

## 2.2. QUALITY INDICATORS

indicators here. For a thorough treatment of translation-function applications, the reader is referred to Chapter 2.3 of *International Tables for Crystallography* Volume B and Chapter 13.3 of the present volume.

**Log-likelihood gain, LLG.** In likelihood-based molecular replacement (McCoy *et al.*, 2007), potential molecular-replacement solutions are evaluated using likelihood, which is defined as the probability  $P$  that the observed diffraction data would have been measured if the orientation (and, usually, position) of the model were correct (Read, 2001). The score is reported in terms of the log-likelihood gain (LLG), which is defined as the logarithm of the likelihood score for the model  $p(F_{\text{obs}}; \text{model})$  minus the logarithm of the likelihood score for a random-atom Wilson distribution  $p_{\text{Wilson}}(F_{\text{obs}})$ . The LLG measures how much better the data can be predicted from the molecular-replacement model than from a collection of random atoms.

$$\text{LLG} = \sum_{hkl} \ln[p(F_{\text{obs}}; \text{model})] - \sum_{hkl} \ln[p_{\text{Wilson}}(F_{\text{obs}})]. \quad (2.2.7.2)$$

**LLG-Z score.** It is important to note that the LLG depends on the quality of the model and the number of reflections, so the absolute values cannot be compared between different molecular-replacement applications. Instead, the quality of a molecular-replacement solution can be judged by the LLG-Z score, which is defined as the number of standard deviations a score is above the mean score in a particular rotation or translation search.

$$\text{LLG-Z} = \text{LLG} - \langle \text{LLG} \rangle / [(\text{LLG} - \langle \text{LLG} \rangle)^2]^{1/2}. \quad (2.2.7.3)$$

The translation function  $Z$  score (TFZ) for the last component placed in a molecular-replacement search is often a good indicator of the confidence that can be placed in the solution. If TFZ is greater than 8 and there is no translational pseudo-symmetry, the solution is almost always correct.

Detailed descriptions of the background and proper application of molecular-replacement approaches are presented in Chapter 2.3 of *International Tables for Crystallography* Volume B and Chapters 13.2 and 13.3 of the present volume.

## 2.2.8. Quality indicators for refinement

The last step of a structure determination is the refinement of the model against the observed data. Refinement is in principle a mathematical operation that is applied in order to minimize the discrepancy between the observed structure-factor amplitudes  $|F_{\text{obs}}|$  and the calculated ones  $|F_{\text{calc}}|$ .

**Crystallographic  $R$  factor,  $R$ .** The crystallographic  $R$  factor  $R$  is defined as the fractional disagreement between the set of observed structure-factor amplitudes and amplitudes calculated from the structural model. Of course, observed and calculated reflection sets need to be on the same scale.

$$R = \sum_{hkl} \frac{||F_{\text{obs}}| - |F_{\text{calc}}||}{\sum_{hkl} |F_{\text{obs}}|}. \quad (2.2.8.1)$$

**Free  $R$  factor,  $R_{\text{free}}$ .** The free  $R$  factor  $R_{\text{free}}$  is defined in the same way as the crystallographic  $R$  factor, but it is based on a set of reflections that have been excluded from the refinement (Brünger, 1992). The excluded set of reflections is called the *test set*, while the set of reflections used for refinement is called the *working set*. The test set can be chosen randomly or systematically, either in thin resolution shells or to account for the presence of noncrystallographic symmetry, respectively. In order to minimize the impact on the final model, the test set should be as small as possible. Typically, it contains about 5–10% of the

reflections, or at least enough reflections to keep the standard deviation of  $R_{\text{free}}$  below 1%, but there is no need to use more than 2000 reflections (Kleywegt & Brünger, 1996; Brünger, 1997). The standard deviation of  $R_{\text{free}}$  has been empirically estimated to be  $R_{\text{free}}/N^{1/2}$ , where  $N$  is the number of reflections in the test set (Brünger, 1997). Of course, there may be concerns about the impact of excluding 5–10% of reflections on the final model, but a few final cycles of refinement against the recombined full data set should allay them.

**Correlation coefficients  $\text{CC}(F_{\text{obs}}, F_{\text{calc}})$  and  $\text{CC}(I_{\text{obs}}, I_{\text{calc}})$ .**  $\text{CC}(F_{\text{obs}}, F_{\text{calc}})$  and  $\text{CC}(I_{\text{obs}}, I_{\text{calc}})$  are Pearson linear correlation coefficients [see equation (2.2.2.13)] between observed and model-based calculated structure-factor amplitudes or intensities, respectively, that find use from time to time. One advantage of the use of a correlation coefficient instead of an  $R$  factor is that it avoids the problem of scaling the two sets of numbers relative to each other.

## 2.2.9. Quality indicators for the refined model

In MX, the observable-to-parameter ratio is mostly unfavourable. Therefore, structure refinements are carried out with boundary conditions, constraints and restraints. Constraints reduce the number of parameters which need to be refined, while restraints provide additional information to the refinement procedure that increases the number of observables. A refined model, therefore, has to fulfil not only the criterion that the crystallographic  $R$  factor [see equation (2.2.8.1)] is good and that the model fits well to the electron density, but also that it fits well to the restraints used in the refinement procedure.

**Real-space residual, RSR.** The real-space residual, RSR (Jones *et al.*, 1991), quantifies the discrepancies between the electron-density maps  $\rho_1$ , calculated directly from a structural model, and  $\rho_2$ , calculated from experimental data. RSR can take the form of a real-space  $R$  factor RSRF and of a real-space correlation coefficient RSCC.

$$\text{RSRF} = 2 \sum_{xyz} |\rho_1 - \rho_2| / \sum_{xyz} (\rho_1 + \rho_2). \quad (2.2.9.1)$$

The sum  $\sum_{xyz}$  runs over all grid points of the electron-density maps that are close to the model. A big advantage of RSRF is that it can be calculated on a residue-by-residue basis. It therefore gives a local picture of structure quality. It can also be used throughout model building and refinement in order to follow the improvement of the model locally on a per-residue basis.

RSCC is defined as the Pearson linear correlation coefficient [see equation (2.2.2.13)] between  $\rho_1$  and  $\rho_2$ . Everything said about RSRF above applies to RSCC as well.

**R.m.s. deviation from ideal of geometric parameter  $x$ .** The root-mean-square deviation of a set of geometric parameters  $x$  from their ideal values is defined as

$$\text{r.m.s.d}(x) = \left\{ \sum_i [x_i(\text{ideal}) - x_i(\text{observed})]^2 / N \right\}^{1/2}. \quad (2.2.9.2)$$

The sum runs over all  $N$  instances of the geometric parameter occurring in a structure. The geometric parameters  $x$  that are typically considered are bond lengths, bond angles, dihedral angles, chiral volumes, planar groups *etc.* The ideal values for proteins are typically taken from the study of Engh & Huber (1991) and for nucleic acids from Parkinson *et al.* (1996).

**Z score.** A measure of the likelihood that an individual geometric parameter is correct is given by its  $Z$  score. The  $Z$  score is defined as the distance of an individual data point of a distribution from the mean of the distribution expressed in standard

## 2. BASIC CRYSTALLOGRAPHY

**Table 2.2.11.1**

Definitions of the most commonly used quality indicators

Indicator	Details
<p>Optical resolution,</p> $d_{\text{opt}} = [2(\sigma_{\text{Patt}}^2 + \sigma_{\text{sph}}^2)]^{1/2} \quad (2.2.2.2)$	<p><math>\sigma_{\text{Patt}}</math> and <math>\sigma_{\text{sph}}</math> are the standard uncertainties of Gaussians fitted to the origin peak of the Patterson function of the diffraction data set and the origin peak of the spherical interference function, respectively</p>
<p><math>R_{\text{merge}}</math> (merging <math>R</math> factor),</p> $R_{\text{merge}} = \frac{\sum_{hkl} \sum_i  I_i(hkl) - \langle I(hkl) \rangle }{\sum_{hkl} \sum_i I_i(hkl)} \quad (2.2.2.4)$	<p><math>\langle I(hkl) \rangle</math> is the mean of the several individual measurements <math>I_i(hkl)</math> of the intensity of reflection <math>hkl</math></p>
<p><math>R_{\text{meas}}</math> or <math>R_{\text{r.i.m}}</math> (redundancy-independent merging <math>R</math> factor),</p> $R_{\text{r.i.m.}} = \frac{\sum_{hkl} \{N(hkl)/[N(hkl) - 1]\}^{1/2} \sum_i  I_i(hkl) - \langle I(hkl) \rangle }{\sum_{hkl} \sum_i I_i(hkl)} \quad (2.2.2.5)$	<p><math>\langle I(hkl) \rangle</math> is the mean of the <math>N(hkl)</math> individual measurements <math>I_i(hkl)</math> of the intensity of reflection <math>hkl</math></p>
<p><math>R_{\text{p.i.m.}}</math> (precision-indicating merging <math>R</math> factor),</p> $R_{\text{p.i.m.}} = \frac{\sum_{hkl} \{1/[N(hkl) - 1]\}^{1/2} \sum_i  I_i(hkl) - \langle I(hkl) \rangle }{\sum_{hkl} \sum_i I_i(hkl)} \quad (2.2.2.6)$	<p><math>\langle I(hkl) \rangle</math> is the mean of the <math>N(hkl)</math> individual measurements <math>I_i(hkl)</math> of the intensity of reflection <math>hkl</math></p>
<p><math>R_{\text{anom}}</math> (anomalous <math>R</math> factor),</p> $R_{\text{anom}} = \frac{\sum_{hkl}  I(hkl) - I(\bar{h}\bar{k}\bar{l}) }{\sum_{hkl} I(hkl)} \quad (2.2.2.12)$	<p><math>\langle I(hkl) \rangle</math> is the mean intensity of the Friedel mates of the reflection <math>hkl</math>, or <math>\frac{1}{2}[I(hkl) + I(\bar{h}\bar{k}\bar{l})]</math></p>
<p><math>R_{\text{Cullis}}</math> (Cullis <math>R</math> factor for isomorphous-replacement applications),</p> $R_{\text{Cullis}} = \frac{\sum_{hkl}  F_{\text{PH}} -  F_{\text{P}} + F_{\text{H}}  }{\sum_{hkl}  F_{\text{PH}} - F_{\text{P}} } \quad (2.2.5.1b)$	<p><math>F_{\text{P}}</math>, <math>F_{\text{PH}}</math> and <math>F_{\text{H}}</math> are the structure factors for the protein, the heavy-atom derivative and the heavy atoms alone, respectively</p>
<p><math>\text{PP}_{\text{iso}}</math> (phasing power for isomorphous-replacement applications),</p> $\text{PP}_{\text{iso}} = \frac{\sum_{hkl}  F_{\text{H}} }{\sum_{hkl}  F_{\text{PH}} -  F_{\text{P}} + F_{\text{H}}  } \quad (2.2.5.3)$ <p>or, as in the program <i>SOLVE</i>,</p> $\text{PP}_{\text{iso}} = \frac{(\sum_{hkl}  F_{\text{H}} ^2)^{1/2}}{(\sum_{hkl} ( F_{\text{PH}} -  F_{\text{P}} + F_{\text{H}}  )^2)^{1/2}} \quad (2.2.5.4)$	<p><math>F_{\text{P}}</math>, <math>F_{\text{PH}}</math> and <math>F_{\text{H}}</math> are the structure factors for the protein, the heavy-atom derivative and the heavy atoms alone, respectively</p>
<p><math>R</math> (crystallographic <math>R</math> factor) and <math>R_{\text{free}}</math> (free <math>R</math> factor),</p> $R = \frac{\sum_{hkl}   F_{\text{obs}}  -  F_{\text{calc}}  }{\sum_{hkl}  F_{\text{obs}} } \quad (2.2.8.1)$	<p><math>R_{\text{free}}</math> is defined as the crystallographic <math>R</math> factor but for a subset of reflections that have been excluded from refinement</p>
<p>RSRF (real-space <math>R</math> factor),</p> $\text{RSRF} = 2 \frac{\sum_{xyz}  \rho_1 - \rho_2 }{\sum_{xyz} (\rho_1 + \rho_2)} \quad (2.2.9.1)$	<p><math>\rho_1</math> and <math>\rho_2</math> are the electron-density maps calculated from the structural model and from the experimental data, respectively</p>
<p>R.m.s.d.'s of geometric parameters <math>x</math>,</p> $\text{r.m.s.d.}(x) = \left\{ \frac{\sum_i [x_i(\text{ideal}) - x_i(\text{observed})]^2}{N} \right\}^{1/2} \quad (2.2.9.2)$	<p><math>x_i</math> are the individual values, ideal or observed, of the geometric parameter <math>x</math> and the sum is over all <math>N</math> <math>x_i</math> observed. The geometric parameters <math>x</math> may be bond lengths, bond angles, dihedral angles, chiral volumes, deviations from planarity <i>etc.</i></p>
<p>DPI (diffraction-component precision index),</p> $\text{DPI} = [N_{\text{atom}}/(N_{\text{hkl}} - N_{\text{para}})]^{1/2} R d_{\text{min}} C^{-1/3} \quad (2.2.10.1)$	<p><math>N_{\text{atom}}</math> is the number of atoms in the structure, <math>N_{\text{hkl}}</math> is the number of reflections, <math>N_{\text{para}}</math> is the number of refined parameters, <math>R</math> is the crystallographic <math>R</math> factor, <math>d_{\text{min}}</math> is the nominal resolution and <math>C</math> is the fractional completeness of the data set</p>

## 2.2. QUALITY INDICATORS

deviations. In the case described here, the mean values of the distribution are the ideal values taken from Engh & Huber (1991) and Parkinson *et al.* (1996).

$$Z(x_i) = [x_i(\text{observed}) - x_i(\text{ideal})]/\sigma(x). \quad (2.2.9.3)$$

Ideally, the  $Z$  score should be 0. A parameter that exhibits a  $Z$  score of less than  $-4$  or greater than  $+4$  is highly unlikely and calls for attention.

**Root-mean-square  $Z$  score, r.m.s.- $Z$ .** Although r.m.s.d. values [see equation (2.2.9.2)] are still popular for use in judging the quality of refined macromolecular models, a much more useful statistic is the r.m.s. value of a distribution of  $Z$  scores or the r.m.s.- $Z$  score.

$$\text{r.m.s.-}Z(x) = \sum_i [Z(x_i)^2/N]^{1/2}. \quad (2.2.9.4)$$

The sum runs over all  $N$  instances of the geometric parameter  $x$  occurring in a structure. A very useful property of  $Z$  scores is that the r.m.s. values of  $Z$ -score distributions should always be 1. Significant deviations from the ideal value indicate potential problems. R.m.s.  $Z$  scores are widely used, for instance, in the program *WHAT\_CHECK* (Hooft *et al.*, 1996).

**R.m.s.d. (NCS).** The root-mean-square deviation from crystallographic symmetry between two molecules related by non-crystallographic symmetry (NCS) can be calculated from a superposition of the two molecules. It is defined as

$$\text{r.m.s.d. (NCS)} = \left( \sum_i d_i^2/N \right)^{1/2}. \quad (2.2.9.5)$$

The sum runs over  $N$  equivalent atom pairs with  $d_i$  being the distance between the two equivalent atoms after superposition.

### 2.2.10. Error estimation for the refined model

An important quality indicator for a refined model is the coordinate uncertainty. Short of full-matrix inversion, which is the standard procedure in small-molecule crystallography but which is applicable only in exceptional cases for macromolecules, some methods have been devised for estimating of the overall coordinate uncertainty.

**Error estimation according to Luzzati.** For most macromolecular structure determinations, atomic standard uncertainties are not available. However, Luzzati (1952) devised a method of estimating the average positional error of a structure. Under the assumption that the atomic positional errors follow a normal distribution, the average error can be estimated by comparing a plot of the crystallographic  $R$  factor [see equation (2.2.8.1)] versus the reciprocal resolution (or  $2 \sin \theta/\lambda$ ) with pre-computed theoretical curves for different average errors. A more recent – and probably better – approach is to use the free  $R$  factor instead of the crystallographic  $R$  factor.

**SigmaA- ( $\sigma_A$ )-type error estimation.** A slightly better estimate of the average positional error of a structure can be obtained by plotting the natural logarithm of the parameter  $\sigma_A$  versus  $(\sin \theta/\lambda)^2$  (Read, 1986). The slope of a straight line fitted to the plot provides an estimate of the average positional error of the structure. The parameter  $\sigma_A$  assumes normally distributed positional errors and takes model incompleteness into account as well.

**Diffraction-component precision index, DPI.** The diffraction-component precision index DPI is an empirical parameter describing the overall coordinate uncertainty of a structure (Cruickshank, 1999a,b). For an atom with an isotropic displacement parameter of average value ( $B_{\text{avg}}$ ), it is defined as

$$\text{DPI} = [N_{\text{atom}}/(N_{\text{hkl}} - N_{\text{para}})]^{1/2} R d_{\text{min}} C^{-1/3}, \quad (2.2.10.1)$$

where  $N_{\text{atom}}$  is the number of atoms included in the refinement,  $N_{\text{hkl}}$  is the number of reflections included in the refinement,  $N_{\text{para}}$  is the number of refined parameters,  $R$  is the crystallographic  $R$  factor,  $d_{\text{min}}$  is the nominal resolution of the data included in the refinement and  $C$  is the data completeness. The free  $R$  factor  $R_{\text{free}}$  is sometimes used instead of the crystallographic  $R$  factor  $R$  to calculate the DPI. In this case  $(N_{\text{hkl}} - N_{\text{para}})$  is replaced by  $N_{\text{free}}$ , which is the number of reflections used for  $R_{\text{free}}$  calculation.

The Cruickshank formula for DPI has been recast into several other forms, including

$$\text{DPI} = \sigma(x, B_{\text{avg}}) = 0.18(1 + sc)^{1/2} V_M^{-1/2} R_{\text{free}} d_{\text{min}}^{5/2} C^{-5/6} \quad (2.2.10.2)$$

(Blow, 2002), where  $sc$  is the solvent fraction ( $= N_{\text{solv}}/N_{\text{atom}}$ , where  $N_{\text{solv}}$  is the number of atoms that are solvent) and  $V_M$  is the Matthews parameter (Matthews, 1968). The utility of this latter formula in guiding the design of the data-collection experiment to achieve a specified target coordinate uncertainty has been demonstrated (Fisher *et al.*, 2008).

### 2.2.11. The most commonly used quality indicators

A summary of the most commonly used quality indicators and their definitions is presented in Table 2.2.11.1 for ready reference.

We gratefully acknowledge the contributions of Kevin Cowtan (York, UK), Kay Diederichs (Konstanz, Germany), Phil Evans (Cambridge, UK), John Helliwell (Manchester, UK), Randy Read (Cambridge, UK), George Sheldrick (Göttingen, Germany) and Tom Terwilliger (Los Alamos, USA).

### References

- Arnberg, L., Hovmöller, S. & Westman, S. (1979). *On the significance of 'non-significant' reflexions*. *Acta Cryst.* **A35**, 497–499.
- Blow, D. M. (2002). *Rearrangement of Cruickshank's formulae for the diffraction-component precision index*. *Acta Cryst.* **D58**, 792–797.
- Blow, D. M. & Crick, F. H. C. (1959). *The treatment of errors in the isomorphous replacement method*. *Acta Cryst.* **12**, 794–802.
- Blundell, T. L. & Johnson, L. N. (1976). *Protein Crystallography*. New York: Academic Press.
- Brünger, A. T. (1992). *Free R value: a novel statistical quantity for assessing the accuracy of crystal structures*. *Nature (London)*, **355**, 472–475.
- Brünger, A. T. (1997). *Free R value: cross-validation in crystallography*. *Methods Enzymol.* **277**, 366–396.
- Collaborative Computational Project, Number 4 (1994). *The CCP4 suite: programs for protein crystallography*. *Acta Cryst.* **D50**, 760–763.
- Cowtan, K. (1999). *Error estimation and bias correction in phase-improvement calculations*. *Acta Cryst.* **D55**, 1555–1567.
- Cowtan, K. & Main, P. (1998). *Miscellaneous algorithms for density modification*. *Acta Cryst.* **D54**, 487–493.
- Cowtan, K. D. & Main, P. (1996). *Phase combination and cross validation in iterated density-modification calculations*. *Acta Cryst.* **D52**, 43–48.
- Cruickshank, D. W. J. (1999a). *Remarks about protein structure precision*. *Acta Cryst.* **D55**, 583–601.
- Cruickshank, D. W. J. (1999b). *Remarks about protein structure precision*. *Erratum*. *Acta Cryst.* **D55**, 1108.
- Cullis, A. F., Muirhead, H., Perutz, M. F., Rossmann, M. G. & North, A. C. T. (1961). *The structure of haemoglobin. VIII. A three-dimensional Fourier synthesis at 5.5 Å resolution: determination of the phase angles*. *Proc. R. Soc. London Ser. A*, **265**, 15–38.
- Diederichs, K. (2006). *Some aspects of quantitative analysis and correction of radiation damage*. *Acta Cryst.* **D62**, 96–101.
- Diederichs, K. (2010). *Quantifying instrument errors in macromolecular X-ray data sets*. *Acta Cryst.* **D66**, 733–740.