

21.1. VALIDATION OF PROTEIN CRYSTAL STRUCTURES

occupancies, unusual bond lengths or angles, unusual torsion angles or deviations from planarity (*e.g.* for the peptide plane), unusual chirality (*e.g.* for the C^α atom of every residue type except glycine), unusual differences in the temperature factors of chemically bonded atoms, unusual packing environments (Vriend & Sander, 1993), very short distances between non-bonded atoms (including symmetry mates), large positional shifts during refinement, unusual deviations from noncrystallographic symmetry (Kleywegt & Jones, 1995*b*; Kleywegt, 1996) *etc.*

21.1.3.3. Global statistics

The crystallographic *R* value used to be the major global quality indicator until it was realised that it can easily be fooled, especially at low resolution (Brändén & Jones, 1990; Jones *et al.*, 1991; Brünger, 1992*a*; Kleywegt & Jones, 1995*b*). The free *R* value, introduced by Brünger (1992*a*, 1993), has been shown to be much more reliable and harder to manipulate (Kleywegt & Brünger, 1996; Brünger, 1997). It is excellently suited for monitoring the progress of refinement, for detecting major problems with model or data and for helping reduce over-fitting of the data (which occurs if many more parameters are refined in a model than is warranted by the information content of the crystallographic data). Moreover, the free *R* value can be used to estimate the coordinate error of the final model (Kleywegt *et al.*, 1994; Kleywegt & Brünger, 1996; Brünger, 1997; Cruickshank, 1999).

In addition, the average or r.m.s. values for many of the local statistics, their minimum or maximum values or the percentage of outliers can be quoted and used to obtain an impression of the overall quality of the model and the overall fit of the model to the data.

21.1.4. Fixing errors

The object of model rebuilding is generally twofold: (1) to make the model as complete and detailed as the data will allow one to do confidently (*e.g.* to add previously unmodelled loops, ligands, water molecules *etc.*) and (2) to remove errors. At first glance, it may not seem all that important to fix each and every side chain and to correct all peptide O atoms that are pointing in the wrong direction, but one should keep in mind that an error in the scattering factor (atom type or charge), position or *B* factor of even a single atom will be detrimental to the entire model. Particularly in the early stages of model rebuilding and refinement, one often finds that after an extensive round of rebuilding followed by more refinement, the density improves dramatically and new features become clear. One should also keep in mind that incorrect features of a model may be very persistent and become 'self-fulfilling prophecies', a phenomenon known as 'model bias' (Ramachandran & Srinivasan, 1961; Read, 1986, 1994, 1997; Hodel *et al.*, 1992). This is particularly relevant in cases where unbiased phase information (*e.g.* SIRAS, MIR or MAD phases, or phases obtained after NCS or multiple-crystal averaging) is not available.

For error detection to be effective, it is best not to approach the rebuilding process in a haphazard way (Kleywegt & Jones, 1997). *O* users can employ a program called *OOPS* (Kleywegt & Jones, 1996*a*) to carry out this task in a systematic yet convenient fashion. This program uses information calculated by *O* (*e.g.* pep-flip and real-space fit values) and retrieves or derives other information from a PDB file of the current model (*e.g.* temperature factors, Ramachandran plot, changes with respect to the previous model). Moreover, results from a coordinate-based quality check by the *WHAT IF* program (Vriend, 1990; Hoofst *et al.*, 1996) can be included. In all, several dozen quality indicators can be used and plots and statistics for many of these can be produced by the program. The program's most useful feature, however, is that it will

generate *O* macros that when executed in *O* will take the crystallographer on a journey to all the residues that may require attention because they are outliers for one or more quality criteria. This makes the rebuilding process often faster and certainly more efficient and focused than a residue-by-residue walk through the model. In addition, it teaches inexperienced crystallographers to recognize and diagnose common model errors.

If a residue is an outlier for a certain criterion, the crystallographer has to inspect the local density and the structural context and decide the course of action. If the residue is in a region of the model in which many residues are outliers for many criteria, there may be something seriously wrong locally (for instance, there could be a register error), possibly because the density is poor. If there is poor density for several residues in a row, the crystallographer might consider leaving these residues out of the model for the next refinement round or cutting off the side chains at the C^β atoms. Sometimes local errors are correlated, such as a pep-flip error in one residue and a Ramachandran violation for its C-terminal neighbour, or a residue with a non-rotamer conformation and high temperature factors in conjunction with a poor real-space fit. *O* contains many tools to manipulate individual residues and atoms (Jones *et al.*, 1991; Jones & Kjeldgaard, 1994, 1997; Kleywegt & Jones, 1997), *e.g.* to flip a peptide plane, to replace a side chain by a rotamer conformation, to change side-chain torsion angles in order to optimize the fit to the density, to move groups of atoms, to use real-space refinement on a single residue or a zone of residues, to 'mutate' a residue to alanine *etc.* Together they constitute a toolbox with which many problems, once recognized, can be fixed relatively effortlessly (Kleywegt & Jones, 1997).

21.1.5. Preventing errors

As with everything else, when it comes to building a model of a protein, prevention of errors is the best medicine. Some general guidelines can be given (Dodson *et al.*, 1996; Kleywegt & Jones, 1997).

(1) Try to obtain the best possible set of data and the best possible set of phases for those data. If the structure has noncrystallographic symmetry (or if multiple crystal forms are available), use electron-density averaging to remove model bias and to reduce phase errors (Kleywegt & Read, 1997). In the absence of noncrystallographic symmetry, use maps that are biased by the model as little as possible [*e.g.* σ_A -weighted (Read, 1986) or omit maps (Bhat & Cohen, 1984; Bhat, 1988; Hodel *et al.*, 1992)]. If experimental phase information is available, keep and consult the experimental map(s). Experimental phases can also be used throughout the refinement process to alleviate or prevent some problems.

(2) Use databases to construct the initial model (or new parts of the model; Jones *et al.*, 1991; Kleywegt & Jones, 1998). All the crystallographer needs to do is to roughly place the C^α atoms in the density. The model-building program can then 'recycle' well refined high-resolution structures to place the main-chain atoms. Similarly, side-chain conformations should initially be chosen from the set of preferred rotamers for each residue type, perhaps in combination with a rigid-body rotation of the entire residue around its C^α atom and/or with minor adjustment of the torsion angles of long side chains (arginine, lysine *etc.*).

(3) After every cycle of refinement, carry out a critical analysis of the quality of the current model. This entails the calculation of properties such as those discussed in Section 21.1.3 and the inspection of the residues that are outliers for any of them, as described in Section 21.1.4. Be conservative during rebuilding, especially when the model is incomplete and possibly full of errors.

(4) Design a refinement protocol that is appropriate for the available data. If NCS restraints do not give a significantly better

21. STRUCTURE VALIDATION

free R value than NCS constraints, then use constraints. If NCS restraints are to be employed, then use the experimental map to design a suitable NCS-restraint scheme (Kleywegt, 1999). Avoid the temptation to model alternative conformations in low-resolution maps or to place putative solvent molecules in every local maximum of the $(F_o - F_c, \alpha_c)$ difference map. In other words, be conservative and remember that the maxim 'where freedom is given, liberties are taken' is highly applicable to refinement programs (Hendrickson & Konnert, 1980; Kleywegt & Jones, 1995b).

(5) Adopt methodological advances as soon as they become available. Several innovations have only been slowly accepted by the mainstream (*e.g.* the use of databases in building and rebuilding, the use of the free R value, the use of electron-density averaging in molecular-replacement cases, bulk-solvent modelling). The most prominent recent development is the use of likelihood-based refinement programs (Bricogne & Irwin, 1996; Pannu & Read, 1996; Murshudov *et al.*, 1997; Adams *et al.*, 1997; Pannu *et al.*, 1998). These programs produce better models and maps and considerably reduce over-fitting (as assessed by the difference between the free and conventional R values).

(6) Most importantly, the crystallographer should be hyper-critical towards the fruits of his or her own labour. Every intermediate model is a hypothesis to be shot down (Jones & Kjeldgaard, 1994). The crystallographer should be more critical than the supervisor, the supervisor more critical than the referee and the referee more critical than the casual reader. It goes without saying that the reader, casual or not, should have access to model coordinates, experimental data and electron-density maps!

21.1.6. Final model

Once the refinement is finished [*i.e.* once the $(F_o - F_c, \alpha_c)$ difference map is featureless (Cruickshank, 1950) and parameter shifts in further refinement cycles are negligibly small], three tasks remain: validation of the final model, description and analysis of the structure, and deposition of the model coordinates and the crystallographic data with the Protein Data Bank (Bernstein *et al.*, 1977).

Until a few years ago, validation of the final model typically entailed calculating the conventional R value, r.m.s. deviations from ideal values of bond lengths and angles, average temperature factors, and a Luzzati-type estimate of coordinate error. Kleywegt & Jones (1995b) showed that these statistics are not necessarily even remotely related to the actual quality of a model. Based on these criteria, a backwards-traced protein model was of higher apparent quality than a carefully refined correct model. After this, the realisation sunk in that the best validation criteria are those that assess aspects of the model that are 'orthogonal' to the information used during model refinement and rebuilding. For instance, the main-chain φ and ψ torsion angles are usually not restrained during refinement; this makes the Ramachandran plot such a powerful validation tool (Kleywegt & Jones, 1996b, 1998). Other examples of useful independent tests include the profile method of Eisenberg and co-workers (Lüthy *et al.*, 1992), the directional atomic contact analysis method of Vriend & Sander (1993) and the threading-potential method of Sippl (1993).

In general, all quality checks provide necessary, but in themselves insufficient, indications as to whether or not a model is essentially correct. A truly good model should make sense with respect to what is currently known about physics, chemistry, crystallography, protein structures, statistics and (last, but not least) biology and biochemistry (Kleywegt & Jones, 1995a). A good model will typically score well on most if not all validation criteria, whereas a poor one will score poorly on many criteria. The same is

true at the level of residues: a poor or erroneous region in a model will be characterized by violations of many residue-level quality criteria (Kleywegt & Jones, 1997).

21.1.7. A compendium of quality criteria

In this section, some of the quality and validation criteria that have been used by macromolecular crystallographers are summarized (for more detailed information, the reader is referred to the primary literature). When judging how useful or powerful these criteria are in a certain case, one should keep in mind that any criterion that has been used explicitly or implicitly during model refinement (*e.g.* geometric restraints) or rebuilding (*e.g.* rotamer libraries) does *not* provide a truly independent check on the quality of the model.

Many, but not all, of the criteria discussed below pertain specifically to protein models. Comparatively little work has been performed on the validation of nucleic acid models, although there are indications that there is a need for such procedures (Schultze & Feigon, 1997). The situation would appear to be even worse for hetero-entities (*e.g.* ligands, inhibitors, cofactors, covalent attachments, saccharides, metals, ions; van Aalten *et al.*, 1996; Kleywegt & Jones, 1998).

21.1.7.1. Data quality

Although many quality and validation criteria have been developed for assessing coordinate sets of protein models, comparatively few criteria are available for assessing the quality of the crystallographic data.

21.1.7.1.1. Merging R values

Possibly the most common mistake in papers describing protein crystal structures is an incorrectly quoted formula for the merging R value (calculated during data reduction),

$$R_{\text{merge}} = \frac{\sum_h \sum_i |I_{h,i} - \langle I_h \rangle|}{\sum_h \sum_i I_{h,i}},$$

where the outer sum (h) is over the unique reflections (in most implementations, only those reflections that have been measured more than once are included in the summations) and the inner sum (i) is over the set of independent observations of each unique reflection (Drenth, 1994). This statistic is supposed to reflect the spread of multiple observations of the intensity of the unique reflections (where the multiple observations may derive from symmetry-related reflections, different images or different crystals). Unfortunately, R_{merge} is a very poor statistic, since its value increases with increasing redundancy (Weiss & Hilgenfeld, 1997; Diederichs & Karplus, 1997), even though the signal-to-noise ratio of the average intensities will be higher as more observations are included (in theory, an N -fold increase of the number of independent observations should improve the signal-to-noise ratio by a factor of $N^{1/2}$). At high redundancy, the value of R_{merge} is directly related to the average signal-to-noise ratio (Weiss & Hilgenfeld, 1997): $R_{\text{merge}} \simeq 0.8/\langle I/\sigma(I) \rangle$.

Diederichs & Karplus (1997) have suggested a number of alternative measures that lack most of the drawbacks of R_{merge} . Their statistic R_{meas} is similar to R_{merge} , but includes a correction for redundancy (m),

$$R_{\text{meas}} = \frac{\sum_h [m/(m-1)]^{1/2} \sum_i |I_{h,i} - \langle I_h \rangle|}{\sum_h \sum_i I_{h,i}}.$$

Another statistic, the pooled coefficient of variation (PCV), is defined as

$$\text{PCV} = \frac{\sum_h \{[1/(m-1)] \sum_i (I_{h,i} - \langle I_h \rangle)^2\}^{1/2}}{\sum_h \langle I_h \rangle}.$$