

## 3. METHODOLOGY

entry can be located without any ambiguity and the best structure for the problem at hand can be used to start the Rietveld refinement.

### 3.7.5. Pearson's Crystal Data (PCD/LPF) (with Pierre Villars and Karen Cenzual)

#### 3.7.5.1. General information

The Pearson's Crystal Data database (PCD; Villars & Cenzual, 2013) is an outgrowth of the (Linus) Pauling File (LPF; Villars *et al.*, 1998; <http://www.paulingfile.com>), which was designed to combine crystal structures, phase diagrams and physical properties under the same computer framework to form a tool useful for materials design. PCD is the result of a collaboration between Material Phases Data Systems (Vitznau, Switzerland) and ASM International (Materials Park, Ohio, USA). The retrieval software was developed by Crystal Impact (Bonn, Germany). As suggested by the name, Pearson's Crystal Data is a follow-up product to *Pearson's Handbook: Crystallographic Data for Intermetallic Phases* (Villars & Calvert, 1985, 1991; Villars, 1997). However, in contrast to the latter, it also covers oxides and halides, which represent about 80% of the compounds with more than four chemical elements.

The 2016/2017 release of Pearson's Crystal Data contains more than 288 000 data sets for more than 165 300 different chemical formulae, representing over 53 000 distinct chemical systems. To achieve this, the editors have processed over 93 500 original publications; recent literature is surveyed in a cover-to-cover approach, including about 250 journal titles. Over 153 000 database entries contain refined atom coordinates, as well as isotropic and/or anisotropic displacement parameters when published, whereas more than 72 000 data sets contain atom coordinates corresponding to the structure prototype assigned by the authors of the original publication or by the database editors. Approximately 15 000 data sets contain only crystallographic data such as the lattice parameters and possibly a space group.

When available in the original publications, each data set contains comprehensive information on the sample-preparation and experimental procedure, as well as on the stability of the phase with respect to temperature, pressure and composition. The presence of plots (cell parameters or diffraction patterns) in the original paper is indicated, and over 30 000 descriptions of the variation of the cell parameters as a function of temperature, pressure or composition are proposed. Roughly 18 300 experimental diffraction patterns are reported.

The Linus Pauling File was designed as a phase-oriented, fully relational database system. This required the creation of a 'distinct phases' table, with internal links between the three parts of the database. In practice, this means that the senior editors have evaluated the distinct phases existing in the system for every chemical system using all information available in the LPF. Each structure entry in Pearson's Crystal Data has been linked to such a distinct phase, which allows a rapid overview of a particular chemical system.

#### 3.7.5.2. Evaluation procedure

Extensive efforts have been made to ensure the quality and reliability of the crystallographic data. Pearson's Crystal Data is checked for consistency by professional crystallographers, assisted by an original software package, *ESDD* (*Evaluation, Standardization and Derived Data*), containing more than 60

different modules (Cenzual *et al.*, 2000). The checking is carried out progressively, level by level. The following checks are made.

Individual database fields:

- (i) order of magnitude of numerical values;
- (ii) Hermann–Mauguin symbols, Pearson symbols;
- (iii) consistency of journal CODEN, year, volume, first page, last page;
- (iv) formatting of chemical formulae;
- (v) neutrality of oxides and halides;
- (vi) spelling.

Consistency within individual data sets:

- (i) atom coordinates, Wyckoff letters, site multiplicities;
- (ii) chemical elements in different database fields;
- (iii) computed, published values (cell volume, density, absorption coefficient, *d*-spacings);
- (iv) Pearson symbol, space group, cell parameters;
- (v) Bravais lattice, Miller indices;
- (vi) site symmetry, anisotropic displacement parameters.

Particular crystal-structure checks:

- (i) interatomic distances, sum of atomic radii;
- (ii) geometry of functional groups;
- (iii) search for overlooked symmetry elements;
- (iv) composition from refinement, chemical formula.

Consistency within the database:

- (i) comparison of cell-parameter ratios for isotypic entries;
- (ii) comparison of atom coordinates for isotypic entries with refined coordinates;
- (iii) comparison of densities;
- (iv) thorough search for duplicates, also considering translated references.

Wherever possible, misprints have been corrected based on arguments explained in remarks; as a result, more than 13 000 crystallographic data sets are accompanied by at least one erratum. In other cases remarks drawing the attention to discrepancies or unexpected features have been added.

The *ESDD* software package also produces derived data such as the Niggli reduced cell, equivalent isotropic displacement parameters, density and formula weight.

#### 3.7.5.3. Standardized crystallographic data

The crystallographic data in Pearson's Crystal Data are presented as published, respecting the original site labels, but are also standardized following the method proposed by Parthé and Gelato (Parthé & Gelato, 1984, 1985; Parthé *et al.*, 1993). This second presentation of the same data has been further adjusted so that compounds crystallizing with the same prototype structure (isotypic compounds) can be easily compared. It is prepared in a three-step procedure as follows.

- (i) The crystallographic data are checked for the presence of overlooked symmetry elements. Whenever it is possible to describe the structure in a higher-symmetry space group, or with a smaller unit cell, without any approximations, this is performed.
- (ii) In the next step, the crystallographic data are standardized using the program *STRUCTURE TIDY* (Gelato & Parthé, 1987).
- (iii) The resulting data are compared with the standardized data of the type-defining data set and, if relevant, adjusted using an *ESDD* module based on the program *COMPARE* (Berndt, 1994).

For data sets with no published coordinates, the cell parameters are standardized following the criteria defined for the unit-cell

### 3.7. CRYSTALLOGRAPHIC DATABASES

and space-group setting. For data sets with unknown space group, the cell parameters have been standardized assuming the space group of lowest symmetry in agreement with the Pearson symbol, e.g.  $P222$  for  $oP^*$  or  $o^{**}$ .

Standardized data are described with respect to the standard settings described in *International Tables for Crystallography* Volume A, with the following additional restraints: inversion centre at the origin, unique  $b$  axis and 'best' cell for monoclinic structures (Parthé & Gelato, 1985), triple-hexagonal cell for rhombohedral structures or Niggli reduced cell for triclinic structures. As a consequence, they can easily be incorporated into any program handling crystallographic data. The systematic standardization of the crystallographic data also greatly simplifies the classification of crystal structures into different prototypes.

A conversion tool to standardize cell parameters and/or compute the Niggli reduced cell is included in the software of Pearson's Crystal Data.

#### 3.7.5.4. Consequent prototype assignment

The prototype is a well known concept in inorganic chemistry, where a large number of compounds often crystallize with very similar atom arrangements. The compilation *Strukturbericht* started to catalogue crystal structures into types named by codes such as A1, B1 or A15. These notations are still in use; however, today prototypes are generally referred to by the name of the compound for which this particular kind of atom arrangement was first identified, i.e. Cu, NaCl and  $Cr_3Si$  for the types enumerated above. Pearson's Crystal Data uses a longer notation which also includes the Pearson symbol and the space-group number: Cu, $cF4,225$ , NaCl, $cF8,225$  and  $Cr_3Si,cP8,223$ . In a few cases several prototypes correspond to the same code, for example several polytypes of  $CdI_2$  have the same notation. A similar situation occurs for the wrong and the correct structure proposals for FeB, which have the same Pearson code and space group. In these cases a letter is added after the type-defining compound, for example the correct FeB type will be referred to as FeB-b, $oP8,62$ .

Each prototype is defined on a particular PCD database entry. In principle, this data set represents a recent refinement of the structure of the type-defining compound, but no effort has been made to find or use the most recent determination.

All of the data sets with published coordinates in Pearson's Crystal Data have been classified into prototypes following the criteria defined in *TYPIX* (Parthé *et al.*, 1993, 1994). According to this definition, isotypic compounds must crystallize in the same space group and have similar cell-parameter ratios; the atoms should occupy the same Wyckoff positions in the standardized description and have similar positional coordinates. If all of these criteria are fulfilled, the atomic environments should be similar. Note that  $H^+$  (protonic hydrogen) is ignored in the assignment of the prototype as well as in the Wyckoff sequence, Pearson symbol/code and atomic environments. Isopointal substitution variants are usually distinguished; however, no distinction is made between structures with fully and partly occupied atom sites. At present, 29 470 prototypes are represented.

When possible, a prototype has also been assigned to data sets without published atom coordinates. The prototype is often stated in the publication; in other cases the editors have assigned it. The editor will have added the exact space-group setting to which the cell parameters refer when this was not published. It is important to note that a prototype has been assigned at two different levels. The first is intimately related to the published

data (entry level), whereas the second is assigned at the phase level and may, in some cases, be inconsistent with the crystallographic data listed below.

For partly investigated structures, the available structural information is given using a similar way, for example the complete Pearson symbol may be replaced by  $t^{**}$  (tetragonal) or  $cI^*$  (cubic body-centred) and the place of the type-defining compound is occupied by an asterisk.

#### 3.7.5.5. Assigned atom coordinates

In order to give an approximate idea of the actual structure, a complete set of positional coordinates and site occupancies is proposed for data sets where a prototype could be assigned but the atom coordinates were not determined. The coordinates of the type-defining entry are proposed as a first approximation. The atom distribution is inserted by an *ESDD* module that compares the chemical formula of the type-defining entry with the chemical formula of the isotypic compound where the chemical elements have been reordered by the editor so that the first element is expected to occupy the same atom sites as the first element in the type-defining formula, and so on. Depending on the character of the prototype, substitutions and/or vacancies are either distributed over all atom sites occupied by the corresponding element or are expected to occur selectively on particular atom sites.

For this category of database entries, structure drawings, diffraction patterns and interatomic distances have also been computed. The structural portion of the database is thus more extensive than the primary literature.

#### 3.7.5.6. External links

When relevant, the database entries contain links to external data sources, including ASM International Alloys Phase Diagrams Centre Online, SpringerMaterials (The Landolt-Börnstein Database incorporating Inorganic Solid Phases PAULING FILE Multinaries Edition – 2010 in SpringerMaterials) and the original publication (through <https://www.crossref.org/>). A (static) reference to the Powder Diffraction File entry number is provided for database entries that are included in the PDF4+ product.

#### 3.7.5.7. Retrievable database fields

In addition to bibliographic (e.g. a particular institute) and chemical (e.g. sulfates) searches, many characteristics of the experiment and data processing (e.g. single crystal, neutron diffraction, range of temperature or reliability factors) or additional studies (e.g. pressure-dependence studies, magnetic structure) can be used as search criteria. Published crystal data, standardized crystal data and the Niggli reduced cells can be searched, as well as crystallographic classifications such as crystal class, Pearson symbol, Pearson code, Wyckoff sequence, structure prototype or structure class. Such searches can be very valuable in identifying a structural model for a new composition and saving the work of an *ab initio* structure determination.

The Quick Search pane includes commonly used searches on chemical elements (including cations in a particular oxidation state for oxides and halides), the number of elements and functional groups. The chemical selection (and/or/not) can be combined with selection on structure prototypes, space-group numbers and symbols or the crystal system. Retrieval on cell parameters (with ranges) and bibliographic information is also

### 3. METHODOLOGY

possible and the desired level of structural studies (*e.g.* complete) may be specified.

Many more searches can be carried out in the complete Search Dialog. Particularly useful are searches on atomic environments and interatomic distances. The atomic environment is defined by the coordination number, the geometry of the coordination polyhedron and the identities of the central and peripheral atoms. Searches on the number of different atomic environment types in the same structure can also be carried out. Specifying a pair of elements makes it possible to select a range of interatomic distances to be included in the search. These histograms are also useful in assessing the reasonableness of a particular distance in a Rietveld refinement.

#### 3.7.5.8. Particular software features

All searches use the 'Perpetual Restraining' feature, which updates the selection set in real time as a new query is introduced, so that the progress of the search scheme can easily be monitored. The complete Search Dialog offers a large variety of features that make the retrieval and presentation of information extremely flexible.

The Chemical System Matrix View makes it easy to locate phases in binary, ternary, quaternary and pseudo-quaternary systems. The Phases List View collects a selection set into its 'distinct' phases. From the individual database entry it is particularly easy to find all database entries with the same prototype structure and plot the unit-cell volume as a function of selected atomic radii. The standard display of an entry includes a short summary about the phase, structural, bibliographic, experimental and editorial data, as well as a structure drawing, a powder pattern and a table of interatomic distances.

The software for producing structure drawings offers the visualization of atomic environments (coordination polyhedra), the statistics of interatomic distances and the calculation of selected distances and angles. Four different models are available (ball and stick, wires, sticks and space-filling), with on-the-fly rotation controlled by the mouse. The nearest-neighbour histogram of a selected atom is compared with a statistical plot containing all distances in the database involving the same chemical elements, and the atomic environments can be instantaneously modified by clicking on the nearest-neighbour histogram.

Powder-diffraction patterns can be computed for any user-defined wavelength and the visualization includes a tool for zoom-in/out tracking. Patterns based on published lists of interplanar spacings can also be visualized. It is further possible to export database entries as CIF files, tables (*e.g.* powder-diffraction pattern, distances and angles) or graphics (*e.g.* structure drawings; BMP, GIF, JPG, PNG, TIFF or Diamond documents), and individually tailored dossiers can be designed and printed.

#### 3.7.6. Metals data file (CRYSTMET)

CRYSTMET (White *et al.*, 2002) began as a database of critically evaluated crystallographic data for metals, including alloys, intermetallics and minerals, and has grown to include inorganic compounds in general. It was started in 1960 by Cromer and Larson at Los Alamos National Laboratory, and its development was continued by the National Research Council of Canada. In 1996, the production and dissemination was transferred to Toth Information Systems.

CRYSTMET contains chemical, crystallographic and bibliographic data, together with comments regarding experimental details for each study. Using these data, a number of associated data files are generated, with the major one being a file of calculated powder patterns. Entry into CRYSTMET is *via* a number of search screens, including chemistry, bibliographic information, unit cell and reduced cell, powder patterns (using the positions of the strongest peaks as input), formula, structure type, Pearson symbol and space group. The results of queries reside in sets, which can be further manipulated using logical operations.

The results are displayed as a series of screens, which include crystallographic data, distances and angles, and the powder pattern. There is some ability to customize the calculation of the powder pattern of an entry; the calculation is performed for Debye–Scherrer geometry. Included on the Results tabs is a direct interface to the *MISSYM* program (Le Page, 1987, 1988), which searches the reported structure for additional symmetry elements. This is a very useful tool for detecting missed symmetry.

#### 3.7.7. Protein Data Bank (PDB)

The Protein Data Bank is described in Chapter 24.1 of *International Tables for Crystallography* Volume F (Berman *et al.*, 2011). Current information is available on the web at <https://www.wwpdb.org/>.

##### 3.7.7.1. Powder diffraction by proteins

Although powder-diffraction techniques had been applied to proteins as long ago as 1936 (Wyckoff & Corey, 1936; Corey & Wyckoff, 1936), and proof-of-principle experiments had been carried out (Rotella *et al.*, 1998, 2000), real progress in protein powder crystallography began with the work of Von Dreele (Von Dreele, 1998, 1999, 2003; Von Dreele *et al.*, 2000).

Progress in powder crystallography on macromolecules has been reviewed by Margiolaki & Wright (2008) and is also discussed in Chapter 7.1 of this volume. Notable studies include the characterization of the binding of *N*-acetylglucosamine oligosaccharides to hen egg-white lysozyme (Von Dreele, 2007*a*) and determination of the second SH3 domain of ponsin (Margiolaki *et al.*, 2007).

As with all powder diffraction, peak overlap ultimately limits the information available. Multi-pattern strategies to overcome the overlap problem have been investigated by Von Dreele (2007*b*). Multiple-pattern resonant-diffraction experiments have enabled study of the binding of  $\text{PtBr}_6^{2-}$  ions to lysozyme (Helliwell *et al.*, 2010). A bootstrap approach has been used to determine the structure of bacteriorhodopsin to 7 Å resolution (Dilanian *et al.*, 2011). Parametric resonant-scattering experiments have been used to determine the secondary structures of lysozyme derivatives (Basso *et al.*, 2010). Powder-diffraction experiments have also been used to gain insight into the general features of a nonstructural protein 3 (nsp3) macro domain (Papageorgiou *et al.*, 2010).

The structure of a five-residue peptide has been determined *ab initio* using laboratory powder data (Fujii *et al.*, 2011). We can expect further useful results at this interface between small-molecule and protein powder crystallography.

As is typical in other areas of science, powder diffraction has proven to be useful in more practical features of protein processing. It has been used to identify insulin (Norrman *et al.*, 2006) and GB1 (Frericks Schmidt *et al.*, 2007) polymorphs and