
8	T.	-1	(0,0,0) (1/2,1/2,0) (1/2,0,0) (0,1/2,0) (1/2,0,1/2) (0,1/2,1/2) (0,0,1/2) (1/2,1/2,1/2)
4	0	4.	(1/4,1/4,z) (1/4,1/4,-z+1/2) (3/4,3/4,-z) (3/4,3/4,z+1/2)
4	d	-4.	(1/4,3/4,0) (3/4,1/4,0) (1/4,3/4,1/2) (3/4,1/4,1/2)
4	c	222 .	(1/4,3/4,3/4) (3/4,1/4,3/4) (3/4,1/4,1/4) (1/4,3/4,1/4)
2	b	422	(1/4,1/4,3/4) (3/4,3/4,1/4)
2	a	422	(1/4.1/4.1/4) (3/4.3/4.3/4)

of the point by its relative coordinates (in fractions or decimals Variable parameters (x,y,z) are also accepted

6 u | 🛛 I

Help

ments, please m

	x = y = z = Show	
	If you want to see the Wyckoff position in other setting, click here	
Server		For com admini

'Interpretative' data

Variable parameters in the experimental set-up or numerical modelling and interpretation

Commentary

ound I an

146 Searbox et al. + DARC. H. N.D. 82H.



N.O.J polyhedron [2.291 (3) Å, 2.442 (3) Å] and three elongated O(carboxylate) [2.437 (3)-2.703 (4) Å] bends (Table 1). The 12. in a chelating, bridging mode (μ_2, η^2, η^1) , sug



Acta Cryst. (2019). E75. 1145-1148

CIF dictionaries (COMCIFS)

- Crystallographic Core (coreCIF) 1991
- Crystallographic Restraints 2011
- Crystallographic Powder Diffraction (pdCIF) 1997
- Modulated and Composite Structures (msCIF) 2002
- Multipole Electron Density (rhoCIF) 2003
- Crystallographic Twinning 2014
- Magnetic Structures (magCIF) 2016
- Lattice topology (topoCIF) 2018
- Crystallographic Symmetry (symCIF) 2001
- Diffraction Images (imgCIF) 2000
- High pressure under development
- Crystallographic Macromolecular Structure (mmCIF) 1997

CIF dictionaries (wwPDB)

- Crystallographic Macromolecular Structure (mmCIF) 1997
- PDB Exchange Dictionary (PDBx/mmCIF) 1997 and ongoing
- Integrative/Hybrid (I/H) methods 2017
- 3DEM Extension Dictionary 2004
- NMRSTAR Dictionary 2013
- Biological Small Angle Scattering– 1998
- Model Archive Extension Dictionary 2018
- BIOSYNC Extension Dictionary 2000
- NMR Exchange Format Dictionary 2016

A coherent information flow



Early drivers for CIF

- Requirement to input raw data from many diffractometers
- Exchange data between software packages in the solution/refinement pipeline
- Provide a mechanism for electronic publication of articles
- Help referees to check the structural model
- Improve the accuracy of data capture by databases

Standardisation of raw data input

Championed by Howard Flack for point detectors

J. Appl. Cryst. (1992). 25, 455-459

DIFRAC, single-crystal diffractometer output-conversion software. By H. D. FLACK, Laboratoire de Cristallographie, University of Geneva, 24 quai Ernest-Ansermet, CH-1211 Genève 4, Switzerland and E. BLANC and D. SCHWARZENBACH, Institut de Cristallographie, University of Lausanne, BSP Dorigny, CH-1015 Lausanne, Switzerland

(Received 25 October 1991; accepted 9 January 1992)

Abstract

Software is described that converts single-crystal diffractometer output files as produced by maufacturers' software into a standardized instrument-independent form consisting of a clear, complete, well documented record of the sample and the diffraction measurements performed upon it. Information not already available in the manufacturers' diffractometer files is obtained in an interactive questionand-answer session. The software is written in a modular way. Available modules can deal with Enraf-Nonius CAD-4, Philips PW1100 and Siemens P2₁ single-crystal diffractometers and produce output in CIF or SCFS format.

Introduction

Users of several models of four-circle single-crystal diffractometers may well have been struck by the diversity of form and content of the data files generated by diffractometer-manufacturers' software. The form of the file can create difficulties in its transfer to other computing equipment and further renders the corresponding computer programs specific to a certain type or types of diffractometer. The difficulties in the content of the files manifest themselves principally by **the paucity of the available information** necessitating additional input to the datatreatment software (e.g. type of radiation, wavelength of radiation, scan width, ...). The problem has become aggravated in recent times by the rapid development in electronic data exchange. ... The advent of machinereadable submission for publication and data-base and supplementary-material deposition further highlights **the problem of missing data**.

Standardisation of raw data input

- Championed by Howard Flack for point detectors
- Relatively unenthusiastic uptake by vendors
- CIFs produced by commercial vendors often faulty
- imgCIF/CBF a second-wave attempt for diffraction images
- Perhaps slightly more effort from vendors
- 'mini-CBFs' showed significant diversity
- Reluctance to record essential experimental metadata

Data exchange between software packages

- STAR File modelled loosely on *Xtal* internal data structures
- Particular early adopters were collaborative packages (NRCVAX, CCP4) or general-purpose analysis programs (PLATON)
- Adoption by 'market leader' SHELXL crucial for success
- Reluctant uptake of mmCIF because of existing PDB format
- In general, has stood up well to the test of time
 - Generic presentations like XML not popular (but watch JSON!)
 - Limitations of PDB format eventually led to mmCIF as standard
 - Adoption of HDF5/NeXus in image capture driven by throughput needs
 - Relational nature of DDL dictionaries has become more important

Electronic publication

- Adoption by IUCr Journals (1991) very important
- Mandatory data format for Acta C by 1996
- Notes for Authors stipulate and validate mandatory metadata by listing required CIF data items
- Most journals reporting crystal structures now require CIFs as supporting information
- For online publication, availability of CIF makes the publication interactive

Technical peer review

- IUCr journals always required reviewers to check chemical structures
- CIF allowed ease of input to popular checking programs (*PARST, MISSYM, UNIMOL, CREDUC*)
- This led to semi-automated validation
 - IUCr checking reports
 - checkCIF dedicated tests
 - Consolidation with *PLATON*
- Elimination of 'Marshing'
- Subsequent ability to validate structure factors and redo refinement

Improved quality of structures in databases

- Improved quality of small-unit-cell structures from IUCr journals
- Convergence of checking criteria with CCDC
- Adoption of *checkCIF* for CCDC direct depositions
- wwPDB database structure isomorphic to mmCIF
- PDB requirement for structure factors in CIF format
- Validation Reports developed alongside extension CIF dictionaries for novel techniques

- CIF originally developed as tagged file format with one-dimensional loops
- York mmCIF workshop (1993) identified similarities with relational database structure
- DDL2 developed 1994/5 as a fully relational description of CIF data items
- The mmCIF/DDL2 schema is the basis for the relational database schema used by the wwPDB
- Latest iteration of DDLm retains this relational nature

_computational

COMMENTARY

Overhauling the PDB

Amanda C Schierz, Larisa N Soldatova & Ross D King

The Brookhaven Protein Data Bank was once a pioneering database, but its organization of structural data is now outdated and in need of an upgrade.

Although structural biology was once a Since the early 1970s, structural biology has been at the forefront of the development dards for the preservation and sharing of of standards for the preservation and sharing problems. scientific data and for database development, of scientific data. The PDB was set up in 1971 this lead has been lost. The main data standard (more than 10 years before the first sequence ing process was not a success and has led to used by PDB, mmCIF, does not meet state-ofthe-art standards in biology for ontologies; Structural biology journals were among the first this has serious repercussions for the analysis to require submission of data to international are important: they can result directly in and sharing of data, and has led to problems in database design. The main database, the reengineered Brookhaven Protein Data Bank (PDB), one of the first standards internationally agreed incomplete updating of the database due to although intended to be a relational database, upon for reporting experimental results. does not conform to relational database design standards and principles; this diminishes the become a 'relational' database based on the bility of designing intelligent analysis tools to quality of information stored in PDB and has serious implications for data-storage capacity. the PDB at the Research Collaboratory for knowledge discovery. If structural biology is to realize the full poten-2002 tial of its wealth of stored data and regain its lead in data standards and databases, then at the European Bioinformatics Institute both mmCIF and PDB need to be extensively (Hinxton Hall, UK), the Protein Data Bank including FMA (Foundational Model of reengineered. In this article, we detail some of the more serious faults we have found in the mmCIF and PDB. We then sketch a way to

databases, such as EMBL-Bank and GenBank). databases as a condition of publication, and the mmCIF dictionary³, published in 1997, was

mmCIF schema⁴. As wwPDB, it now includes analyze the data in PDB and aid the process of Structural Bioinformatics (RCSB, Brookhaven), the Macromolecular Structure Database (MSD) (Osaka University, Japan) and more recently Anatomy), which details the key concepts the Biological Magnetic Resonance Data Bank. and relations in anatomy⁵, and FuGO (The

Since the early 1970s, structural biology Our discussion here focuses mainly on the RCSB PDB, as we consider it to have the most

Our analysis indicates that the reengineerproblems of data repetition, redundancy, inconsistency and integrity. These problems incorrect answers to queries or more indirectly lead to erroneous results through the the formation of uncontrolled redundancies. The PDB was recently reengineered to They also seriously inhibit the future possi-

Why mmCLF is not a true ontology Several biological ontologies are now in use,



Global schema map of the entire PDB relational database; from Schierz, A. C., Soldatova, L. N. & King, R. D. (2007). Nature Biotechnology, 25, 437-442

The Brookhaven Protein Data Bank was once a pioneering database, but its organization of structural data is now outdated and in need of an upgrade.

mmCIF **RDBMS**

 Critique: too many tables in the database empty or nearly empty

 \rightarrow poorly defined schema

- *Reality*: each table represents a category of information needed to describe the structure completely
- *Problem*: depositors unwilling to provide that information (*i.e.* populate the tables)
- Most of the poorly populated tables capture experimental *metadata*

Metadata capture in macromolecular reports

← → C ① Not secure journals.iucr.org/f/services/structuralcommunications/		
Acta Crystallographica Section F STRUCTURAL BIOLOGY COMMUNICATIONS	← → C ① O Not secure journals.iucr.org/f/services/structuralcommunications/mmcifreqditems.html	
home archive editors for authors for readers submit subscribe open access		
data for structural and crystalling	Acta Crystallographica Section F	
	F STRUCTURAL BIOLOGY COMMUNICATIONS	
This page gives a list of recommended items for inclusion in structural and crystallization communications in Acta Crystallographica Section F.		
Items that are given in a magenta typeface are mandatory.	home archive editors for authors for readers submit subscribe open access	
The items will be published in tables that provide information on the sample and its treatment (including crystallization); data collection and struct refinement details.	data for structural and crystallization communications	
An online tool is available for preparing these tables from an mmCIF.	This page gives a list of recommended items for inclusion in structural and crystallization communications in Acta Crystallographica Section F.	
Click here for more details on these individual items and information on supplying the recommended information in mmCIF format. If your structure has been solved using NMR, click here for a separate set of recommendations.	The recommendations are tabulated below alongside the data names from the PDB mmCIF exchange dictionary available from the Protein Data Bank. These data items are provided in the mmCIF data sets created by the Protein Data Bank when a structure is deposited.	
	If you intend to submit to Acta Crystallographica Section F, you are recommanded to use the POB_EXTRACT utility available from the Protein Data Bank to extract a much as	
	possible of table minimized in our residue and organise of most continuous contraction tables and tables are not and tables for inclusion in our at tables are not and tables and tables for inclusion in your at tables.	
1. Sample information	The list below also includes examples to show how particular data will be organised in the mmCIF and how they will be arranged in the journal article.	
1.1. Macromolecule and source information	Click here for a more compact summary of the recommendations.	
Structure name Component molecules Additional molecular identifiers Biological functional unit (BFU) or macromolecular assembly, numbers and types of chains Mass of BFU (Da) Macromolecule sequence and chemical configuration Sequence database reference code Polymers (one-letter code sequence) or Polymer sequence as list of residues Ligand, cofactor, ions, solvent Mutations Post-translational modifications Formula weight of entity (Da) Source organism Scientific name Strain Details	 1: Sample information 1. Macromolecule and source information 1.2. Macromolecule production 1.3. Crystal data 1.4. Crystal data 2. Data collection and structure solution statistics 2.2.1. MAD/SAD data and structure solution statistics 2.2.2. MIR/MIRAS/SIR/SIRAS data and structure solution statistics 2.3. Model generation and refinement 4. Model validation 	
Source gene Scientific name	1. Sample information	
Strain Details	_symmetry.cell_setting _symmetry.Int_Tables_number	
	_symmetry.space_group_name_H-M Description	
Except macromolecule production		
PCR protocol Cloning protocol Expression protocol Purification protocol Additional details	1.1. Macromolecule and source information Example 1: complex of <i>E. coli</i> glutamate decarboxylase a with glutarate Example 2: a zinc-induced heterodimer of two isoforms of phospholipase A ₂ Example 3: mutation Example 4: mutation and modification	
	Structure name 🔨struct.title	
1.3. Crystallization 😶 Crystallization method	Component moleculesentity.pdbx_description ENTITY_NAME_SYS ENTITY_NAME_COM	
Temperature (K) Additional details	entity.pdbx_ec Biological functional unit (BFU) or	
	macromolecular assembly, numbersstruct_biol.details and types of chains	
	Mass of BFU (Da) _struct_biol.pdbx_formula_weight _struct_biol.pdbx_formula_weight_method	
	Macromolecule sequence and chemical configuration 0	

Authoring tools to assist metadata capture (*publBio* publisher)

- For MX a particular shortcoming is details of sample preparation and crystallization
- *publBio* offers a helpful interface
- Still a drought of crystallization papers
- Perhaps a role for *IUCrData*?

olBio data collection - Mozilla Firefox Precipitant solution Volume Volume units pH				
Components of the precipitar	t solution Concentra range in the prec	Crystallization screens (data kindly provided by Rigaku)		
polyethylene glycol 8000			Vendor Emerald BioSystems	
ammonium acetate ammonium bromide ammonium chloride Ammonium citrate - ammoniu hydroxide ammonium citrate - citric acid ammonium dihydrogen phosp Emmonium phosphate (mono ammonium phosphate monot ammonium phosphate, mono ammonium fluoride	m bhate basic) basic re p basic ra tra	precipitant solution, cli tion OR concentration	Well number 8 cadmium acetate (0.2 M) polyethylene glycol 8000 (20 w/v) Click any of the above components to add them to your table [the 'target' input box will be highlighted when you hover over the above item - click in any of the component fields (Name, Concentration) to change the 'target']	
ammonium iormate ammonium iodide ammonium nitrate ammonium sulfate		ervoir solution	Concentration units	
Ammoniumsulfate				

The checkCIF paradigm

Three aspects to *checkCIF* validation:

1. Completeness (the 'metadata' problem)

The metadata problem

- checkCIF for small structures achieves this reasonably well
 - For IUCr journals: published requirements in Notes for Authors
 - For other journals: in principle, customised 'request lists'; in practice, *de facto* acceptance of IUCr recommendations with *ad hoc* pick and choose
 - Possibility of better uptake through automated LIMS systems
- For protein structures
 - Reluctance in specifying mandatory data items
 - Inertia in design of standard experimental protocols in some facilities
 - Automated LIMS and other systems would help
- For raw data
 - CommDat sponsoring 'checkCIF for raw data' initiative

The checkCIF paradigm

Three aspects to *checkCIF* validation:

- 1. Completeness (the 'metadata' problem)
- 2. Internal self-consistency (relational methods)

Internal self-consistency

- IUCr *checkCIF* tests perform wide range of consistency checks
- Many of these also in *PLATON*
- Relational nature of DDL2 ensures category integrity within mmCIF
- DDLm offers a generic approach to specifying integrity relationships in dictionaries
- *dREL* a specific implementation of DDLm methods

The checkCIF paradigm

Three aspects to *checkCIF* validation:

- 1. Completeness (the 'metadata' problem)
- 2. Internal self-consistency (relational methods)
- 3. Comparison with related structures (the 'knowledge' problem)

The knowledge problem

- *PLATON* includes huge amount of crystallographic and chemical knowledge
- *Mogul* offers opportunity to compare new structures with existing curated ones in CSD
- PDB Validation Reports build on existing knowledge bases
- These tools are tuned to molecular geometry and other chemical properties
- Prospect of interrogating any ensemble of CIFs for any other defined properties through a non-specific API

Pain points for CIF in the 21st century

- Reluctance/difficulty in supplying experimental metadata
- Reticence in revising core dictionaries
- Lax approach to detailed standards in some software packages
- Lack of resources in developing CIF2/DDLm applications
- Lack of resources for validation software maintenance and development
- None are critical many show some (but slow) progress
- Important to keep momentum and train new generation

66

There is a traditional hierarchy of components of understanding: **data**, **information**, **knowledge**, **wisdom** (the DIKW model).

Crystallographic **information** is the component that bridges the gap between the raw experimental **data** and the global **knowledge** bases represented by the Protein Data Bank, Cambridge Structural Database, International Centre for Diffraction Data, Inorganic Crystal Structure Database, Crystallography Open Database etc.

This school will teach students to respect their raw data, extract the most reliable information they can, and disseminate that information in a complete and verifiable manner. In this way they will contribute to the sum total of scientific knowledge with rigour and integrity.

Crystallographic wisdom is outside the scope of this course.



Crystallographic information in the FAIR era

