3.7. CRYSTALLOGRAPHIC DATABASES

lot-to-lot variations in lyophilized protein formulations (Hirakura *et al.*, 2007), and has been explored for use in structure-based generic assays (Allaire *et al.*, 2009).

### 3.7.7.2. *Calculation of protein powder patterns (with Kenny Ståhl)*

The Powder Diffraction File contains a few experimental powder patterns of proteins. These include silk fibroin protein (00-054-1394), tubulin (00-036-1547 and 00-036-1548), insulin (00-060-1360 through 00-060-1368), tomato bushy stunt virus (00-003-0001) and tobacco mosaic virus (00-003-0003 and 00-003-0004). Patterns have not yet been calculated from the structures in the Protein Data Bank because the calculated intensities generally fit poorly to those in experimental patterns.

Protein structures in the PDB do not generally contain H-atom positions, and the contributions from the disordered solvent in the solvent channels (which is the major source of the discrepancy) is not described (Hartmann *et al.*, 2010). The conventional Lorentz factor tends to infinity when approaching $2\theta = 0°$. Differences in data-collection temperatures and solvent content between powder and single-crystal specimens often mean that the lattice parameters differ. The relatively poor scattering from the protein and the large scattering from the mother liquor and sample holder result in significant background contributions to experimental powder patterns.

Optimization of the lattice parameters is generally straightforward and is important because most protein crystal structures are determined at low temperatures, while powder data are collected under ambient conditions. Protein crystals contain 30–80% disordered solvent. The solvent contribution to the diffraction pattern is most important for the low-angle powder data. In conventional protein crystallography several correction models have been developed (Moews & Kretsinger, 1975; Phillips, 1980; Jiang & Brünger, 1994), but the flat bulk-solvent model is the simplest one which yields a realistic correction (Jiang & Brünger, 1994; Hartmann *et al.*, 2010). This model includes two parameters: $k_{sol}$, which defines the level of electron density in the solvent region, and $B_{sol}$, which defines the steepness of the border between the solvent and macromolecular regions. These parameters are typically refined in contemporary software and cluster around $k_{sol} = 0.35$ e $Å^{-3}$ and $B_{sol} = 46$ $Å^2$ (Fokine & Urzhumtsev, 2002).

The flat bulk-solvent correction can be applied using *phenix.pdbtools* (Adams *et al.*, 2010), which requires a PDB coordinate file and values of $k_{sol}$ and $B_{sol}$ as input. Average values can be used, but refined values or values from the Electron Density Server (EDS; Kleywegt *et al.*, 2004) can improve the results. The bulk-solvent correction is highly anisotropic, and both parameters affect the anisotropy.

The ideal H-atom positions can be calculated using *phenix.pdbtools*. The solvent and hydrogen contributions to the pattern can be significant (Fig. 3.7.13).

The Lorentz factor $L$ describes the fraction of a reflection that is in the diffracting condition. For Bragg–Brentano and Debye–Scherrer geometries it is given by

$$L = \frac{1}{\sin 2\theta} \frac{1}{\sin \theta}. \qquad (3.7.3)$$

This equation assumes ideal crystals, resulting in infinitesimally small reciprocal-lattice points. The true size of the lattice points depends on the crystallite size and imperfections (strain). This smearing needs to be included in the Lorentz factor at low angles. A revised Lorentz factor for protein powder diffraction has been derived (Hartmann *et al.*, 2010),

$$L_{rev} = \frac{1}{\sin 2\theta} \frac{1}{\sin \theta} \frac{\sin^2 \theta}{(\sin^2 \theta + \lambda^2 \eta^2 / 12)}, \qquad (3.7.4)$$

in which $\eta$ reflects the distribution of scattering-vector amplitudes. For Guinier geometry these equations become more complex (Hartmann *et al.*, 2010). Fig. 3.7.14 shows that the Lorentz factor has a smaller effect than the solvent and H atoms, but that it is still significant. By applying these corrections it should be possible for the ICDD editorial staff to calculate useful powder patterns from PDB entries that could be included in the Powder Diffraction File.

Separating the background from the diffraction pattern is not straightforward (Frankaer *et al.*, 2011). Estimation of the background is greatly assisted by a correct calculated pattern. The calculated pattern can be scaled to the experimental data using *PROTPOW* (http://www.kemi.dtu.dk/english/Research/PhysicalChemistry/Protein_og_roentgenkrystallografi/Protpow).

Ståhl *et al.* (2013) have demonstrated that existing search/match procedures can be used to identify proteins using their powder patterns, and that powder patterns calculated from Protein Data Bank coordinates with proper care can be added to a database and included in the search/match procedure. Several problems can be foreseen when including large amounts of protein data into the Powder Diffraction File. It may be worthwhile including powder patterns with several levels of solvent correction, rather than just an average value. Asymmetry from instrumental effects and specimen transparency, which can affect the peak positions, needs to be taken into account. The use of an average thermal expansion coefficient may be sufficient to account for the differences in lattice parameters between low-temperature single-crystal structures and powder patterns measured under ambient conditions.
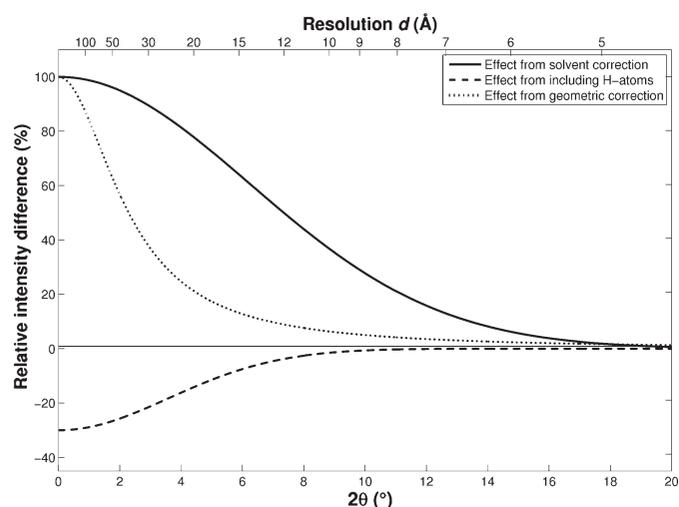


**Figure 3.7.13**
Overview of the trends from the different corrections. The effects are shown as the relative intensity difference $(I_{non-corr} - I_{corr})/I_{non-corr}$ plotted as functions of the scattering angle $2\theta$ (using Cu $K\alpha_1$) and resolution $d = \lambda/(2\sin\theta)$. The curves are based on average corrections of lysozyme and insulin data. $I_{non-corr}$ is the raw intensity from a calculated pattern which has only been Lorentz corrected. The geometric correction curve was calculated using $\eta = 0.045$ $Å^{-1}$. From Hartmann *et al.* (2010).

## 3.7.8. Crystallography Open Database (COD) (with Saulius Gražulis)

The Crystallography Open Database (COD) project (http://www.crystallography.net/cod/; Gražulis *et al.*, 2009, 2012) aims at collecting in a single open-access database all organic, inorganic and organometallic structures, except for the structures of biological macromolecules, which are available in the Protein Data Bank (Berman *et al.*, 2003, 2011). The database was founded by Armel Le Bail, Lachlan Cranswick, Michael Berndt, Luca Lutterotti and Robert M. Downs in February 2003 as a response to Michael Berndt's letter published on the Structure Determination by Powder Diffractometry (SDPD) mailing list (Berndt, 2003). Since December 2007, the main database server has been maintained and new software has been developed by Saulius Gražulis and Andrius Merkys at the Institute of Biotechnology of Vilnius University (VU). Currently, the database includes more than 376 000 entries describing structures of small molecules and small-to-medium-sized unit-cell materials as published in IUCr journals and other major crystallographic and peer-reviewed journals, as well as contributions by crystallographers from major laboratories. Most of the mineral data are obtained from the American Mineralogist Structure Database (Rajan *et al.*, 2006) and are donated by its maintainer and COD co-founder Robert M. Downs.
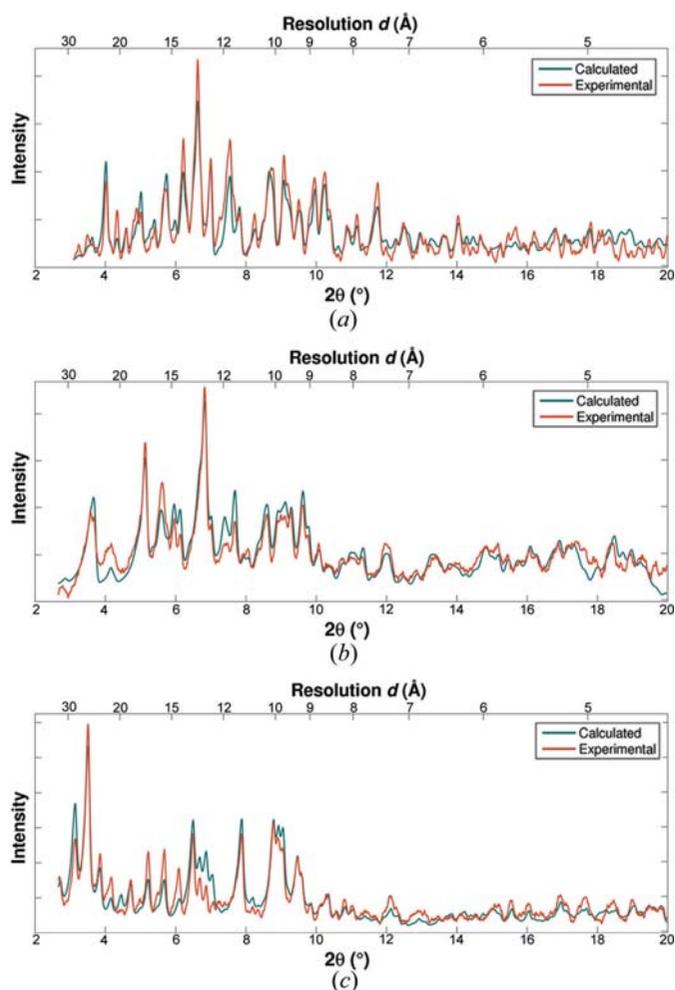
The database is an internet resource (Fig. 3.7.15) with data-search and download capabilities designed by Armel Le Bail and Michael Berndt. In addition, registered users may deposit new data, whether from previous publications or as personal communications, using the deposition web site designed at VU by



(*a*)



(*b*)

**Figure 3.7.15**
(*a*) The website and search interface of the Crystallography Open Database (COD) permits searches of crystallographic data by a range of parameters and unrestricted retrieval of the found data. (*b*) Data can be viewed online in the interactive *Jmol* applet (Hanson, 2010, 2013) or downloaded for further processing either one record at a time or in bulk.



**Figure 3.7.14**
Calculated and experimental powder patterns for (*a*) lysozyme, (*b*) trigonal insulin and (*c*) cubic insulin. The calculated patterns (blue) are corrected for bulk-solvent and geometrical effects using the revised Lorentz factor. From Hartmann *et al.* (2010).

Saulius Gražulis, Justas Butkus and Andrius Merkys. The deposition software performs rigorous checks of syntax and semantics.

The COD website allows searching on COD numerical identifier, unit-cell parameters, chemical composition and bibliographic data. Substructure searches using SMILES and SMARTS strings have been implemented. The free software package *OpenBabel* (O'Boyle *et al.*, 2011; Hutchison, 2007) is used for both the CIF-to-SMILES transformation and the actual search.

The retrieved records can be viewed online or downloaded for further processing. For massive data mining, COD permits downloads and updates of the whole database using Subversion, Rsync or http protocols. The ease of access to the COD data and its open nature has spurred the use of this resource for software testing (Grosse-Kunstleve & Gildea, 2011), teaching (Moeck, 2004) and research (First & Floudas, 2013). Multiple mirrors around the globe (Quirós-Olozábal, 2006; Gražulis, 2007; Moeck, 2007*a*; Chateigner, 2010) ensure data preservation, provide off-site backups, offer improved search interfaces (Moeck, 2007*b*) and increase reliability.

For the powder-diffraction community, the COD is interesting not only as an archive of structures solved by powder-diffraction methods, but also as a possibility for use in search/match procedures to identify crystalline compounds. Recently, the development of an open full-pattern search/match internet tool was launched by the COD developers. It allows phase quantifications from X-ray, neutron and electron powder patterns (with high- or medium-resolution instruments) provided that the structures are already in the COD. This tool is particularly suited to nanocrystalline powders, in which severe line broadening appears, precluding phase identification from only peak positions (Lutterotti *et al.*, 2012). COD-derived databases are also offered for software produced by several diffractometer vendors (Rigaku, 2011; PANalytical, 2012*a*,*b*; Bruker, 2013). In addition to the COD, searches and matches can be performed against its sister database, the PCOD, which contains structures predicted by the *GRINSP* program (Le Bail, 2005) and hypothetical zeolites (Pophale *et al.*, 2013). The power of such an approach is demonstrated by PCOD entry 3102887 (formulated as $SiO_2$). It was recently identified as corresponding structurally to a new phosphorus(V) oxonitride polymorph $\delta$-PON (Baumann *et al.*, 2012).

### 3.7.9. Other internet databases

Other useful databases include the following:

(i) The American Mineralogist Crystal Structure Database (http://rruff.geo.arizona.edu/AMSamcsd.php).

(ii) The Mineralogy Database (http://webmineral.com).

(iii) MinCryst (http://database.iem.ac.ru/mincryst/index.php).

(iv) The International Zeolite Association Database of Zeolite Structures (http://www.iza-structure.org/databases).

(v) The Incommensurate Structures Database (http://webbdcrista1.ehu.es/incstrdb/).

(vi) The Nucleic Acid Database (http://ndbserver.rutgers.edu).

**References**

Adams, P. D., Afonine, P. V., Bunkóczi, G., Chen, V. B., Davis, I. W., Echols, N., Headd, J. J., Hung, L.-W., Kapral, G. J., Grosse-Kunstleve, R. W., McCoy, A. J., Moriarty, N. W., Oeffner, R., Read, R. J., Richardson, D. C., Richardson, J. S., Terwilliger, T. C. & Zwart, P. H. (2010). *PHENIX: a comprehensive Python-based system for macromolecular structure solution. Acta Cryst.* D**66**, 213–221.

Allaire, M., Moiseeva, N., Botez, C. E., Engel, M. A. & Stephens, P. W. (2009). *On the possibility of using polycrystalline material in the development of structure-based generic assays. Acta Cryst.* D**65**, 379–382.

Allen, F. H., Cole, J. C. & Verdonk, M. L. (2011). *The relevance of the Cambridge Structural Database in protein crystallography. International Tables for Crystallography*, Vol. F, 2nd ed., edited by E. Arnold, D. M. Himmel & M. G. Rossmann, pp. 736–748. Chichester: Wiley.

Allmann, R. & Hinek, R. (2007). *The introduction of structure types into the Inorganic Crystal Structure Database ICSD. Acta Cryst.* A**63**, 412–417.

Barr, G., Dong, W. & Gilmore, C. J. (2009). *PolySNAP3: a computer program for analysing and visualizing high-throughput data from diffraction and spectroscopic sources. J. Appl. Cryst.* **42**, 965–974.

Barr, G., Gilmore, C. J. & Paisley, J. (2004). *SNAP-1D: a computer program for qualitative and quantitative powder diffraction pattern analysis using the full pattern profile. J. Appl. Cryst.* **37**, 665–668.

Basso, S., Besnard, C., Wright, J. P., Margiolaki, I., Fitch, A., Pattison, P. & Schiltz, M. (2010). *Features of the secondary structure of a protein molecule from powder diffraction data. Acta Cryst.* D**66**, 756–761.

Baumann, D., Sedlmaier, S. J. & Schnick, W. (2012). *An unprecedented $AB_2$ tetrahedra network structure type in a high-pressure phase of phosphorus oxonitride (PON). Angew. Chem. Int. Ed.* **51**, 4707–4709.

Behrens, H. & Luksch, P. (2006). *A bibliometric study in crystallography. Acta Cryst.* B**62**, 993–1001.

Belsky, A., Hellenbrandt, M., Karen, V. L. & Luksch, P. (2002). *New developments in the Inorganic Crystal Structure Database (ICSD): accessibility in support of materials research and design. Acta Cryst.* B**58**, 364–369.

Bergerhoff, G. & Brandenburg, K. (1999). *Typical interatomic distances: inorganic compounds. International Tables for Crystallography*, Vol. C, edited by E. Prince, pp. 770–781. Dordrecht: Kluwer Academic Publishers.

Bergerhoff, G. & Brown, I. D. (1987). *Crystallographic Databases*, edited by F. H. Allen, G. Bergerhoff & R. Sievers. Chester: International Union of Crystallography.

Berman, H. M., Henrick, K., Kleywegt, G., Nakamura, H. & Markley, J. (2011). *The Worldwide Protein Data Bank. International Tables for Crystallography*, Vol. F, 2nd ed., edited by E. Arnold, D. M. Himmel & M. G. Rossmann, pp. 827–832. Chichester: Wiley.

Berman, H., Henrick, K. & Nakamura, H. (2003). *Announcing the Worldwide Protein Data Bank. Nature Struct. Mol. Biol.* **10**, 980.

Berndt, M. (1994). Thesis, University of Bonn, Germany. Updates by O. Shcherban, SCC Structure-Properties Ltd, Lviv, Ukraine.

Berndt, M. (2003). *Open crystallographic database – a role for whom?* http://www.cristal.org/SDPD-list/2003/msg00025.html.

Bigelow, W. C. & Smith, J. V. (1964). *Two new indexes to the Powder Diffraction File. ASTM Spec. Tech. Publ.* STP372, 54. https://doi.org/10.1520/STP48334S.

Boldyrev, A. K., Mikheev, V. I., Dubinina, V. N. & Dovalev, G. A. (1938). *X-ray determination tables for minerals, Ft. I. Ann. Inst. Mines Leningrad*, **11**, 1–157.

Boles, M. O., Girven, R. J. & Gane, P. A. C. (1978). *The structure of amoxycillin trihydrate and a comparison with the structures of ampicillin. Acta Cryst.* B**34**, 461–466.

Bravais, A. (1866). *Etudes Cristallographiques*. Paris: Gathier Villars.

Bruker-AXS (2013). Crystallography Open Database for DIFFRAC.EVA. https://www.bruker.com/products/x-ray-diffraction-and-elemental-analysis/x-ray-diffraction/xrd-software/eva/cod.html.

Bruno, I. J., Cole, J. C., Edgington, P. R., Kessler, M., Macrae, C. F., McCabe, P., Pearson, J. & Taylor, R. (2002). *New software for searching the Cambridge Structural Database and visualizing crystal structures. Acta Cryst.* B**58**, 389–397.

Bruno, I. J., Cole, J. C., Kessler, M., Luo, J., Motherwell, W. D. S., Purkis, L. H., Smith, B. R., Taylor, R., Cooper, R. I., Harris, S. E. & Orpen, A. G. (2004). *J. Chem. Inf. Comput. Sci.* **44**, 2133–2144.

**references**