Introduction to the CIF Ontology

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

N IUCr



This talk introduces the 2023 COMCIFS Dictionary Writing Workshop at the Melbourne Congress. It serves two main functions: (1) to explain the historical development of CIF dictionaries to define in a machine-actionable manner the contents of data files covering various aspects of crystallography and related structural sciences; (2) to demonstrate some of the more complex types of information that can be handled with this approach. A simple controlled vocabulary is adequate to describe individual objects (*e.g.* the volume of a unit cell). However, to describe more complex objects, the relationships between the simple 'building blocks' need to be characterised. The great strength of the CIF approach is that the building blocks and their relationships can be handled within the same formalism – a formalism, moreover, that allows the way the relationships are described to be defined also in the same way. We use the term 'CIF ontology' here to cover the canonical dictionaries maintained by the IUCr and wwPDB. The purpose of the workshop is to help participants to develop these further, but also to create their own extension dictionaries (in structural science or, indeed, any other knowledge domain).



The 1991 CIF paper of Hall, Allen and Brown was not the IUCr's first attempt to standardise information interchange. It built on the 'Standard Crystallographic File Structure' (SCFS) of David Brown in the 1980s. Like SCFS, the first iteration of CIF had two components: a specific file format, and a list of identifiers ('tags') that should be used by all compliant software.

chemical formula			
_chemical_name	'C13 H12 05'		
; 3-(2,5-dihydro-4-h	ydroxy-5-oxo-3-pheny1	-2-furyl)propionic acid	
; _publication_title			
; Structure of WF-36	81, 3-(2,5-Dihydro-4-	hydroxy-5-oxo-3-pheny1-2-fury	1)propionic Acid.
;			
cell a	18,757(8)		
_cell_b	7.282(2)		
_cell_c	17.511(8)		
cell_alpha	90		
 _cell_beta	91.20(3)		
cell_gamma	90		
_cell_volume	2391(3)		
_symmetry_space_grou	p '-C 2yc'		
<pre>loopsymmetry_pos_</pre>	in_xyz		
'x,y,z'	'-x,-y,-z'	'-x,y,1/2-z'	
'x,-y,1/2+z'	'1/2+x,1/2+y,z'	'1/2-x,1/2-y,-z'	
'1/2-x,1/2+y,1/2-z'	'1/2+x,1/2-y,1/2+z'		October 22, 1989

The SCFS file format was very much in the spirit of Fortran, defining specific formatted data types and a certain ordering of presentation. The CIF format was much more fluid, allowing arbitrary whitespace to separate the identifying tags, now called 'data names', from their values. Further, there was no restriction on the order of presentation, except that closely related data values that would naturally be tabulated together were collected in a looped structure. (But even here, table columns could be permuted in any order.)

Each data name is recognized because it has a leading underscore. The value follows, separated by white space (the amount of white space does not matter). If the value itself *contains* white space, it should be wrapped in quote marks or surrounded by semicolons in the first column. Looped values are laid out as if in a table, where the data names are listed together after the loop_ keyword and the corresponding values follow in strict rotation.

data_crystal_structu	re		
_chemical_formula	'C13 H12 05'		
; 3-(2,5-dihydro-4-hy	ydroxy-5-oxo-3-phenyl	-2-furyl)propionic acid	
;			
_publication_title			
; Structure of WF-36	81, 3-(2,5-Dihydro-4-	hydroxy-5-oxo-3-pheny1-2-fur	yl)propionic Acid.
,			
_cell_a	18.757(8)		
_cell_b	7.282(2)		
_cell_c	17.511(8)		
_cell_alpha	90		
_cell_beta	91.20(3)		
_cell_gamma	90		
_cell_volume	2391(3)		
	p '-C 2yc'		
_symmetry_space_grou	in_xyz		
<pre>loopsymmetry_pos_</pre>		'-x v 1/2-z'	
<pre>_symmetry_space_group loopsymmetry_pos_ 'x,y,z'</pre>	'-x,-y,-z'	~, , , , +/ = =	
<pre></pre>	'-x,-y,-z' '1/2+x,1/2+y,z'	'1/2-x,1/2-y,-z'	

Each data name could occur only once in a file (strictly, in a *data block*; files could be partitioned into separate data blocks, each representing a standalone description of, say, a crystal structure or a diffraction experiment).



In the original paper, each data name was defined, alongside information about the values associated with it (*e.g.* type, numeric range constraints, permitted code terms *etc.*) in a layout similar to that of a conventional linguistic dictionary.



Notice how, in the manner of a lexicographic dictionary, the definition includes a number of attributes explaining how the term should be used. These attributes were initially somewhat *ad hoc*, but were stored in a CIF-format file that was the precursor of the dictionary paradigm we use today. This informal set of attribute definitions is sometimes called 'DDL0'.



By 1997, the attributes themselves were formally described using the same attribute set, and this was the first 'official' dictionary definition language, DDL1. Some of the early attributes permitted by the original CIF paper (*e.g.* ability to generate variant data names according to the physical units used) were omitted from DDL1. The decision was taken to enforce a single physical unit to be associated with each defined physical quantity. This did not exclude the possibility of defining new data names where the only difference was in the associated physical unit, but the practice was discouraged.

f	urther refined (20	13) as a relational o	ontology
	<pre>save_cell.length_a</pre>		
	_definition.id _alias.definition_id _name.category_id _name.object_id	'_cell.length_a' '_cell_length_a' cell length_a	
	_definition.update _description.text ; The length of each cell	2014-06-08 axis.	
	; _type.purpose _type.source _type.container _type.contents _enumeration.range _units.code	Measurand Recorded Single Real 1.: angstroms	
	save_		
8 IUCT	Introd	uction to the CIF Ontology	000000000000000000000000000000000000000

The similarity of the category/included data names structure that was established in DDL1 to the structure of a relational database became apparent to the working group developing the mmCIF dictionary for describing protein and nucleic acid structures. They developed a variant DDL (designated DDL2) that was almost completely isomorphous to a relational database schema. This is the basis for the current wwPDB implementation, is fully mature, and continues to be extended as needed to accommodate the specific requirements of the PDB and structural biology community. It won't be discussed here (but details can be found in the handout for the last COMCIFS workshop at https://www.iucr.org/resources/cif/comcifs/workshop-2017)

However, the relational approach was a significant enhancement to the original DDL, and in 2013 a new language, DDLm, was introduced that was capable of expressing the relational attributes that were in DDL2, but that also allowed the inclusion of algorithmic methods to express or calculate relationships between different data items. This new DDLm would also be able to describe the attribute behaviours of DDL1, and so was potentially universal in scope. A new structuring of data names was encouraged through a **_CATEGORY.item** construction (*i.e.* using a dot separator); but legacy files with different data names (all underscores) could be supported through an aliasing mechanism.



Some more historical background and a description of how CIF became established among different structural science communities is found in a relatively recent *IUCr Newsletter* article, which is also reprinted in the workshop booklet.

The CORE of the CIF ontology

- Scope of the original CIF specification: single-crystal structures
 - Description of the diffraction experiment (X-ray, neutron, electron)
 - Crystal characteristics (size, shape, composition, space group)
 - Diffracted beam intensities (point detectors, area detectors)
 - Structure solution and refinement (methodology, software)
 - Structure factors (allowing reprocessing, replicability)

ILC

- 3D structure (molecular geometry, displacement factors)
- Publication (authors, citations, database entries, commentary)

So we turn to an overview of the core CIF dictionary that was developed for the 1991 CIF article, and the extensions that have subsequently been built under the aegis of the IUCr and the wwPDB. These canonical dictionaries we are referring to collectively as the *CIF Ontology*.

Introduction to the CIF Ontology

Definitions organised by CATEGORY Category EXPTL Category EXPTL CRYSTAL FACE Example _exptl_crystal_appearance.general _exptl.crystals_number _exptl_crystal_face.index_h 'opaque' _exptl.method _exptl_crystal_face.index_k _exptl.method_details _exptl_crystal_face.index_1 loop_ _exptl.special_details _exptl_crystal_face.diffr_chi _exptl_crystal_face.index_h _exptl.transmission_factor_max _exptl_crystal_face.diffr_kappa _exptl_crystal_face.index_k _exptl.transmission_factor_min _exptl_crystal_face.diffr_phi _exptl_crystal_face.index_1 etc. __exptl_crystal_face.diffr_psi _exptl_crystal_face.perp_dist _exptl_crystal_face.perp_dist etc. -1 -1 0 0.725 -1 -1 1 0.764 Category EXPTL CRYSTAL APPEARANCE -1 0 0 0.548 -1 0 1 0.597 0 -1 -1 0.616 _exptl_crystal_appearance.general 0 0 -1 0.391 _exptl_crystal_appearance.hue 1.355 0 1 -1 _exptl_crystal_appearance.intensity Example from Kaminsky, W. (2007). From CIF to virtual morphology using the WinXMorph program. J. Appl. Cryst. 40, 382-385 **IIIC**

The essential characteristic of the CIF dictionaries is that related items are gathered together into categories. If you are storing CIF information in a relational database, then each category can simply be considered as a table. However, some items will only occur once (at least within the scope of the subject area covered by that dictionary). It is possible to display everything in tabular format, but in the case of these 'degenerate' or 'scalar' categories, it can seem more natural to list the items where the data value immediately follows the declared data name. This could, for example, allow items within such a category to appear at several locations within a file (potentially useful if data is accumulated at different stages along a workflow). Categories that are expected to have single-instance data names are indicated with the attribute **_definition.class 'Set'**, while the more typical multi-value (tabular) categories are identified with **_definition.class 'Loop'**.

Category hierarchy of core CIF dictionary	
 HEAD category : CIF_CORE (establishes 'root' of definition family) 6 'themes' DIFFRACTION EXPTL FUNCTION MODEL PUBLICATION STRUCTURE 	
12 Introduction to the CIF Ontology	

There is an organisational hierarchy amongst categories in the core and other CIF dictionaries, though relationships between categories are expressed through foreign keys and are not dependent on this notional hierarchy. There are also parent-child category relationships, which express 'projections' or instances when sub-categories are separated out into different tables to escape sparse presentations. Ultimately, relationships are all expressed through inter-category pointers independently of the notional organisational hierarchy. Nevertheless, for understanding large dictionaries it is helpful to use this hierarchy as a way of collecting together categories that are thematically related, and we use the notion of 'theme' to structure our discussion of the core dictionary.

Core categories by 'theme': PUBLICATION DISPLAY AUDIT In practice, DISPLAY COLOUR AUDIT_AUTHOR much broader than the term JOURNAL AUDIT_AUTHOR_ROLE 'publication' JOURNAL COEDITOR AUDIT_CONFORM suggests: JOURNAL DATE AUDIT_CONTACT_AUTHOR

JOURNAL INDEX

PUBL AUTHOR

PUBL BODY

PUBL REQUESTED

PUBL SECTION

PUBL

JOURNAL TECHEDITOR

PUBL CONTACT AUTHOR

PUBL_MANUSCRIPT_INCL_EXTRA

PUBL MANUSCRIPT

essentially

metadata

overall

research

relating to the

project, and so

common to

any scientific

programme

IUCr

AUDIT_LINK

CITATION

AUDIT_SUPPORT

CITATION_AUTHOR

CITATION_EDITOR

DATABASE_CODE

DATABASE_RELATED

COMPUTING

DATABASE

The 'publication' group is, as indicated on the slide, really a broad collection of generic metadata that could apply to any information collection. It is possible that these categories might even be separated out into a distinct ('common' ?) dictionary.

Introduction to the CIF Ontology

Core categories by 'theme': EXPTL Experimental CHEMICAL work carried CHEMICAL_CONN_ATOM out prior to CHEMICAL_CONN_BOND diffraction CHEMICAL_FORMULA measurements EXPTL ABSORPT EXPTL CRYSTAL EXPTL_CRYSTAL_APPEARANCE EXPTL_CRYSTAL_FACE SPACE GROUP SPACE_GROUP_GENERATOR SPACE_GROUP_SYMOP SPACE_GROUP_WYCKOFF 14 Introduction to the CIF Ontology IUCr

Categories in this 'exptl' group are associated with experimental work carried out prior to diffraction measurements. Many of these are specific to crystallographic experiments (chemical characterisation, crystal preparation, space-group determination), and it is clear that any other experimental science would have a parallel requirement to record preparation and set-up information about an experiment.



These categories cover the data collected during a diffraction experiment. They also include much experimental metadata, and numerical data that could be, in the nomenclature favoured by the IUCr Committee on Data, 'raw' (collected straight from the instrument) or 'processed' (in the crystallographic context, typically structure factors).

Other structure- related information, including refinement strategy	Core categories by 'theme': STRUCTURE ATOM ATOM_ANALYTICAL ATOM_ANALYTICAL_MASS_LOSS ATOM_ANALYTICAL_SOURCE ATOM_SITE ATOM_SITE ATOM_SITES CATOM_SITES CATOM_SITES CARTN_TRANSFORM ATOM_SITES_FRACT_TRANSFORM ATOM_TYPE ATOM_TYPE SCAT REFINE REFINE_LS REFINE_LS_CLASS	
16	Introduction to the CIF Ontology	J.

These categories generally provide information about the structural model constructed after solution and refinement (in CommDat terminology, 'derived' data).



Broadly speaking, these cover the description of the chemical characteristics of the species investigated that are largely independent of its crystalline state.



This is a purely internal category that allows computational methods to be included for validating or evaluating relationships between data items.



We will present a series of diagrams showing how categories are developed amongst the various dictionaries that have extended the core ontology since 1991.



We begin with the macromolecular CIF dictionary, which rapidly grew into a largescale undertaking and has become the responsibility of the Worldwide Protein Data Bank for future maintenance and extension. We will not therefore discuss it in any detail here, but present a couple of slides to demonstrate the richness and complexity of relationships that it is capable of expressing. The STRUCT_CONN category records details about the connections between portions of a macromolecular structure. These can be hydrogen bonds, salt bridges, disulfide bridges and so on. The STRUCT_CONN_TYPE records define the criteria used to identify these connections. This example describes two disulfide bonds between portions of the folded protein backbone, as well as two hydrogen bonds (not illustrated).

mmCIF structure – altern	ative conformations
<pre>_atom_sites_alt_ens.id _atom_sites_alt_ens.details 'Ensemble 1' ; The inhibitor binds to the enzyme in two, roughly twofold symmetric, alternative conformations. This conformational ensemble includes the more populated conformation of the inhibitor (id=1) and the amino acid side chains that correlate with this inhibitor conformation. 'Ensemble 2' ; The inhibitor binds to the enzyme in two, roughly twofold symmetric, alternative conformations. This conformational ensemble includes the less populated conformation of the inhibitor (id=2) and the amino acid side chains that correlate with this inhibitor conformation. ; loop_ _atom_sites_alt_gen.ens_id _atom_sites_alt_gen.alt_id 'Ensemble 1' 1 'Ensemble 1' 1 'Ensemble 2' 2</pre>	Attenative conformations in an HIV-1 protease structure (Fitzgerald <i>et al.</i> , 2900) described with data items in the ATOM_SITES_ALT, ATOM_SITES_ENS and 2000).
21 Introduction to the CIF	Ontology

This example demonstrates how alternative conformations in the same structure can be handled. There are two conformations in which the inhibitor binds to the enzyme in a HIV-a protease structure. Reference: Fitzgerald, P.M., McKeever, B.M., VanMiddlesworth, J.F., Springer, J.P., Heimbach, J.C., Leu, C.T., Herber, W.K., Dixon, R.A., Darke, P.L. (1990). Crystallographic analysis of a complex between human immunodeficiency virus type 1 protease and acetyl-pepstatin at 2.0-Å resolution. *J. Biol. Chem.* **265**, 14209-14219.



The original mmCIF dictionary included all the existing core categories, and introduced many new ones. (Some, such as those describing bibliographic citations, were more detailed than existing treatments in the core, but were developed as completely distinct categories.) The mmCIF dictionary was thus a proper superset of the core.



With the introduction of DDLm methods, some new categories have been added to the core dictionary that are not present in mmCIF.

Description of restraints and constraints that can be applied in a structure refinement program	'Simple' extens RESTR RESTR_ANGLE RESTR_DISTANCE RESTR_DISTANCE_MIN RESTR_EQUAL_ANGLE RESTR_EQUAL_ANGLE RESTR_EQUAL_DISTANCE_CLASS RESTR_EQUAL_DISTANCE_CLASS RESTR_EQUAL_TORSION RESTR_EQUAL_TORSION_CLASS RESTR_PARAMETER RESTR_PARAMETER RESTR_PARAMETER_CLASS RESTR_PLANE RESTR_PLANE RESTR_PLANE RESTR_PLANE RESTR_RIGID_BODY RESTR_RIGID_BODY RESTR_TORSION RESTR_U_RIGID RESTR_U_RIGID RESTR_U_SIMILAR	<pre>sion dictionaries - ress</pre>	traints
24		Introduction to the CIF Ontology	

We now consider the individual extension dictionaries that have been developed with COMCIFS supervision. We will not go through them in great detail, but in each case we will try to illustrate the novel descriptive features that the dictionary is interested in developing. The restraints dictionary describes the types of geometric or ADP restraints that have been applied during a refinement strategy. Because of the dynamic way in which refinement software can arbitrarily restrain parameters, this cannot be done with complete generality. Nevertheless, a series of descriptions has been developed that provide insight into the refinement strategy. The rigid-bond restraint assumes that for atoms bound to one another the amplitude of motion along the direction of the bond is similar. The similar ADP restraint assumes that atoms close to one another would move similarly with respect to direction and amplitude. The isotropy restraint assumes that the atomic motion is approximately spherical. From Müller (2009).



The categories in the restraints dictionary all begin **RESTR_** and so are distinct from any in the core.

	Core	Restraints RESTR	
An aside Strictly speaking, from the ontological (<i>i.e.</i> conceptual) point of view, the extension dictionaries such as cif_restraints.dic are proper supersets of the core, as they import all its definitions	(mony others)		
26 IUCT	Introduction to the CIF Ontology		2023

A brief visual diversion to reinforce that these diagrams are sets of discrete category names, rather than an attempt to indicate the ontological completeness of an assembly of dictionaries. Ontologically, an extension dictionary such as cif_restraints.dic includes all the definitions in the core (and any other dictionaries it subsumes through **_import.get** statements or implicitly, by inheritance).



So, to reiterate, we draw the separate dictionaries as separate sets (unless they do share categories, as will be seen later).



And we'll include the mmCIF for good measure, as we build up our collection.



The twinning dictionary actually contains only these three categories (apart from the TWIN_GROUP parent) which are distinct from any categories in the core. Note that the TWIN_REFLN list may contain numeric values that also occur in the core REFLN list (*e.g.* measured and calculated squared structure-factors associated with diffraction peaks) but are partitioned differently, and must therefore be indexed with different category keys. The illustration shows two macroscopic twin forms of the mineral staurolite. The example listings are for a different structure.



The detailed electron density around individual atoms can be described by a multipole expansion method [Hansen, N. K. & Coppens, P. (1978). *Acta Cryst*. A**34**, 909–921]. In this formalism, local axes must be defined around each atom of interest, and the spherical harmonic coefficients according to the specific method of Hansen & Coppens may be enumerated. This is therefore a very specialised application. There is growing interest in returning to the topic of electron density to create a dictionary that can accommodate a wider range of modern quantum crystallographic formalisms.

	'Simple' extens	sion dictionaries – topology
Describes topological and connectivity properties of periodic nets	TOPOL TOPOL_ATOM TOPOL_ENTANGL TOPOL_LINK TOPOL_NET TOPOL_NODE TOPOL_TILING	<pre>_topol_net.id 1 _topol_net.z_number 2 _topol_net.overall_topology_RCSR 'dia' _topol_node.id 1 _topol_link.id 1 _topol_link.node_id_1 1 _topol_link.node_id_2 1 _topol_link.symop_id_2 13 _topol_link.type g1 loop</pre>
31		Introduction to the CIF Ontology

The topology dictionary is largely concerned with lattice descriptions rather than the physical structure of a crystal. However, it is clearly of interest to relate topological nets with an underlying distribution of atoms, and the TOPOL_ATOM category provides a link to a structure described in the core dictionary by requiring that values of **_topol_atom.atom_label** match corresponding values of **_atom_site.label** in a core category.



So all of these dictionaries contain their own discrete categories and can sit apart from the core in our diagrammatic representation.

	Extensior	n dictionaries – pov	wder	
	PD_AMORPHOUS	PD_MEAS_INFO_AUTHOR	20 20	10 ₂₀ (7) 40 50
	PD_BACKGROUND	PD_MEAS_OVERALL	_pd_meas_number_of_points	5500
Describes powder	PD_BLOCK	PD_PEAK	_pd_meas_2theta_range_min _pd_meas_2theta_range_max	5.000
data collection,	PD_BLOCK_DIFFRACTOGRAM	PD PEAK OVERALL	_pd_meas_2theta_range_inc	0.01
intensity profiles	PD_CALC_COMPONENT	PD PHASE	_diffrn_radiation_probe	X-ray
time of flight	PD_CALC_OVERALL	PD PHASE BLOCK	_pd_proc_ls_profile_function	'Pseudo-Voigt
time-oj-jiight,	PD_CALIB	PD PHASE MASS	_pd_proc_ls_background_function 'Manual background combined with 20 I	egendre polynoms
structure	PD_CALIBRATION	PD PREF ORIENT	_pd_proc_ls_pref_orient_corr	none 0.0586
refinement, phase	PD CALIB DETECTED INTENSITY	PD PREF ORIENT MARCH DOLLASE	_pd_proc_1s_prof_wR_factor	0.0785
mixturac atc	PD CALIB D TO TOF	PD PREF ORTENT SPHERICAL HARMONICS	_pd_proc_Is_prot_wk_expected	0.0580
mixtures, etc.	PD_CALTB_INCIDENT_INTENSITY	PD_PREP	_pd_proc_2theta_range_min pd_proc_2theta_range_max	5 59,99
			_pd_proc_2theta_range_inc	0.01
			loop_	
			_pd_proc_point_id _pd_proc_2theta_corrected	
		PD_PROC_OVERALL	_pd_proc_intensity_net _pd_calc_intensity_net	
		PD_QPA_CALIB_FACTOR	_pd_proc_ls_weight	
		PD_QPA_EXTERNAL_STD	2 5.0600 209.87 200.60	5 14.55 5 14.49
	PD_CALC	PD_QPA_INTENSITY_FACTOR	3 5.0700 212.65 200.3 4 5.0800 203.98 200.0	14.58 14.28
	PD_MEAS	PD_QPA_INTERNAL_STD	5 5.0900 185.85 199.80	13.63
	PD_PROC	PD_QPA_OVERALL	7 5.1100 192.20 199.3 7 5.1100 190.23 199.2	13.79
	PD_DIFFRACTOGRAM	PD_SPEC	Example Devider CIE and final Distantic data	f
	PD_INSTR	NREFLN	at room temperature.	i priuse ii oj griseoju
	PD_INSTR_DETECTOR			
33		Introduction to the CIF Ontology		

The powder dictionary is another large extension dictionary, and covers the rather diverse range of experimental types and structural characterisation met with in powder diffraction. Since a sample analysed by powder diffraction often contains distinct phase mixtures, powder CIFs are often organized with multiple data blocks, each data block describing one candidate phase. This dictionary introduced mechanisms for defining the relationship of each such data block to others present in the same file. There is an overlap with one core category (REFLN), as indicated by the pilcrow symbol, but I found it difficult to draw an overlap with the core diagram in a way that didn't interfere with the other overlaps I wanted to show.

Extension dictionaries – image ARRAY DATA **¶ DIFFRN MEASUREMENT** Description of raw ARRAY DATA EXTERNAL DATA DIFFRN MEASUREMENT AXIS image data ARRAY ELEMENT SIZE **¶ DIFFRN_RADIATION** collected in a ARRAY INTENSITIES **¶ DIFFRN_REFLN** diffraction ARRAY STRUCTURE DIFFRN SCAN experiment ARRAY STRUCTURE LIST DIFFRN SCAN AXIS ARRAY STRUCTURE LIST AXIS DIFFRN SCAN COLLECTION ARRAY STRUCTURE LIST SECTION DIFFRN SCAN FRAME AXIS DIFFRN SCAN FRAME AXIS DIFFRN DATA FRAME DIFFRN SCAN FRAME MONITOR Example: Reciprocal space of the **¶ DIFFRN DETECTOR** plastic phase of cyclohexane, at MAP DIFFRN DETECTOR AXIS 255 K, viewed in a projection normal MAP SEGMENT to [011]. DIFFRN DETECTOR ELEMENT VARIANT DIFFRN FRAME DATA 34 Introduction to the CIF Ontology IUCr

The core dictionary did allow for 'raw' data from point detectors, which were common at the time CIF was developed. However, area detectors were already beginning to become available, and a separate image dictionary was developed to handle diffraction images from a range of detector types. It was hoped that the image data could be standardised (using the ARRAY_DATA and related categories), but detector manufacturers proved resistant to conforming to such a standard, and imgCIF (and the associated binary CBF files) had relatively little uptake. However, latterly with growing community interest in archiving raw diffraction images, imgCIF has re-emerged as the natural vehicle at least for describing the experimental metadata, and it is used in IUCr journals to ensure reusability of data sets discussed in the *Raw Data Letters* section of the journal *IUCrData*. Increasingly high data capture rates make creation of CBF files on the fly increasingly difficult, but the ARRAY_DATA categories are still available for providing a standard image format if needed or desired for long-term archiving of selected data sets.



Because the image dictionary extended the description of raw data beyond the provisions already present in the core dictionary, it extended several core categories as well as introducing many new ones.

Extension dictionaries – modulated structures

Description of modulated and structures as superposition of distinct structural models or projection into 3space of higherdimensional symmetries

36

IUCr

ATOM SITES AXES ATOM SITES DISPLACE FOURIER ATOM SITES MODULATION ATOM SITES ORTHO *incommensurate* ATOM_SITES_ROT_FOURIER ATOM SITE DISPLACE FOURIER ATOM SITE DISPLACE FOURIER PARAM ATOM SITE DISPLACE LEGENDRE ATOM_SITE_DISPLACE_ORTHO ATOM SITE DISPLACE SPECIAL FUNC ATOM SITE DISPLACE XHARM ATOM SITE FOURIER WAVE VECTOR ATOM_SITE_OCC_FOURIER ATOM SITE OCC FOURIER PARAM ATOM SITE OCC LEGENDRE ATOM_SITE_OCC_ORTHO ATOM_SITE_OCC_SPECIAL_FUNC ATOM_SITE_OCC_XHARM ATOM SITE PHASON ATOM SITE ROT FOURIER ATOM_SITE_ROT_FOURIER_PARAM

ATOM_SITE_ROT_LEGENDRE ATOM_SITE_ROT_ORTHO ATOM_SITE_ROT_SPECIAL_FUNC ATOM_SITE_ROT_XHARM ATOM_SITE_U_FOURIER ATOM_SITE_U_FOURIER_PARAM ATOM_SITE_U_LEGENDRE ATOM_SITE_U_ORTHO ATOM_SITE_U_XHARM CELL_SUBSYSTEM CELL_SUBSYSTEMS CELL_WAVE_VECTOR CELL_WAVE_VECTORS ¶ DIFFRN_REFLN **¶ DIFFRN REFLNS** ¶ DIFFRN_STANDARD_REFLN ¶ GEOM_ANGLE etc. **¶ REFINE ¶ REFLN ¶ REFLNS**

¶ SPACE_GROUP_SYMOP



Another challenge was addressed by the modulated and composite structures dictionary (first available in 2002). The need was to extend the core description of crystal structures that had conventional space-group symmetries into the growing area of aperiodic structures. There are two main approaches: to describe distinct structural components that overlap in a non-periodic manner; and to use superspace groups to represent the aperiodicities as projections into 3-space of higherdimensional symmetries. Spatial modulations were described by a variety of harmonic modulation functions, and superspace group nomenclatures and symmetry operations were introduced.

The resultant dictionary overlaps with and extends several categories already existing in the core. And, because the characteristics of the diffraction spots are changed by the aperiodicity, there is also overlap with the DIFFRN_REFLN category which is now also common to imgCIF.

On top of the description of the atomic positions and motions in a crystal structure, there is now also growing interest in magnetic properties, and the magnetic CIF dictionary under development adds descriptions of magnetic moments and symmetries to the core description.

As magnetic properties are found in both periodic and aperiodic structures, so there is overlap with categories in both the core and the modulated structures dictionaries.

Summary: data definitions in CIF dictionaries

Managed by COMCIFS

- Crystallographic Core (coreCIF) 1991 and ongoing
- 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) 2021
- Crystallographic Symmetry (symCIF) 2001
- Diffraction Images (imgCIF) 2000

40

IUCr

- High pressure under development
- Crystallographic Macromolecular Structure (mmCIF) – 1997

Managed by 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

Introduction to the CIF Ontology

