

Championing data standards in chemical crystallography with CIF

Kamil Filip Dziubek

Institute of Mineralogy and Crystallography University of Vienna, Austria

IUCr Representative to CODATA
IUCr Committee on Data (CommDat) member













International Union of Crystallography: a scientific union

The International Union of Crystallography (IUCr) is a non-profit scientific union serving the worldwide interests of crystallographers and other scientists employing crystallographic methods

Paul Ewald of the Poly-

technic Institute of Brooklyn, President of the International Union of Crystallography (a scientific union, not a labor union)

Elizabeth A. Wood, *Crystals and Light: An Introduction to Optical Crystallography*, D. Van Nostrand, Princeton 1964.

- 55 National Committees make up the General Assembly
- The scientific work is conducted through Commissions and Committees
- The **Committee on Data (CommDat)** work with the IUCr's Commissions, including the Commission on Journals, having a coordinating and advisory role regarding data
- CommDat exist alongside, and have a formal relationship with, the Committee for the Maintenance of the CIF Standard (COMCIFS)
- The **Committee for the Maintenance of the CIF Standard (COMCIFS)** oversees the development of the CIF and reports to the Executive Committee of the IUCr
- Drafts of new definitions and official dictionaries are submitted to COMCIFS for technical validation and ultimately for approval





Simon Coles CommDat chair



James Hester COMCIFS chair



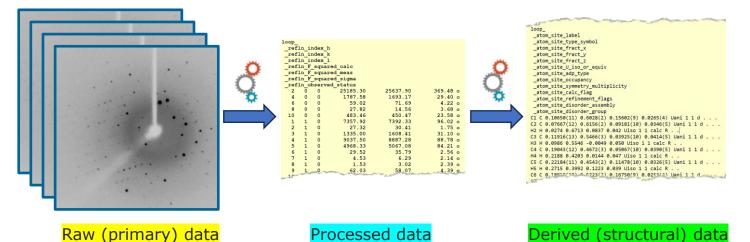




Flavors and data cycle of crystallographic data

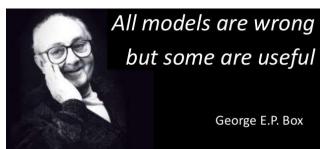
Crystallographic data lifecycle

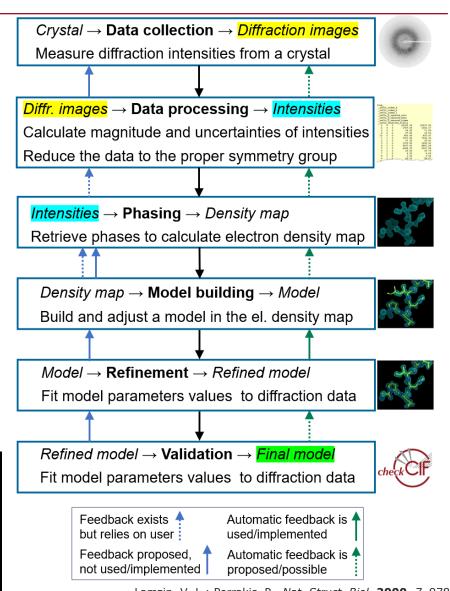
- Raw data images collected in a diffraction experiment
- Processed data (unit-cell parameters and indexed Bragg reflections with associated integrated intensities)
- Crystal structure (symmetry, unit-cell parameters, atomic coordinates, ADPs)



Crystal structures are models!

- The actual raw data in X-ray crystallography are diffraction data leading to electron density
- Atom positions that make up a crystal structure are a model that attempts to explain the raw data in as complete a fashion as possible

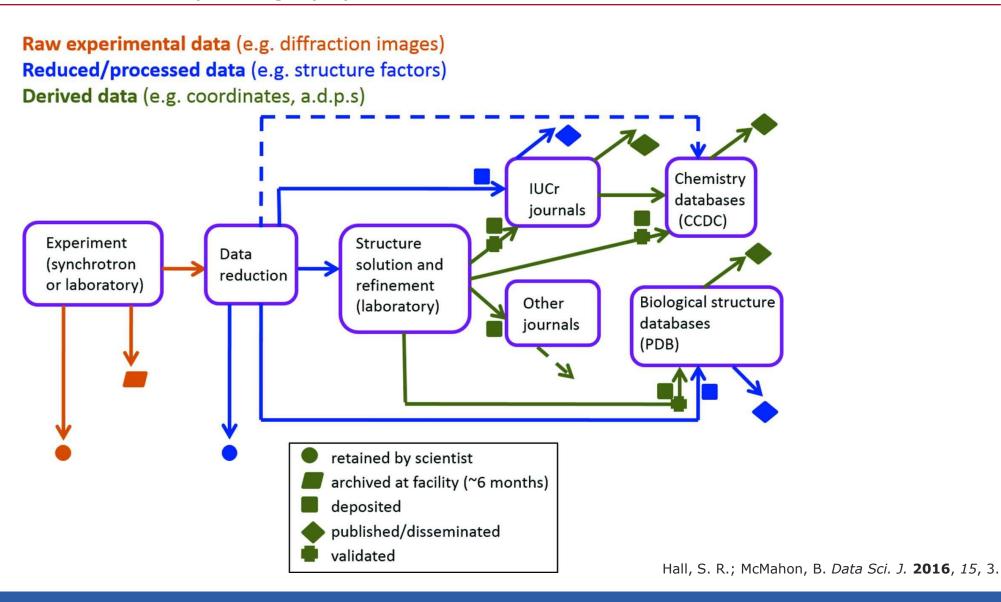






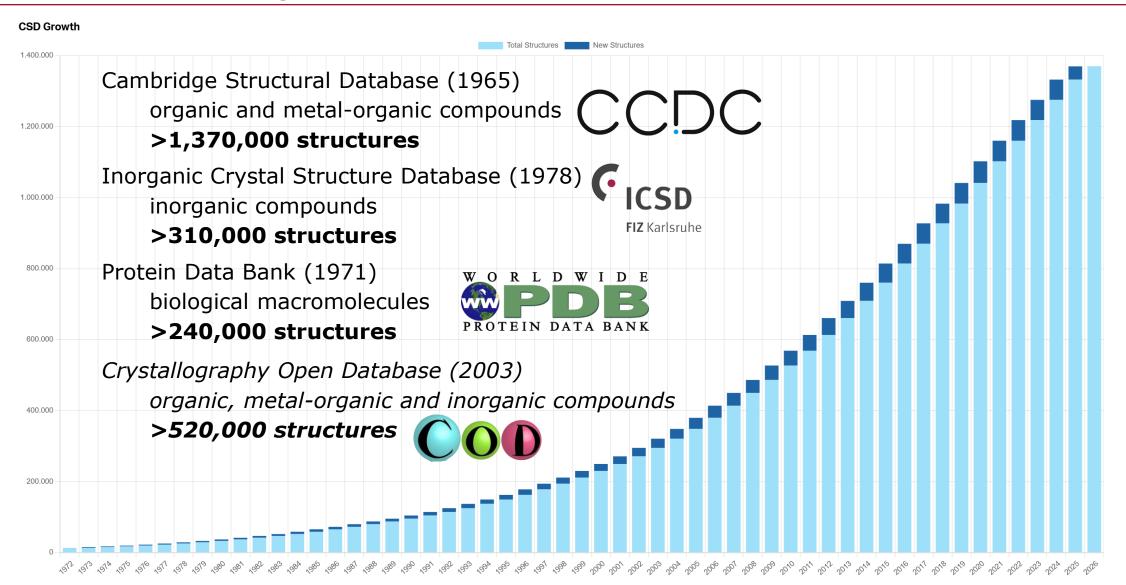


Coherent information flow in crystallography





Databases in the data deluge era



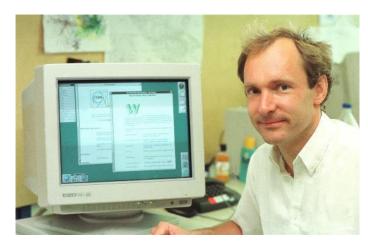
In 1991...



...USSR passes into history



...Operation Desert Storm began



...World Wide Web publicly debuts (HTML, XML and Java were yet to come!)

Acre Cruz, C991), A47, 653-683

International Union of Crystallography

Commission on Crystallographic Data
Commission on Journals
orking Party on Crystallographic Information

he Crystallographic Information File (CIF): a New Standard

BY SYDNEY R. HALL

ystallography Contro, University of Western Australia, Nodlands 6009, Australia

FRANK H. ALLEN

Crystallagraphic Data Center, University Chemical Laboratory, Lengheld Road, Combridge CB2 LEW, England

Institute for Moverials Resourch, McMaster University, Hamilton, Ontario L85 4M1, Canada

Abstract

The apecification of a new sandard Crystallogophic Information Fase (CIFF) is described, fits development is based in the Self-Defining Text Authors and Retrieval (STAR) exceedure [Hail (1991), J. Garn, Job. Comput. 87: 31, 305–313). The CIFF is a poursel, finishes and easily excensible free format archive fac; if is human and machine consider and one to childrel by a simple.



...IUCr introduced the Crystallographic Information File

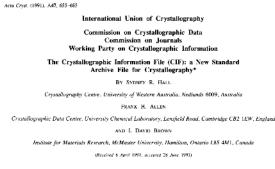




What is CIF and how is it used?

Crystallographic Information Framework (CIF) is a data exchange format for crystallographic and related structural science data

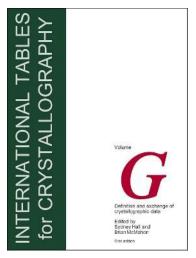
- The CIF format and extensible CIF dictionaries are defined, adopted and curated by the IUCr
- Until 2005 the acronym stood for 'Crystallographic Information File', but the name was modified in recognition of its application across crystallography and all related structural science fields
- CIF offers domain-specific ontology (collection of data identifiers, attributes and relationships)
- Simple syntax well suited to archive and exchange purposes (human- and machine-readable)
- Adopted by major databases (CSD, ICSD, COD, RRUFF, Mindat) and required as ESI in most journals
- Can also be adopted for the hosting structure factors / reflection intensity data
- imgCIF/CBF can handle binary data (including raw data images) and is isomorphous to CIF
- Canonical information on CIF contained in International Tables Vol. G



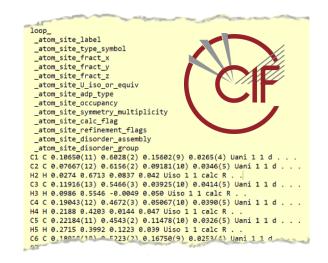
Abstract

Introduction

The specification of a new standard Crystallographic laformation File (CIF) is described. Its development is based on the Sel-Defining Text Archive and Retrieva (STAR) procedure [Hall (1991). J. Chen. Ing. Comput. Sci. 33, 326–333]. The CIF is a general flexible and easily extensible free-format archive file; it is human flexible and easily extensible free-format archive file; it is human flexible and easily extensible free-format archive file; it is human flexible and easily extensible free-format archive file; it is human flexible and easily extensible CF of the design of the electronic data requirements are well defined. Problems of duta ex-







loop_					
refln index h					
refln index k					
refln index l					
refln F squared calc					
_refln_F_squared_meas					
_refln_F_squared_sigma					
refln_observed_status					
_ 2	0	0	25185.30	25637.90	369.48 o
4	0	0	1787.58	1693.17	29.40 o
6	0	0	59.02	71.69	4.22 o
8	0	0	27.82	14.56	3.68 o
10	0	0	483.46	450.47	23.58 o
1	1	0	7357.92	7392.33	96.02 o
2	1	0	27.32	30.41	1.75 o
3	1	0	1335.00	1608.41	31.10 o
4	1	0	9037.50	8887.28	88.78 o
5	1	0	4968.33	5067.08	84.21 o
6	1	0	29.52	35.79	2.56 o
7	1	0	4.53	6.29	2.14 o
8	1	0	1.53	3.02	2.39 o
9	1.	_0 _	62.03	58.07	4.39 o
1 man and a second seco					



Wide applicability of STAR

CIF is based on the Self-Defining Text Archive and Retrieval (STAR) procedure

- STAR language has the capacity to handle the more complex data representations found in various disciplines such as imaging, quantum chemistry, botany, etc.
- Non-CIF STAR applications include
 - NMR imaging field (Biological Magnetic Resonance Data Bank)
 - Molecular Information File (MIF)
 - QCHEM ontology for quantum chemistry
 - Botanical ontologies (Florabase, Western Australian Herbarium), e.g. plant naming hierarchies lend themselves to the use of 'definition' methods to connect these hierarchies

326

J. Chem. Inf. Comput. Sci. 1991, 31, 326-333

The STAR File: A New Format for Electronic Data Transfer and Archiving

SYDNEY R. HALL

Crystallography Centre, University of Western Australia, Nedlands 6009, Australia

Received October 2, 1990

A new type of format is proposed for the computer archiving and electronic transmission of text and numerical data. The Self-defining Text Archive and Retrieval (STAR) File uses standard ASCII text to specify both the data structure and the information. The syntax of this file is simple, and it may be easily interpreted visually or by computer. The STAR format is the basis for the Crystallographic Information File (CIF), which has been adopted by the International Union of Crystallography for the submission of data and text to crystallographic journals and data bases.



Sydney R. Hall





Extensibility of CIF: dictionaries

CIF can host not only structural data

- Formally, the CIF approach makes no distinction between 'data' and 'metadata' and is adaptable and extensible to any domain of structural science
- CIF can host more than just data (symmetry, unit-cell parameters, atomic coordinates, ADPs); it can also host structural, experimental, and processing details, as well as provenance, authorship, journal details, etc.
- CIF dictionaries provide a formal taxonomy of crystallographic terms and ideas that can be discipline-specific
 - Core dictionary (coreCIF) "set of data names designed to cover the requirements of archiving and exchanging raw and processed data and derived structural results"
 - Powder dictionary (pdCIF)
 - Modulated and composite structures dictionary (msCIF)
 - Electron density dictionary (rhoCIF)
 - Twinning dictionary
 - Magnetic dictionary (magCIF)
 - Topology dictionary (topoCIF)
 - Macromolecular dictionary (mmCIF)
 - Symmetry dictionary (symCIF)
 - Image dictionary (imgCIF)
 - Restraints dictionary







Data flavors: the raw, the cooked and the medium-rare

- Q: Which of the data types (1) raw, (2) processed or (3) derived data can be described using the CIF format?
- A: All of them! However, raw data requires imgCIF/CBF ontologies, the particular formats describing diffraction images.
- Q: Do the major structural databases accept raw diffraction data?
- A: Currently not, but most large-scale photon facilities host and make raw diffraction data publicly available after the embargo period. The in-house diffractometer data can be uploaded to free data-sharing platforms (e.g. Zenodo).
- Q: And how about the processed data (structure factors)?
- A: From 2011 the Cambridge Structural Database (CSD) and the Inorganic Crystal Structure Database (ICSD) strongly encourage the inclusion of structure factor data along with CIFs, in line with the recommendations by the IUCr. This service is also available in the Crystallography Open Database (COD). Besides, many journals accept CIF files with structure factors as Electronic Supplementary Information.





Raw diffraction data reuse: the good, the bad and the challenging A Satellite Workshop to the XXVI IUCr Congress

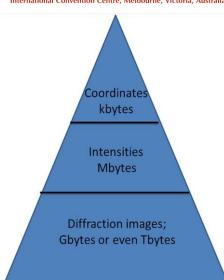
Organised by Loes Kroon-Batenburg, Selina Storm, John R. Helliwell and Brian McMahon under the auspices of the IUCr Committee on Data

The raw, the cooked and the medium-rare: unmerged diffraction data as a rich source of opportunities for data re-use and improvements in methods and results

G. Bricogne, C. Flensburg, R. H. Fogh, P. A. Keller, I. J. Tickle and C. Vonrhein

Tuesday 22 August 2023

International Convention Centre, Melbourne, Victoria, Australia





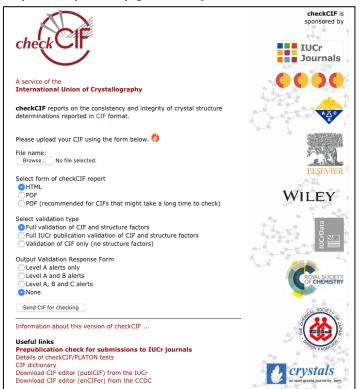


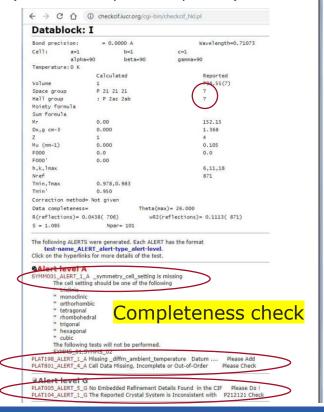
The element of trust: data validation

checkCIF is the IUCr service that reports on the completeness, consistency and integrity of CIF files

- Began as IUCr journals in-house suite for validation and consistency checks
- Consolidated with Ton Spek's PLATON (tool for analyzing geometrical calculations, e.g. bond lengths, angles, torsions and various structural tests, e.g. for missing symmetry, voids in the lattice) and now offered as public service

Adopted by many journal publishers and databases (checkCIF reports required!)







PLAT150_ALERT_1_A Volume as Calculated Differs from that Given ...
PLAT198_ALERT_1_A Missing_diffrn_ambient_temperature Datum ..
PLAT701_ALERT_1_A Bond Calc 1.384(5), Rep 1.368(5), Dev..

PLAT726_ALERT_2_G H...A Calc 1.74000, Rep 1.71000 Dev..

PLAT726_ALERT_2_G H...A Calc 1.78000, Rep 1.74000 Dev...

1.555 4.745 ...

1.555 2.665

1.555 1.555



Anthony L. Spek

PLAT701_ALERT_1_A Bond Calc 1.527(5), Rep 1.496(4), Dev.. 6.20 Sigma 1.555 1.555 # 13 Check PLAT702_ALERT_1_A Angle Calc 121.9(3), Rep 122.9(3), Dev.. O2 -C8 -O1 1.555 1.555 # 13 Che 1.555 1.555 1.555 # 13 Check PLAT707_ALERT_1_A D...A Calc 2.732(4), Rep 2.676(4), Dev.. 1.555 4.745 PLAT707_ALERT_1_A D...A Calc 2.747(3), Rep 2.689(3), Dev.. 1.555 2.665 # 19 Check PLAT701_ALERT_1_B Bond Calc 1.382(4), Rep 1.371(4), Dev... PLAT701_ALERT_1_B Bond Calc 1.395(4), Rep 1.385(5), Dev. C5 -C4 1.555 1.555 # 9 Check PLAT702_ALERT_1_B Angle Calc 123.6(3), Rep 122.9(3), Dev. O2 -C8 -C7 1.555 1.555 1.555 # 14 Check PLATO46 ALERT 1 C Reported 7 MW and D(calc) are Inconsistent PLAT068_ALERT_1_C Reported F000 Differs from Calcd (or Missing)... PLAT242_ALERT_2_C Low 'MainMol' Ueq as Compared to Neighbors of PLAT340_ALERT_3_C Low Bond Precision on C-C Bonds 0.0045 Ang. PLAT355_ALERT_3_C Long O-H (X0.82,N0.98A) O3 PLAT701_ALERT_1_C Bond Calc 1.195(5), Rep 1.186(4), Dev.. C8 -O2 1.555 1.555 # 11 Check PLAT702_ALERT_1_C Angle Calc 120.9(3), Rep 120.3(3), Dev.
C3 -C2 -C1 1.555 1.555 1.555 # 4 Char 1.555 1.555 1.555 # 4 Check PLAT702_ALERT_1_C Angle Calc 118.5(3), Rep 117.9(3), Dev.. C6 -C5 -C4 1.555 1.555 1.555 # 10 Check PLAT702_ALERT_1_C Angle Calc 120.0(3), Rep 120.5(3), Dev.. 1.555 1.555 1.555 # 10 Check C4 -C5 -C7 1.555 1.555 1.555 # 12 Check PLAT702_ALERT_1_C Angle Calc 119.1(3), Rep 119.5(3), Dev... PLAT702_ALERT_1_C Angle Calc 120.7(3), Rep 121.1(3), Dev.. C3 -C4 -C5 1.555 1.555 1.555 # 19 Check Alert level G PLATOOS_ALERT_5_G No Embedded Refinement Details Found in the CIF PLAT007_ALERT_5_G Number of Unrefined Donor-H Atoms . PLAT721_ALERT_1_G Bond Calc 1.04000, Rep 1.02990 Dev... PLAT721_ALERT_1_G Bond Calc 0.92000, Rep 0.89950 Dev.. 1.555 1.555 PLAT721_ALERT_1_G Bond Calc 1.01000, Rep 0.98500 Dev. O3 -H3A 1.555 1.555 # 18 Check PLAT721_ALERT_1_G Bond Calc 0.97000, Rep 1.555 1.555 ... PLAT722 ALERT_1_G Angle Calc 119.00, Rep 117.90 Dev... 1.555 1.555 1.555 # 17 Check PLAT722_ALERT_1_G Angle Calc 119.00, Rep 117.40 Dev... -Н7В PLAT725_ALERT_2_G D-H Calc 0.97000, Rep 0.95000 Dev -H1A 1.555 1.555 # 19 Check

19 Check

A typo in a cell parameter (b = 9.3092 Å instead of 9.0392 Å) generates inconsistencies and anomalies in bond geometry

Consistency check



Three prongs of the *check*CIF approach

FAIR principles do not, in themselves, cover the crucial aspects of intrinsic data quality and FAIR data are not per se high quality data! Three C's of checkCIF*:

- **Completeness**: Including minimal metadata (complete crystal description, details of the diffraction experiment, structure solution and refinement strategy), required for reproducibility
- **Correctness**: Integrity and internal self-consistency (e.g. space group must agree with lattice parameters)
- **Context**: Comparison with similar structures in the wider universe of knowledge (*e.g.* tests for missed symmetry, missed twinning, solvent accessible voids, and mis-assigned atom types)





*McMahon, B. The vital role of Crystallographic Information Files in chemical and biological crystallography to underpin the databases' validation reports, Data Science Skills in Publishing Workshop, 2019, Vienna.





checkCIF reports and alerts

The automated *check*CIF report contains three types of alerts:

- ALERT level A = In general: serious problem
- ALERT level B = Potentially serious problem
- ALERT level C = Check and explain





A. L. Spek, Acta Crystallogr. E 2020, 76, 1-11.

Try to eliminate the problems, but do not change the CIF or alter your refinement just to make *check*CIF alerts vanish!

Sometimes alerts cannot be eliminated, but instead have valid reasons for their presence. In these cases validation reply form (vrf) statements can be inserted into a CIF file, which acknowledge and account for checkCIF alerts.

```
# start Validation Reply Form
vrf ATOM007 Sb3N5 35 GPa
PROBLEM: atom site aniso label is missing
RESPONSE: Due to incompleteness of the high pressure dataset for a sample
in a diamond anvil cell ADPs of all atoms were refined
in isotropic approximation.
vrf PLAT029 Sb3N5 35 GPa
PROBLEM: diffrn measured fraction theta full value Low .
                                                               0.462 Why?
RESPONSE: This is a high pressure data set, a substantial part of
a reciprocal space was shaded by the diamond anvil cell.
vrf PLAT911 Sb3N5 35 GPa
PROBLEM: Missing FCF Refl Between Thmin & STh/L=
                                                    0.600
                                                                  80
Report
RESPONSE: Certain part of the reflections is missing due to shading
by the diamond anvil cell.
# end Validation Reply Form
```

```
The following ALERTS were generated. Each ALERT has the format test-name_ALERT_alert-type_alert-level.

Click on the hyperlinks for more details of the test.
```

Alert level A

ATOM007_ALERT_1_A _atom_site_aniso_label is missing Unique label identifying the atom site.

Author Response: Due to incompleteness of the high pressure dataset for a sample in a diamond anvil cell ADPs of all atoms were refined in isotropic approximation.

PLAT029_ALERT_3_A _diffrn_measured_fraction_theta_full value Low . 0.462 Why?

Author Response: This is a high pressure data set, a substantial part of a reciprocal space was shaded by the diamond anvil cell.

风 Alert level B

PLAT911_ALERT_3_B Missing FCF Refl Between Thmin & STh/L= 0.600 80 Report

Author Response: Certain part of the reflections is missing due to shading by the diamond anvil cell.





CIF and *check*CIF: Editors' requirements





Data sharing

Guidance for best practice and reproducibility of experimental data.

Recommended repositories

A data repository is an external storage space for researchers to deposit datasets associated with their research. Data should be submitted to a discipline-specific, community-recognised repository where possible, or alternatively to an institutional repository, or a generalist repository if no subject discipline repository is available for the given data type.

The choice of repository is the author's decision, provided it is in line with institutional or funder guidelines. The exception to this is small molecule crystal data, which must be deposited with the Cambridge Crystallographic Data Centre (CCDC).

Data availability statements

 Crystallographic data for [compound number] has been deposited at the [name of repository, such as CCDC / ICSD / PBD] under [accession number] and can be obtained from [URL of data record, format https://doi.org/DOI].





CIF and *check*CIF: Editors' requirements

X-Ray crystallography

These guidelines provide details for the presentation of single crystal and powder diffraction data; they apply to submissions to any of our journals.



Small molecule single crystal data

Authors should present their crystal data in a CIF (Crystallographic Information File) format and deposit this with the <u>Cambridge Crystallographic Data Centre</u> (CCDC) before submission. Data will be held in the CCDC's confidential archive until publication of the article, but it will be made accessible to reviewers and the publisher assigned to review the data.

At the point of publication, any deposited data will be made publicly available through the CCDC Access Structures service. In addition, organic and metal-organic structures will be curated into the Cambridge Structural Database, and inorganic structures will be curated into the Inorganic Crystal Structures Database (FIZ Karlsruhe).

Upon deposition, each data set is assigned a Digital Object Identifier (DOI), so that the crystal structure is unambiguously identified and registered.

Include CCDC or ICSD numbers in the manuscript prior to submission as part of a Data Availability Statement. During submission authors will be asked to cite CCDC or ICSD reference numbers; CIFs should not be submitted with the manuscript. Any revised CIFs should be deposited directly with the CCDC before the revised manuscript is submitted to us.

CheckCIF

In addition, authors are required to provide a checkCIF report for their reported crystal data. The checkCIF report can be obtained via the International Union of Crystallography's (IUCr) free checkCIF service, or as part of the CCDC deposition process. Any 'level A' alerts in the report should be explained in the submission details for the article or an explanation

provided within the submitted CIF. Authors should submit the checkCIF reports to the Royal Society of Chemistry along with the manuscript files.

If the editor deems it necessary during the peer-review process, the crystallography associated with the manuscript may undergo specialist crystallographic assessment, in which case a report will be provided along with the other reports from reviewers. Any points raised in this assessment should be attended to and all revised CIFs should be deposited with the CCDC prior to uploading the revised manuscript.

For recommended information to include in your CIF, please see the CCDC CIF <u>deposition</u> guidelines. If SQUEEZE or MASK procedures are used, this should be noted in the CIF file.

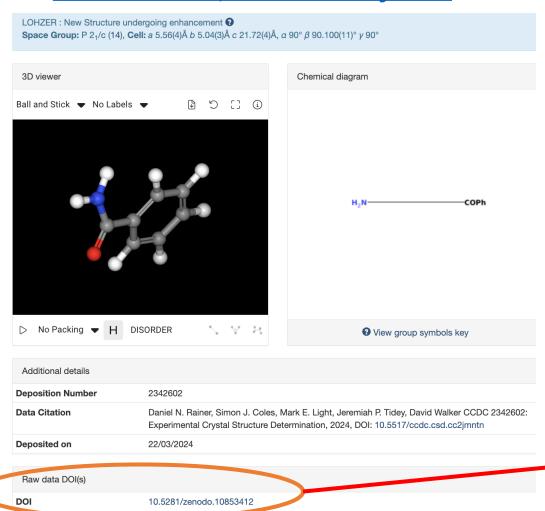
We encourage authors to include hkl data in the deposited CIF file. Alternatively authors can submit hkl data and the structure files (.fcf) separately during deposition with the CCDC. Raw data accompanying a structure should be made available by the authors for the review process, on request.

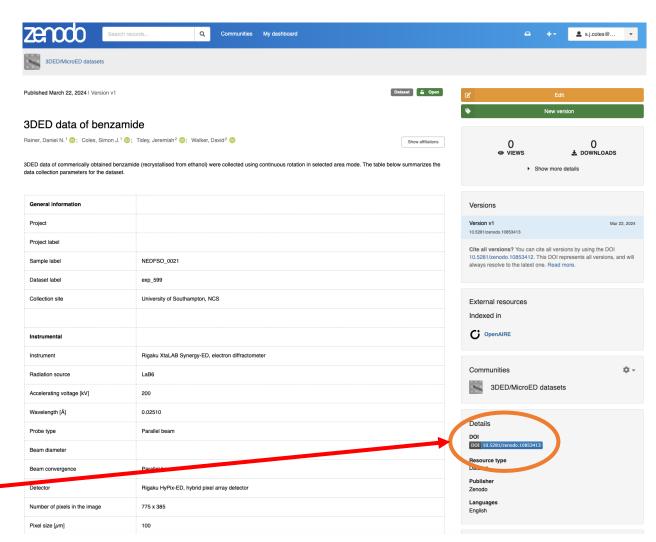




Raw data: cross-referencing between CCDC and external databases (e.g. Zenodo)

DOI: 10.5517/ccdc.csd.cc2jmntn





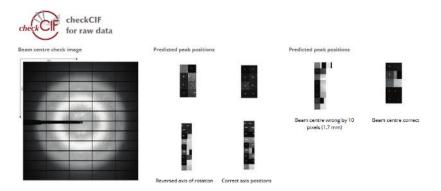




Continuing FAIRification of raw data within CIF

A brief history of raw data deposition

- Launch of *IUCrData* for peer-reviewed short structure reports (2016)
- A Gold Standard for metadata to describe an MX raw data (2020)
- Raw Data Letters, a new section of *IUCrData* for authors to describe unprocessed diffraction images (2022)
- checkCIF for raw data (2022)



research papers

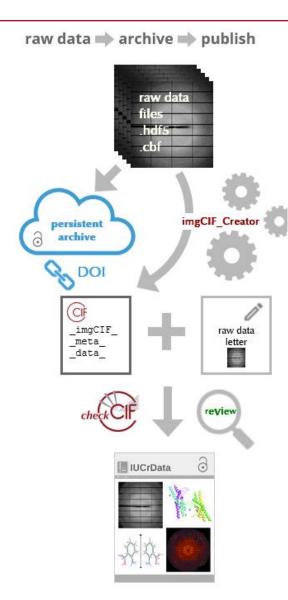


Accepted 26 June 2020

Gold Standard for macromolecular crystallography diffraction data

Herbert J. Bernstein, a* Andreas Förster, Asmit Bhowmick, Aaron S. Brewster,

Sandor Brockhauser, de,f Luca Gelisio, David R. Hall, Filip Leonarski, Valerio Mariani, Gianluca Santoni, Clemens Vonrheink and Graeme Winterh



EDITOR



Loes Kroon-Batenburg Bijvoet Centre for Biomolecular Research, Utrecht University, The Netherlands

chemical crystallography, macromolecular crystallography, raw data archiving

CO-EDITORS



Miguel Aranda University of Malaga, Spain

Rietveld method, powder diffraction



Elena Boldyreva Novosibirsk State University, Russia

high pressure



Aaron Brewster Lawrence Berkeley National Lab, USA

serial crystallography computational methods development



Simon Coles University of Southampton, UK

chemical crystallography, structural chemistry, databases



John R. Helliwell University of Manchester, UK

macromolecular crystallography, data science





New challenge: Crystal Structure Prediction standards and CIF dictionary

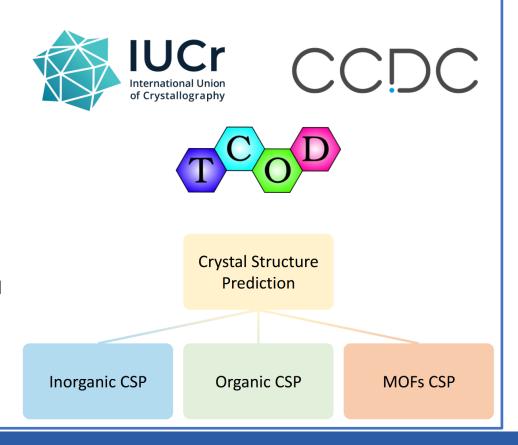
Number of theoretically calculated crystal structures far exceeds that of experimentally determined structures!

- In one year a draft of the final dictionary will be communicated to the IUCr and the finalized version will be officially published
- In the meantime feedback, input and approval from community is required
- Available at: github.com/COMCIFS/Structure Prediction Dictionary

CSP Data Standards

Scope of the project:

- An inclusive initiative that establishes community standards for predicted crystal structures
- Enables effective publication of results and easy search and retrieval across resources
- Supports levels of **reproducibility** deemed appropriate by different communities
- Flexible structure for future expansion as
 CSP methods are continuing to evolve





James Hester



Nicholas Francia



Ian Bruno





Potential of CIF to host non-crystallographic (spectroscopic) data

research papers





Received 23 November 2018 Accepted 28 March 2019

Edited by A. Borbély, Ecole National Supérieure des Mines, Saint-Etienne, France

Keywords: Raman spectroscopy: open databases; combined Raman-X-ray diffraction; DDLm dictionary; CIF2.

Raman Open Database: first interconnected Raman-X-ray diffraction open-access resource for material identification

Yassine El Mendili, a* Antanas Vaitkus, b Andrius Merkys, b Saulius Gražulis, b Daniel Chateigner, Fabrice Mathevet, Stéphanie Gascoin, Sebastien Petit, A Jean-François Bardeau, Marco Zanatta, Maria Secchi, Gino Mariotto, Arun Kumar, d Michele Cassetta, d Luca Lutterotti, e Evgeny Borovin, e Beate Orberger, f Patrick Simon, Bernard Hehlenh and Monique Le Gueni

a Normandie Université, CRISMAT-ENSICAEN, UMR6508 CNRS, Université de Caen Normandie, 6 Boulevard Maréchal Juin, 14050 Caen, France, ^bVilnius University, Institute of Biotechnology, Sauletekio av. 7, LT-10257 Vilnius, Lithuania, ^cInstitut des Molécules et Matériaux du Mans, UMR6283 CNRS, Le Mans Université, Avenue Olivier Messiaen, 72085 Le Mans, France, ^dDepartment of Computer Science, University of Verona, Strada Le Grazie 15, 37134 Verona, Italy, ^eDepartment of Industrial Engineering, University of Trento, via Sommarive 9, 38123 Trento, Italy, ^fGEOPS-Paris Sud, Université Paris-Saclay, UMR8148 (CNRS-UPS), Bâtiment 504, 91405 Orsay, France, 8CEMHTI, UPR CNRS 3079, Université d'Orléans, 1D Avenue de la Recherche Scientifique, 45071 Orléans Cedex 2, France, hLaboratoire Charles Coulomb, UMR5521 CNRS, Université Montpellier 2, Place Eugène Bataillon, 34095 Montpellier Cedex 5, France, and ⁱERAMET, 1 Avenue Albert Einstein, 78190 Trappes, France. *Correspondence e-mail: yassine.elmendili@ensicaen.fr

Detailed crystallographic information provided by X-ray diffraction (XRD) is complementary to molecular information provided by Raman spectroscopy. Accordingly, the combined use of these techniques allows the identification of an unknown compound without ambiguity. However, a full combination of Raman and XRD results requires an appropriate and reliable reference database with complete information. This is already available for XRD. The main objective of this paper is to introduce and describe the recently developed Raman Open Database (ROD, http://solsa.crystallography.net/rod). It comprises a collection of high-quality uncorrected Raman spectra. The novelty of this database is its interconnectedness with other open databases like the Crystallography Open Database (http://www.crystallography.net/cod and Theoretical Crystallography Open Database (http://www.crystallography.net/tcod/). The syntax adopted to format entries in the ROD is based on the worldwide recognized and used CIF format, which offers a simple way for data exchange, writing and description. ROD also uses JCAMP-DX files as an alternative format for submitted spectra. JCAMP-DX files are compatible to varying degrees with most commercial Raman software and can be read and edited using standard text editors.



Journal of

Towards data format standardization for X-ray

Synchrotron Radiation

ISSN 0909-0495

Received 28 March 2012 Accepted 26 August 2012 absorption spectroscopy

B. Ravel, a* J. R. Hester, b V. A. Solé and M. Newvilled

^aNational Institute of Standards and Technology, Gaithersburg, MD 20899, USA, ^bBragg Institute, ANSTO, Locked Bag 2001, Kirrawee DC, NSW 2232, Australia, ^cEuropean Synchrotron Radiation Facility, 6 rue Jules Horowitz, BP 220, 38043 Grenoble Cedex 9, France, and defenter for Advanced Radiation Studies, University of Chicago, Building 434A, Argonne National Laboratory, Argonne, IL 60439, USA. E-mail: bravel@bnl.gov

A working group on data format standardization for X-ray absorption spectroscopy (XAS) has recently formed under the auspices of the International X-ray Absorption Society and the XAFS Commission of the International Union of Crystallography. This group of beamline scientists and XAS practitioners has been tasked to propose data format standards to meet the needs of the world-wide XAS community. In this report, concepts for addressing three XAS data storage needs are presented: a single spectrum interchange format, a hierarchical format for multispectral X-ray experiment, and a relational database format for XAS data libraries.

© 2012 International Union of Crystallography Printed in Singapore - all rights reserved

Keywords: XAFS; standardization; data formats.

J. Synchrotron Rad. (2012). 19, 869-874

doi:10.1107/S0909049512036886

q2xafs workshop





Synergies between IUCr and other scientific standard organizations

CIF integrates standards developed within IUPAC

- Wavelength description uses the IUPAC nomenclature
- Enantioexcess is as per IUPAC recommendation
- InChi and InChIKey can be embedded in CIF
- Systematic chemical name and chemical formula follow IUPAC standards
- Topological CIF dictionary includes recommendations of IUPAC on nomenclature for coordination polymers

Liaisons between CIF and CDIF

Physical Sciences Data Infrastructure (PDSI) at Southhampton managed by Simon Coles strives to introduce CDIF (the CODATA Cross Domain Interoperability Framework) that integrates a wide range of resources across chemistry (with crystallography and CIF at the core)



Coalition for the Sustainability of Digital Data Standards in the Chemical Sciences

IUCr (representing the CIF standard) actively participates in the DigSustain meetings from the beginning







DigSustain2: Digital Data Standards Sustainability in Chemical Sciences, Delitzsch, Germany, April 2025



DigSustain3: Coalition for the Sustainability of Digital Data Standards in the Chemical Sciences Planning Meeting, London, UK, November 2025

Stay tuned for more details in the talk of Wendy Patterson!



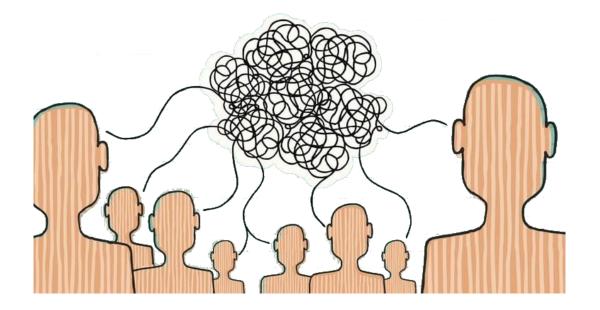
Where do we go from here?

Impact of CIF can be assessed by its popularity

- "The single most important impact of CIF has been the ability to easily access data globally across journals, databases and laboratories"
- CIF is now one of Universal Data Languages
- Peter Murray-Rust called CIF 'the best working data system in the scientific world' (2014)
- "Technology lock-in': foundational technology choice can become so embedded that it stifles future alternatives

Progress and adoption can be measured only indirectly

- Efficiency to submission–deposition processes
- Better understanding of particular data items across laboratories and databases
- Adoption in programming libraries and end-user software
- "Anyone can place their own tag-value items into a CIF, but until their definitions are accepted into the global ontology they can only be for local use"
- Official global ontology is the responsibility of the IUCr COMCIFS group, who moderates and maintain the standard
- Important to keep momentum and train new generations
- The art of understanding and being understood







Thank you for listening!

Acknowledgments







James Hester



John R. Helliwell



COMMITTEE ON DATA (COMMDAT)

MEMBERSHIP

- S. Coles (Chair, UK)
- H.J. Bernstein (USA)
- A. Brink (South Africa)
- I. Bruno (UK)
- . S. Coles (UK)
- K. Dziubek (Austria)
- · A. Goetz (France)
- S. Kabekkodu (USA)
- L.M.J. Kroon-Batenburg (The Netherlands)
- G. Kurisu (Japan)
- W. Minor (USA)
- S. Storm (Germany)
- L. Van Meervelt (Belgium)
- J. Hester (Australia) (COMCIFS liaison)

CONSULTANTS

- S. Androulakis (Australia)
- M.P. Blakeley (France)
- G. Bricogne (UK)
- S. Grazulis (Lithuania)
- J.R. Helliwell (UK)
- B. Matthews (UK)
- A. Sarjeant (USA)
- D. Szebenyi (USA)
- E.F. Weckert (Germany)
- J. Trewhella (Australia)