

## Chapter 2.2. Quality indicators in macromolecular crystallography: definitions and applications

H. M. EINSPAHR AND M. S. WEISS

### 2.2.1. Introduction

The genesis of this chapter was a perceived need for a single location in the volume in which consensus definitions could be found for the many statistical indicators of quality or figures of merit that have been developed to monitor a macromolecular crystallography (MX) experiment and its final outcome, a model. The evolving experiment has generated a rich collection of *R* values, signal-to-noise indicators, correlation coefficients and other figures of merit. As improvements in data collection, processing and other aspects of the experiment continue, we can expect a continued evolution of new or improved indicators of quality with which to monitor the impact of those improvements.

This chapter, then, attempts to provide a comprehensive list of the indicators of quality currently in use and, for each indicator in the list, a precise definition that conforms to consensus interpretations of the literature and current practice. The authors acknowledge that useful indicators may have been missed in the sweep that produced this list, but hasten to point out that the generation of new indicators is an ongoing process with the newest ones often in obscurity for a period before their utility is recognized and adopted by experimenters. There is also a subset of the indicators in the list whose members stand out as universally accepted and crucial indicators that either every experimental description must include or that every experimenter should be familiar with. A summary of these is given at the end of this article in Table 2.2.11.1. The authors also acknowledge that there may be competing definitions for some indicators. Where they occur, these competing definitions will be pointed out and a discussion will be provided which, at a minimum, will attempt to clarify the differences, but will also, where possible, attempt to arbitrate or discriminate among opposing views. It should be noted that the scope of this chapter is primarily concerned with the crystallographic experiment. While quality indicators useful for model refinement are given in Sections 2.2.8 to 2.2.10, the validation of the refined model is covered more extensively in Part 21 of this volume.

Before proceeding with the quality indicators themselves, it is necessary to discuss a small collection of parameters whose values, typically left to the experimenter to fix, impact on virtually all the indicators of quality we present here. The reason for this wide impact is that they are used to determine which reflections are included in the final data set in question. In their simplest form, the effect is a limit, applied as a cutoff in intensity or resolution, beyond which reflections are excluded from consideration.

When a cutoff is applied based on intensities, the limiting value may be simply a fixed number, for example zero, or it may vary from reflection to reflection, for example some multiple of the standard uncertainty of the reflection intensity. Reflections with intensities below limiting values are excluded. Reflections excluded in this manner are often referred to as unobserved.

When the objective is to estimate electron-density distributions, the exclusion of intensity terms below very low limiting

values is unlikely to have any significant negative impact and may offer a positive savings in computation time. This may not be the case, however, when the objective is refinement. Inclusion of reflections of low intensity, even zero, may have important positive results on the quality of the final model. True, some of these intensities may be poorly known or, if less than zero, physically unrealistic, but given the excellent quality of modern data-collection instrumentation and techniques, few intensities, if any, are without some reasonable estimate of their standard uncertainty, so that weighting procedures can be applied to modulate the impact of individual terms in the final result and judgments can be made about exclusion of intensity values that are truly improbable. With the application of proper weighting procedures, there seems to be little justification for exclusion of reflections from consideration based on intensity (except for the truly improbable mentioned above) when computing indicators of quality. As instrumentation and techniques continue to improve, this is a topic that merits continued attention and debate within the MX community in search of consensus best practices for handling weak reflections and for including them in estimators of quality.

While an acceptable estimate of the nominal resolution of a diffraction data set is widely considered to be a high-value indicator of data quality, assignment of a value to this limit is typically left to the experimenter and is thus prone to subjectivity. Some guidelines for estimation have emerged. One would set the nominal resolution based on the percentage of weak reflections above the limit so that reflections with intensities above that limit would be included. An example might be the resolution at which 70% of unique reflections have intensities above zero or above some multiple of their standard uncertainties. Another way to estimate the nominal resolution which has gained wide acceptance is based on the overall signal-to-noise value where, in its most popular expression, the limit value is the resolution at which the mean signal-to-noise ratio in the outer resolution shell falls to 2. Each of these estimation methods is susceptible to convolution with limits based on intensity leading to, for example, a limit set as the resolution at which the mean signal-to-noise ratio 'of observed reflections' falls to 2. The guidelines for estimating resolution limits – the methods to be applied, the constraining values such as 70% or 2, and the proper integration with limits applied based on intensity – also merit further attention and discussion by the MX community, the goal being definition of consensus best practices. We take this opportunity to suggest that the widely accepted estimate, the resolution at which the mean signal-to-noise in the outer resolution shell falls to 2, calculated without imposition of a cutoff in intensity, has much to recommend it. True, some susceptibility to subjectivity remains in the definition of resolution-shell ranges and limits, but the impact is minimal and hardly worth the effort. On the other hand, modifying the indicator to remove the shells so that the mean signal-to-noise ratio applied to all data above an appropriately adjusted limit value would extinguish that source of subjectivity.

When a cutoff is applied based on nominal resolution, then considerations similar to those for intensity-based cutoffs apply.

## 2.2. QUALITY INDICATORS

If the objective is to estimate electron-density distributions, exclusion of large numbers of weak intensity terms beyond a limiting resolution is unlikely to be of any significant impact except, perhaps, a favourable one in computation time. If a sparse population of more intense reflections is also excluded, the effect may also be positive by reducing aberrations that might interfere with interpretation. If the objective is refinement, the benefit is less clear, except possibly in computation time. The cost for that may be the exclusion of a few more intense terms that provide positive guidance to refinement. Another positive effect of a sharply defined limiting resolution might be improved estimation of the spherical interference function where it is needed. It appears that, in practice, the consensus is that, in applying a cutoff, the exclusion of a few intense terms is at worst of negligible impact overall.

Because many imposed limits that exclude data are dependent on estimates of standard uncertainty, procedures for estimation of standard uncertainty need to be considered. As suggested earlier, a diffraction data set without a standard uncertainty for each intensity measurement is certainly the exception in current practice. The methods used by data-processing programs to estimate these standard uncertainties may be difficult, even impossible, to discern, but at core they must all be based on counting statistics. With the continuing trend toward data sets of high multiplicity (also referred to as redundancy), however, estimates of standard uncertainties are available from distributions of replicate measurements about means. Both of these estimates have value and express impacts of somewhat different sources of error. It would seem therefore that, where applicable, the best way to accommodate both would be to calculate weighted average values of individual intensities, where the weights are derived from the standard uncertainties from data processing, and then to estimate the standard uncertainties for the weighted averages by application of standard propagation-of-error procedures. This, finally, is the third area we highlight here that deserves focused discussion within the MX community with the objective of defining consensus best practices.

Little remains to be included in these introductory preliminaries. It should be understood that most of these indicators of quality may be cast in terms of either structure factors or intensities. We draw little attention to this as we define individual indicators, except that, where one of the two possibilities seems to dominate in usage, we tend to define that form alone. It is also true that many of these indicators have counterparts in which individual reflection terms are weighted. In current practice, these forms find little use in application to biological macromolecules and we ignore them here. Finally, many of these indicators may be expressed as fractions, that is, as numbers between one and zero. They are also often expressed in the percentage form, that is, as numbers between one hundred and zero. While we confine our use to the former here, we make no statement of preference.

The remainder of this chapter consists primarily of individual sections that reflect the various steps from the crystallographic data-collection experiment to the refinement of the final model. Relevant quality indicators and definitions are given for each of the steps and the most commonly used are collected in Table 2.2.11.1.

### 2.2.2. Quality indicators for diffraction data

Once useful crystalline samples have been obtained, the collection of X-ray diffraction data is the next (and the last) experi-

mental step in a structure determination. Although the greatest care may be taken to collect data of as high quality as possible, there remain circumstances and influences that limit the quality of the data. Over time, many indicators have been defined to describe various aspects of diffraction data quality. The most important ones are discussed here.

**Nominal resolution,  $d_{\min}$ .** The resolution of a diffraction data set describes the extent of measurable data and is calculated by Bragg's law [equation (2.2.2.1)] based on the maximum Bragg angle  $2\theta$  included in the data set for a given data-collection wavelength  $\lambda$ . As discussed above, the nominal resolution is a limit set by the experimenters and is well known to be prone to subjective judgment. A number of suggestions have been made to reduce the subjectivity associated with this limit. One defines the limit as the resolution within which the intensities of a fraction of the unique reflections, for example 70%, are above a threshold, for example zero or three times their standard uncertainties. Another suggested limiting criterion, discussed further below, recommends that the nominal resolution be set as the midpoint of the resolution range of the shell at which the mean signal-to-noise ratio falls below 2.

$$d_{\min} = \lambda/2 \sin \theta. \quad (2.2.2.1)$$

**True resolution,  $d_{\text{true}}$ .** The true resolution of a diffraction data set is defined as the minimum distance between two objects in a crystal that permits their images in the resultant electron-density map to be resolved. Often,  $d_{\text{true}}$  is approximated as  $d_{\min}$ .

To illustrate this crucial distance, represent two equivalent atoms by equal overlapping Gaussians. One might then consider that the distance between them that just permits distinguishing them as individual atoms might be the distance at which the electron-density value at the midpoint between the atoms drops to a value just below that at the positions of the atoms. For a normal distribution, this distance is  $2\sigma$ , twice the standard deviation of the distribution.

Another perspective is provided by the realization that, when a Fourier synthesis is terminated at a resolution cutoff  $d_{\min}$ , successive spheres of negative and positive density of decreasing amplitude surround the maxima of positive density at atomic positions in that synthesis. It has been shown that the distance from the centre of such a maximum to the first zero is  $0.715d_{\min}$  (James, 1948), which is a useful estimate of the limiting distance between distinguishable features in the electron-density map or  $d_{\text{true}}$ . Similar estimates are used in other areas, notably for defining resolution in astronomy. A more recent re-evaluation suggests that a limit of  $0.917d_{\min}$  is a better value, especially when the effects of form factors and atomic displacement parameters are considered (Stenkamp & Jensen, 1984). Add to that the effects of errors in experimental amplitudes and derived phases, and the approximation of  $d_{\text{true}}$  as  $d_{\min}$  seems quite reasonable.

**Optical resolution,  $d_{\text{opt}}$ .** The optical resolution  $d_{\text{opt}}$  is calculated from the standard deviation of a Gaussian fitted to the origin peak of the Patterson function ( $\sigma_{\text{Patt}}$ ) of the diffraction data set and the standard deviation of another Gaussian fitted to the origin peak of the spherical interference function ( $\sigma_{\text{sph}}$ ). This definition is based on Vaguine *et al.* (1999) and is implemented in the program *SFCHECK*. The optical resolution is intended to account for uncertainties in the data, atomic displacement factors, effects of crystal quality and series-termination effects by means of a propagation-of-error-like approach (Blundell & Johnson, 1976; Vaguine *et al.*, 1999). It has been suggested that  $d_{\text{opt}}$  is a better approximation of  $d_{\text{true}}$  than

## 2. BASIC CRYSTALLOGRAPHY

$d_{\min}$  (Weiss, 2001).

$$d_{\text{opt}} = [2(\sigma_{\text{Patt}}^2 + \sigma_{\text{sph}}^2)]^{1/2}. \quad (2.2.2.2)$$

**Completeness,  $C$ .** The completeness  $C$  of a diffraction data set is defined as the fraction of the unique reflections in a given space group to a given nominal resolution  $d_{\min}$  that have been measured at least once during data collection.  $C$  may be given assuming that Friedel symmetry is either applied or not. In the latter case,  $C$  is also referred to as anomalous completeness. In the program *SCALA* (Evans, 2006), the anomalous completeness is defined based on acentric reflections only.

**Effective resolution,  $d_{\text{eff}}$ .** Since any missing reflection of a data set leads to a deterioration of the model parameters (Hirshfeld & Rabinovich, 1973; Arnberg *et al.*, 1979), an effective resolution may be defined based on the nominal resolution  $d_{\min}$  and the cube root of the completeness  $C$  of the data set.

$$d_{\text{eff}} = d_{\min} C^{-1/3}. \quad (2.2.2.3)$$

**Multiplicity (or redundancy),  $N$ .** The multiplicity or redundancy  $N$  of a diffraction data set defines on average how many times a reflection  $hkl$  has been observed during the data-collection experiment including symmetry mates and replicate measurements.  $N$  may be given assuming that Friedel symmetry is either applied or not.

**Merging  $R$  factor,  $R_{\text{merge}}$ .** The merging  $R$  factor of a diffraction data set describes the spread of the individual intensity measurements  $I_i$  of a reflection  $hkl$  around the mean intensity  $\langle I(hkl) \rangle$  of this reflection. Sometimes  $R_{\text{merge}}$  is also referred to as  $R$  factor (observed) (in the program *XDS*; Kabsch, 1988, 1993, 2010; Chapter 11.6), as  $R_{\text{sym}}$  or as  $R_{\text{linear}}$  (in the program *SCALEPACK*; Otwinowski & Minor, 1997). In fractional form, this is

$$R_{\text{merge}} = \frac{\sum_{hkl} \sum_i |I_i(hkl) - \langle I(hkl) \rangle|}{\sum_{hkl} \sum_i I_i(hkl)}, \quad (2.2.2.4)$$

where  $\langle I(hkl) \rangle$  is the mean of the several individual measurements  $I_i(hkl)$  of the intensity of reflection  $hkl$ . The sums  $\sum_{hkl}$  and  $\sum_i$  run over all observed unique reflections  $hkl$  and over all individual observations  $i$  of a given reflection  $hkl$ . It should be noted that alternative definitions of  $R_{\text{merge}}$  exist. In one,  $I_i(hkl)$  in the denominator is replaced by  $\langle I(hkl) \rangle$ , thereby producing an expression that is formally equivalent to the one above. In another,  $I_i(hkl)$  in the denominator is replaced by  $|I_i(hkl)|$  with the suggestion that the denominator is thereby prevented from becoming negative or zero, even in the case of many negative-intensity observations. One should note, however, the counter-intuitive side effect: artificial damping of  $R_{\text{merge}}$  values, that is, reducing expected higher  $R_{\text{merge}}$  values of data sets with more weak reflections.

The usefulness of  $R_{\text{merge}}$  as a quality indicator for diffraction data is limited because it is dependent on the multiplicity of a data set (Diederichs & Karplus, 1997a,b; Weiss & Hilgenfeld, 1997; Weiss, 2001). The higher the multiplicity of a data set, the higher its  $R_{\text{merge}}$  will be, although, based on statistics, the better determined the averaged intensity values should be. Despite these shortcomings,  $R_{\text{merge}}$  is still widely used today.

**Redundancy-independent merging  $R$  factor,  $R_{\text{r.i.m.}}$  or  $R_{\text{meas}}$ .** The redundancy-independent merging  $R$  factor  $R_{\text{r.i.m.}}$  or  $R_{\text{meas}}$  describes the precision of the individual intensity measurements  $I_i$ , independent of how often a given reflection has been measured. Because of its independence of the redundancy (hence its name), it has been proposed that  $R_{\text{r.i.m.}}$  or  $R_{\text{meas}}$  should be used as a substitute for the conventional  $R_{\text{merge}}$  (Diederichs & Karplus,

1997a,b; Weiss & Hilgenfeld, 1997; Weiss, 2001). In fractional form, this is

$$R_{\text{r.i.m.}} = \frac{\sum_{hkl} \{N(hkl)/[N(hkl) - 1]\}^{1/2}}{\sum_i |I_i(hkl) - \langle I(hkl) \rangle| / \sum_{hkl} \sum_i I_i(hkl)}, \quad (2.2.2.5)$$

where  $\langle I(hkl) \rangle$  is the mean of the  $N(hkl)$  individual measurements  $I_i(hkl)$  of the intensity of reflection  $hkl$ . As for  $R_{\text{merge}}$ , the sums  $\sum_{hkl}$  and  $\sum_i$  run over all observed unique reflections  $hkl$  and over all individual observations  $i$  of a given reflection  $hkl$ .

**Precision-indicating merging  $R$  factor,  $R_{\text{p.i.m.}}$ .** The precision-indicating merging  $R$  factor  $R_{\text{p.i.m.}}$  describes the precision of the averaged intensity measurements  $\langle I(hkl) \rangle$  (Weiss, 2001). In fractional form, this is

$$R_{\text{p.i.m.}} = \frac{\sum_{hkl} \{1/[N(hkl) - 1]\}^{1/2}}{\sum_i |I_i(hkl) - \langle I(hkl) \rangle| / \sum_{hkl} \sum_i I_i(hkl)}, \quad (2.2.2.6)$$

where  $\langle I(hkl) \rangle$  is the mean of the  $N(hkl)$  individual measurements  $I_i(hkl)$  of the intensity of reflection  $hkl$ . As with  $R_{\text{merge}}$  and  $R_{\text{r.i.m.}}$  or  $R_{\text{meas}}$ , the sums  $\sum_{hkl}$  and  $\sum_i$  run over all observed unique reflections  $hkl$  and over all individual observations  $i$  of a given reflection  $hkl$ .

**$R$  factor of merged intensities or amplitudes,  $R_{\text{mrgd-I}}$  and  $R_{\text{mrgd-F}}$ .** An alternative precision-indicating merging  $R$  factor, called  $R_{\text{mrgd}}$ , is defined as the  $R$  factor between two or more data sets or between two subsets of a data set created by randomly apportioning the individual intensity measurements between the two subsets (Diederichs & Karplus, 1997a,b).  $R_{\text{mrgd}}$  can be calculated for intensities ( $R_{\text{mrgd-I}}$ ) or structure-factor amplitudes ( $R_{\text{mrgd-F}}$ ). The latter quantity was suggested to present a lower limit for the crystallographic  $R$  factor of a model against the observed data (Diederichs & Karplus, 1997a,b). In fractional form

$$R_{\text{mrgd-I}} = 2 \frac{\sum_{hkl} |\langle I_1(hkl) \rangle - \langle I_2(hkl) \rangle|}{\sum_{hkl} \langle I_1(hkl) \rangle + \langle I_2(hkl) \rangle}, \quad (2.2.2.7)$$

where  $\langle I_1(hkl) \rangle$  and  $\langle I_2(hkl) \rangle$  are the mean intensity values for the individual observations of the reflections  $hkl$ , which have been partitioned into the two subsets 1 and 2. The sums  $\sum_{hkl}$  run over all observed unique reflections.  $R_{\text{mrgd-I}}$  is related to  $R_{\text{p.i.m.}}$  by a constant factor ( $R_{\text{mrgd-I}} = 2^{1/2} R_{\text{p.i.m.}}$ ).

$R_{\text{mrgd-F}}$  is defined analogously to  $R_{\text{mrgd-I}}$  (Diederichs & Karplus, 1997a,b). In the equation, only the intensities are replaced by structure-factor amplitudes.

$$R_{\text{mrgd-F}} = 2 \frac{\sum_{hkl} |\langle F_1(hkl) \rangle - \langle F_2(hkl) \rangle|}{\sum_{hkl} \langle F_1(hkl) \rangle + \langle F_2(hkl) \rangle}. \quad (2.2.2.8)$$

In order to cope with negative-intensity observations, pseudo-amplitudes had to be introduced just for the purpose of calculating  $R_{\text{mrgd-F}}$  ( $F = I^{1/2}$  if  $I \geq 0$  and  $F = -|I|^{1/2}$  if  $I < 0$ ).

*Note.* The approach of comparing randomly partitioned subsets of a given data set is used for a variety of quality indicators. While there is potential for variation in these indicators from one partitioning of the data set to another, an average of several random partitionings should be expected to give a useful estimate. There is also potential for subjectivity,

## 2.2. QUALITY INDICATORS

but the principal value of these indicators is to assist the experimenter in proper analysis and they are less often applied to compare experiments from different laboratories and are seldom published.

**Pooled coefficient of variation, PCV.** The pooled coefficient of variation PCV is the ratio of the sum of the standard deviations to the sum of the reflection intensities (Diederichs & Karplus, 1997a,b). PCV is related to  $R_{\text{meas}}$  or  $R_{\text{r.i.m.}}$  by the factor  $(\pi/2)^{1/2}$ . In fractional form, this is

$$\text{PCV} = \frac{\sum_{hkl} (\{1/[N(hkl) - 1]\}^{1/2} \sum_i |I_i(hkl) - \langle I(hkl) \rangle|^2)^{1/2}}{\sum_{hkl} \sum_i I_i(hkl)}, \quad (2.2.2.9)$$

where  $\langle I(hkl) \rangle$  is the mean of the  $N(hkl)$  individual measurements  $I_i(hkl)$  of the intensity of reflection  $hkl$ .

**Mean signal-to-noise ratio,  $\langle I \rangle / \sigma(I)$ .** The signal-to-noise ratio  $I_i / \sigma(I_i)$  of an individual intensity measurement describes the statistical significance of a measured intensity. As a measure of the overall quality of a data set, the mean signal-to-noise ratio for all reflections is useful as an indication of the robustness of the data, that is, the average intensity as a multiple of the standard uncertainty. In addition, as mentioned above, the mean signal-to-noise ratio for all reflections within the outer resolution shell can be used to define the nominal resolution of a data set. For the data set as a whole or for a resolution shell of that data set, the mean signal-to-noise ratio,  $\langle I \rangle / \sigma(I)$ , is the sum of the signal-to-noise ratios of all individual reflections  $hkl$  within resolution limits divided by the number of individual reflections  $hkl$  within those resolution limits.

In principle there are two ways to define a mean signal-to-noise ratio of a data set (or a given resolution shell). The two ways yield different quantities, although, unfortunately, they are both called the mean signal-to-noise ratio. They differ in the manner in which mean signal-to-noise ratios of individual reflections  $hkl$  are calculated.

- (i)  $\langle I(hkl) \rangle / \sigma[I(hkl)]$ . The mean signal-to-noise ratio of individual reflections  $hkl$  may be calculated as the ratio of the mean intensity  $\langle I(hkl) \rangle$  and the r.m.s. scatter of  $I_i(hkl)$  about that mean. This is a measure of the average significance of individual observations, but it does not take into account the multiplicity or redundancy of the measurements.

$$\begin{aligned} \langle I(hkl) \rangle / \sigma[I(hkl)] \\ = \langle I(hkl) \rangle / [(1/N) \sum_i |I_i(hkl) - \langle I(hkl) \rangle|^2]^{1/2}. \end{aligned} \quad (2.2.2.10)$$

In the program *SCALA* (Evans, 2006), this value is reported as  $I/\sigma$ .

- (ii)  $\langle I(hkl) \rangle / \sigma \langle I(hkl) \rangle$ .  $\langle I(hkl) \rangle$  is the average over all observations of the reflection  $hkl$ , and is sometimes weighted.  $\sigma \langle I(hkl) \rangle$  is the propagation-of-error combination of standard uncertainties assigned at data processing for the individual measurements  $I_i(hkl)$ , that is, a modification of equation (2.2.2.10) in which the term  $|I_i(hkl) - \langle I(hkl) \rangle|$  in the denominator is replaced by  $\sigma_i(hkl)$ , the experimental standard uncertainty for the measurement  $I_i(hkl)$ .

$$\langle I(hkl) \rangle / \sigma \langle I(hkl) \rangle = \langle I(hkl) \rangle / [(1/N) \sum_i \sigma_i(hkl)^2]^{1/2}. \quad (2.2.2.11)$$

An error model<sup>1</sup> is often applied in the denominator here to scale to the r.m.s. scatter in (i) above. In the program *SCALA* (Evans, 2006), this value is reported as  $\text{Mn}(I)/\text{sd}$ .

Both methods of defining the mean signal-to-noise ratio for the reflection  $hkl$  have merit. As suggested for individual intensities in Section 2.2.1, perhaps the best approach would be to calculate weighted averages and weighted standard uncertainties of the  $I(hkl)$  where weights are the experimental standard uncertainties  $\sigma_i(hkl)$  for individual measurements  $I_i(hkl)$ .

**Highest possible signal-to-noise ratio,  $I/\sigma(I)_{\text{asymptotic}}$ .** A relatively recent addition to the collection of diffraction-data quality indicators is the highest possible signal-to-noise ratio of a data set  $I/\sigma(I)_{\text{asymptotic}}$  or ISa (Diederichs, 2010). ISa is calculated from the parameters of the error model used for inflating the standard deviations of the reflections with an intensity-dependent term.<sup>1</sup> Since ISa is practically independent of counting statistics, it was suggested to be a good measure of instrument errors manifesting themselves in the data set, provided the crystal is close to ideal and radiation damage is negligible. Data sets with ISa values of 25 or greater are considered to be very good and amenable to straightforward structure determination, while data sets exhibiting ISa values of 15 or less are considered marginal at best. The calculation of ISa is implemented in *XDS* versions of December 2009 or later (Kabsch, 2010).

**Anomalous  $R$  factor,  $R_{\text{anom}}$ .** The anomalous  $R$  factor  $R_{\text{anom}}$  describes the sum of the differences in intensities of Friedel-related reflections ( $hkl$ ) and ( $\bar{h}\bar{k}\bar{l}$ ) relative to the sum of their mean intensities. In fractional form, this is

$$R_{\text{anom}} = \frac{\sum_{hkl} |I(hkl) - I(\bar{h}\bar{k}\bar{l})|}{\sum_{hkl} \langle I(hkl) \rangle}, \quad (2.2.2.12)$$

where, in this case,  $\langle I(hkl) \rangle$  is the mean intensity of the Friedel mates of the reflections  $hkl$ , or  $\frac{1}{2}[I(hkl) + I(\bar{h}\bar{k}\bar{l})]$ . Here, the sums  $\sum_{hkl}$  run over all unique reflections with one of the indices, typically  $h$ , greater than zero ( $h > 0$ ) for which both Friedel mates have been observed at least once.

The ratio of  $R_{\text{anom}}$  to  $R_{\text{p.i.m.}}$  has been proposed as a possible indicator for the strength of the anomalous signal (Panjikar & Tucker, 2002).

**Anomalous correlation coefficient,  $\text{CC}_{\text{anom}}$ .** The anomalous correlation coefficient  $\text{CC}_{\text{anom}}$  quantifies the linear dependence of observed anomalous differences in two diffraction data sets. These can be data sets, for example, collected at two different wavelengths in a MAD experiment. In cases where only one data set is available, two randomly partitioned half data sets can be created for comparison.

*Note.* The correlation coefficient referred to here and elsewhere in this chapter is invariably the Pearson linear correlation coefficient (Rodgers & Nicewander, 1988):

$$\text{CC} = \frac{\sum (x - \langle x \rangle)(y - \langle y \rangle)}{[\sum (x - \langle x \rangle)^2 \sum (y - \langle y \rangle)^2]^{1/2}}, \quad (2.2.2.13)$$

with  $x$  and  $y$  being, in this case, the anomalous differences  $[I(hkl) - I(\bar{h}\bar{k}\bar{l})]$  or  $[F(hkl) - F(\bar{h}\bar{k}\bar{l})]$  in the two data sets,  $\langle x \rangle$  and  $\langle y \rangle$  are their averages, and the summations are over all reflections

<sup>1</sup> Typically, the experimental standard uncertainties are modified by various correction factors in order to account for errors in the data that arise from other sources than counting statistics. The inflation factors are applied in various ways in different programs. In *SCALA* (Collaborative Computational Project, Number 4, 1994; Evans, 2006) they are called SDFAC, SDB and SDADD. In *SCALEPACK* (Otwinski & Minor, 1997), they are called ERROR SCALE FACTOR and ESTIMATED ERROR. In *D\*TREK* (Pflugrath, 1999), they are called  $E_{\text{mul}}$  and  $E_{\text{add}}$ . In *XDS* (Kabsch, 1988, 1993, 2010; Chapter 11.6) they are called  $a$  and  $b$ .

## 2. BASIC CRYSTALLOGRAPHY

$hkl$  for which observations exist in both data sets across the entire resolution range or within a particular resolution shell.  $CC_{\text{anom}}$  is a reliable indicator of the strength of the anomalous signal. Values above 0.30 are considered good.

**R.m.s. correlation ratio.** This is another statistic based on randomly partitioned data sets, which is calculated by the program *SCALA* (Evans, 2006; Collaborative Computational Project, Number 4, 1994). It is an analysis of the scatterplot of  $[I_1(hkl) - I_1(\bar{h}\bar{k}\bar{l})]$  versus  $[I_2(hkl) - I_2(\bar{h}\bar{k}\bar{l})]$ , where the subscripts 1 and 2 identify the two half data sets. The analysis assumes that the correlation is ideally 1.0. The r.m.s. correlation ratio is defined as the ratio of the r.m.s. widths of the scatterplot distribution along the diagonal and perpendicular to the diagonal. This statistic seems to be more robust than  $CC_{\text{anom}}$  to the presence of outliers. It cannot, however, be applied to analysing the correlations between different data sets.

**Mean anomalous signal-to-noise ratio,  $\langle d''/\sigma(d'') \rangle$ .** The anomalous signal-to-noise ratio of an individual reflection measurement  $d''(hkl)/\sigma[d''(hkl)]$  is defined as the ratio of the observed anomalous intensity difference  $d'' = |I(hkl) - I(\bar{h}\bar{k}\bar{l})|$  and the corresponding estimated standard uncertainty  $\sigma(d'')$  in the measurement of this anomalous difference. The average of the anomalous signal-to-noise ratios for all reflections within a certain resolution range is used as an indicator of utility for phasing. A value of  $(2/\pi)^{1/2} \simeq 0.8$  for mean  $d''/\sigma(d'')$  of a resolution shell, for example, is taken to indicate that no anomalous signal is present (G. Sheldrick & G. Bunkoczi, personal communication).

**Decay  $R$  factor,  $R_d$ .** The decay  $R$  factor  $R_d$  is defined as a pairwise  $R$  factor based on the intensities of symmetry-related reflections occurring on different diffraction images (Diederichs, 2006). An increase in  $R_d$  as a function of difference in image-collection times is a good indicator of radiation damage occurring during data collection. In fractional form, this is

$$R_d = 2 \sum_{hkl} \sum_{m-n} |I_m(hkl) - I_n(hkl)| / \sum_{hkl} \sum_i [I_m(hkl) + I_n(hkl)], \quad (2.2.2.14)$$

where  $I_m(hkl)$  and  $I_n(hkl)$  are the intensities of the reflection  $hkl$  occurring on images  $m$  and  $n$ . The only program in which this is currently implemented is *XDSSTAT*.

**Wilson-plot  $B$  factor,  $B_{\text{Wilson}}$ .** A Wilson plot (Wilson, 1949) is a plot for a contiguous series of resolution shells of the logarithm of the mean intensity in a given resolution shell divided by the sum of the squared atomic form factors for all atoms in the unit cell evaluated at the mean of the resolution limits of the shell. From a least-squares fit of a straight line to the linear part of the Wilson plot, the  $B$  factor  $B_{\text{Wilson}}$  can be derived. Typically, data of lower than 4.5 Å resolution are excluded from the fit. The more meaningful determinations of  $B_{\text{Wilson}}$  come from Wilson plots that are linear all the way to the nominal resolution  $d_{\text{min}}$  and minimize the occurrence of spikes due to ice rings.

$$\ln(\langle I_{\text{obs}} \rangle / \sum_i f_i^2) = -\ln K_{\text{Wilson}} - 2B_{\text{Wilson}}/d^2, \quad (2.2.2.15)$$

where  $\langle I_{\text{obs}}(hkl) \rangle$  is the mean over the intensities of all observed reflections  $hkl$  in a given resolution shell. The sum  $\sum_i$  runs over all atoms in the structure. The parameter  $d$  is the midpoint of the resolution shell over which  $I_{\text{obs}}$  has been averaged.  $K_{\text{Wilson}}$  is an absolute scale factor.

### 2.2.3. Comparing different diffraction data sets

In MX, there are many instances when two or more diffraction data sets have to be compared to each other. An important decision which has to be made is, for instance, whether data sets collected from different crystals are isomorphous enough so that they can be merged together. Another task is the comparison of a native data set and a heavy-atom derivative. The question here is how much of the observed difference is due to non-isomorphism and how much is due to isomorphous replacement.

**Scaling  $R$  factor,  $R_{\text{scale}}$ .** The scaling  $R$  factor  $R_{\text{scale}}$  between two data sets is defined as the difference in the structure-factor amplitudes of one data set relative to the structure-factor amplitudes of the other data set after the two data sets have been scaled to each other.

$$R_{\text{scale}} = \sum_{hkl} |F_1 - F_2| / \sum_{hkl} |F_1|. \quad (2.2.3.1)$$

However, the above formula is asymmetric with respect to data sets 1 and 2. An alternative symmetric formula is

$$R_{\text{scale}} = 2 \sum_{hkl} |F_1 - F_2| / \sum_{hkl} (F_1 + F_2). \quad (2.2.3.2)$$

$R_{\text{scale}}$  is also often given based on intensities rather than structure-factor amplitudes.

**Gradient from a normal probability analysis.** For each reflection  $hkl$  in two data sets that have been scaled to each other, the quantity  $\Delta(\text{real})$  is computed and compared with its expectation value  $\Delta(\text{expected})$  assuming a normal distribution of errors (Howell & Smith, 1992). If  $\Delta(\text{real})$  is plotted against  $\Delta(\text{expected})$ , random differences between the two data sets give rise to a slope of 1.0. Significant differences result in slopes significantly larger than 1.0. Such an analysis is implemented in the program *SCALEIT* (Collaborative Computational Project, Number 4, 1994).

$$\Delta(\text{real}) = (F_1 - F_2) / [\sigma(F_1)^2 + \sigma(F_2)^2]^{1/2}. \quad (2.2.3.3)$$

### 2.2.4. Quality indicators for substructure determination

The first step of a crystal structure determination after processing and analysing the diffraction data is to use the observed isomorphous and/or anomalous differences in order to determine the isomorphous or anomalous substructure. Since such substructures typically consist of rather few atoms, which are relatively far apart, direct-methods approaches have been adopted from small-molecule crystallography and have largely superseded pure Patterson-based methods. The most widely used computer programs for substructure determination are *SHELXD* (Schneider & Sheldrick, 2002; Sheldrick, 2008), *SnB* (Weeks *et al.*, 1993, 1994) and *HySS* (Grosse-Kunstleve & Adams, 2003).

**Correlation coefficient  $CC(\text{all})$ .**  $CC(\text{all})$  measures the correlation of the  $E$  values (analogous to normalized structure factors; Hauptman & Karle, 1953) derived from the observed isomorphous or anomalous differences and those calculated from the substructure model.  $CC(\text{all})$  is another application of the Pearson linear correlation coefficient [see equation (2.2.2.13)] and is calculated for all reflections. In the equation,  $x$  and  $y$  are the  $E$  values derived from the observed and calculated differences.  $CC(\text{all})$  is typically used to assess the successful determination of the isomorphous or anomalous substructure in a macromolecule. A value of  $CC(\text{all}) \geq 0.3$  often indicates that the substructure has been correctly identified (Schneider & Sheldrick, 2002).

**Correlation coefficient  $CC(\text{weak})$ .** Like  $CC(\text{all})$ ,  $CC(\text{weak})$  is the Pearson linear correlation coefficient [see equation (2.2.2.13)]

## 2.2. QUALITY INDICATORS

between the  $E$  values (analogous to normalized structure factors as above) derived from the observed isomorphous or anomalous differences and those calculated from the substructure model. In contrast to CC(all), however, CC(weak) is calculated for the weak reflections only. As above,  $x$  and  $y$  are the  $E$  values derived from the observed and calculated differences. The  $E$ -value cutoff for defining a reflection as weak can be chosen by the user, but a typical threshold value is 1.5, although lower values may be required for low-resolution data (Schneider & Sheldrick, 2002). A value of CC(weak)  $\geq 0.15$  often indicates that the substructure has been correctly identified (Sheldrick, 2010).

**The minimal function,  $R(\varphi)$ .** The minimal function  $R(\varphi)$  is a measure of the mean-square difference between the values of the triplets calculated using a particular set of phases and the expected values of the same triplets as given by the ratio of modified Bessel functions. The minimal function is expected to have a constrained global minimum when the phases are equal to their correct values for some choice of origin and enantiomorph (the minimal principle).

$$R(\varphi) = \sum_{H,K} A_{HK} \left\{ \cos \varphi_{HK} - [I_1(A_{HK})/I_0(A_{HK})] \right\}^2 / \sum_{H,K} A_{HK}, \quad (2.2.4.1)$$

where  $A_{HK} = (2/N^{1/2})|E_H E_K E_{H+K}|$  with  $N$  being the number of atoms in the corresponding primitive unit cell. The minimal function is the function minimized in the program *Shake&Bake*, abbreviated as *SnB* (Weeks *et al.*, 1993, 1994).

### 2.2.5. Quality indicators for phase determination

Once the isomorphous or anomalous substructure has been successfully determined, it can be used as reference point for the calculation of phases. The quality of the resulting phases is dependent on the strength of the isomorphous or anomalous signal and the completeness and correctness of the isomorphous or anomalous substructure.

**Cullis  $R$  factor,  $R_{\text{Cullis}}$ .** The Cullis  $R$  factor (Cullis *et al.*, 1961) for phase determination by isomorphous replacement is defined as the ratio between the lack-of-closure error  $\varepsilon(\varphi_P)$  [equation (2.2.5.1a) below] and the isomorphous difference  $|F_{\text{PH}} - F_P|$ . It is the most useful signal for a usable heavy-atom derivative. Values  $< 0.6$  for centrosymmetric data are excellent, while values  $< 0.9$  are still usable.

$$\varepsilon(\varphi_P) = |F_{\text{PH}} - |F_P + F_H||, \quad (2.2.5.1a)$$

$$R_{\text{Cullis}} = \sum_{hkl} |F_{\text{PH}} - |F_P + F_H|| / \sum_{hkl} |F_{\text{PH}} - F_P|. \quad (2.2.5.1b)$$

**Anomalous Cullis  $R$  factor,  $R_{\text{Cullis,ano}}$ .** The Cullis  $R$  factor for phase determination by anomalous dispersion is defined as the ratio between the lack-of-closure error and the observed anomalous difference  $|F_{\text{PH}}(hkl) - F_{\text{PH}}(\bar{h}\bar{k}\bar{l})|$ . The lack-of-closure error in the anomalous-dispersion case is the difference between the observed anomalous difference and the calculated anomalous difference  $2F_H \sin \alpha_P$ , where  $\alpha_P$  is the protein phase. A value of  $R_{\text{Cullis,ano}} < 1.0$  suggests that a contribution to the phasing from the anomalous data is likely (*MLPHARE* program documentation; Collaborative Computational Project, No. 4, 1994).

$$R_{\text{Cullis,ano}} = \frac{\sum_{hkl} ||F_{\text{PH}}(hkl) - F_{\text{PH}}(\bar{h}\bar{k}\bar{l})| - |2F_H \sin \alpha_P|}{\sum_{hkl} |F_{\text{PH}}(hkl) - F_{\text{PH}}(\bar{h}\bar{k}\bar{l})|}. \quad (2.2.5.2)$$

**Phasing power,  $\text{PP}_{\text{iso}}$ .** The isomorphous phasing power  $\text{PP}_{\text{iso}}$  for phase determination by isomorphous replacement is defined

for a particular pair of native and heavy-atom-derivative data sets as the ratio of  $|F_H|$  and  $\varepsilon(\varphi_P)$ , where  $|F_H|$  is the calculated amplitude of the heavy-atom structure factor and  $\varepsilon(\varphi_P)$  is the lack-of-closure error [equation (2.2.5.1a)].  $F_P + F_H$  is a vector sum of the calculated structure factor  $F_H$  and the structure factor  $F_P$ .

$$\text{PP}_{\text{iso}} = \sum_{hkl} |F_H| / \sum_{hkl} |F_{\text{PH}} - |F_P + F_H||. \quad (2.2.5.3)$$

There is another, slightly different, definition of  $\text{PP}_{\text{iso}}$ , which is implemented in the program *SOLVE*. Here,  $\text{PP}_{\text{iso}}$  is defined as the ratio of the r.m.s. of the  $|F_H|$  values and the r.m.s. of the lack-of-closure errors  $\varepsilon(\varphi_P)$ . For each reflection  $hkl$ , a weighted average of  $\varepsilon(\varphi_P)$  is calculated by integrating  $\varepsilon(\varphi_P)$  over the whole phase circle.

$$\text{PP}_{\text{iso}} = (\sum_{hkl} |F_H|^2)^{1/2} / (\sum_{hkl} \langle |F_{\text{PH}} - |F_P + F_H|| \rangle^2)^{1/2}. \quad (2.2.5.4)$$

*Note.* Owing to the cancelling out of the factor  $N^{1/2}$  in the numerator and denominator, the latter  $\text{PP}_{\text{iso}}$  formula does not appear as a ratio of r.m.s. values at first sight.

**Anomalous phasing power,  $\text{PP}_{\text{ano}}$ .** The anomalous phasing power  $\text{PP}_{\text{ano}}$  for phase determination by anomalous-dispersion methods is defined as the ratio of the sum of calculated anomalous differences  $d''_{\text{calc}}$  and the sum of estimated standard uncertainties  $\sigma(d''_{\text{obs}})$  in the measurement of these anomalous differences:

$$\text{PP}_{\text{ano}} = \sum_{hkl} d''_{\text{calc}} / \sum_{hkl} \sigma(d''_{\text{obs}}). \quad (2.2.5.5)$$

As with  $\text{PP}_{\text{iso}}$  (see above), the program *SOLVE* uses a slightly different definition of  $\text{PP}_{\text{ano}}$ . Here, the anomalous phasing power is defined as the ratio of the r.m.s. of the  $d''_{\text{calc}}$  values and the r.m.s. of  $\sigma(d''_{\text{obs}})$ . For this, a weighted average of  $d''_{\text{calc}}$  is computed by integrating over the whole phase circle for each reflection.

$$\text{PP}_{\text{ano}} = (\sum_{hkl} \langle d''_{\text{calc}} \rangle^2)^{1/2} / [\sum_{hkl} \sigma(d''_{\text{obs}})^2]^{1/2}. \quad (2.2.5.6)$$

*Note.* As above in the  $\text{PP}_{\text{iso}}$  formula, the factors  $N^{1/2}$  in the numerator and denominator cancel out.

**Figure of merit (f.o.m.),  $m$ .** The figure of merit  $m$  is a term used in a number of contexts in X-ray crystallography. In its most common use, it is defined as the weight applied to an individual structure-factor amplitude that, in conjunction with its best phase, gives rise, in a Fourier synthesis, to the electron-density map with the minimum level of noise (Blow & Crick, 1959). Typically,  $m$  is given as an average value over all reflections in the data set or in a given resolution shell.

$$m = \int P(\alpha) \exp(i\alpha) d\alpha / \int P(\alpha) d\alpha = \langle \cos(\Delta\alpha) \rangle, \quad (2.2.5.7)$$

where  $P(\alpha)$  is the probability of the phase  $\alpha$ , initial or refined, being the best phase and  $\Delta\alpha = \alpha_{\text{best}} - \alpha$  is the error in the phase angle at  $\alpha$ . The integration is from 0 to  $2\pi$  and values for  $m$  range from 0 to 1.

### 2.2.6. Quality indicators for density modification and phase improvement

After determination of initial phases, a first electron-density map can be computed. It is expected that this map will contain significant errors and improbable features. Additional information, such as the flatness of the electron density in the solvent region or the similarity of electron-density regions of two or more identical molecules in the asymmetric unit, can be exploited to modify the electron density and hence improve the phases.

## 2. BASIC CRYSTALLOGRAPHY

**Figure of merit (f.o.m.),  $m$ .** The figure of merit  $m$  [see equation (2.2.5.7)] is also used to judge the degree of improvement in the phase values. Again, in its most common use, it is defined as the weight applied to an individual structure-factor amplitude that, in conjunction with its best phase, gives rise, in a Fourier synthesis, to the electron-density map with the minimum level of noise. Typically,  $m$  is given as an average value over all reflections in the data set or in a given resolution shell.

**Density-modification (DM) free  $R$  factor.** The DM free  $R$  factor is defined in the same way as the refinement free  $R$  factor  $R_{\text{free}}$  [see equation (2.2.8.1) and the following paragraph describing  $R_{\text{free}}$ ]. It describes the disagreement between the observed structure-factor amplitudes of a certain set of reflections and the calculated amplitudes after density modification. It is a weak indicator used in the program *DM* (Cowtan, 1999) and is helpful mainly in identifying the correct enantiomorph.

**Density-modification (DM) real-space free residuals.** The DM real-space free residuals are two numbers (real-space free residual for the solvent area and real-space free residual for the protein area) which are calculated by omitting two small spheres of protein and solvent from the density-modification process. The real-space free residual for the solvent area indicates how flat the solvent is in a non-flattened region of solvent and the real-space free residual for the protein area indicates how well the electron density in a non-histogram-matched region of protein obeys the expected histogram. The two residuals can also be combined into the combined real-space free residual by weighted summation, where the weight is defined by the solvent content. The DM real-space free residuals have some value in determining when to stop a density-modification calculation, once no more progress is being made, but limited value otherwise.

**Contrast,  $c$ .** The contrast  $c$  between the r.m.s. electron density in the solvent region and the r.m.s. electron density in the macromolecular region can serve as an indication of the presence of a clearly-defined solvent boundary in the electron-density map. It is defined as the standard deviation of the local r.m.s. electron-density values over the entire asymmetric unit (Terwilliger & Berendzen, 1999; Sheldrick, 2002; Terwilliger *et al.*, 2009). The electron-density values are first squared (optionally after truncating very high and very low density values) and then smoothed using a moving local sphere typically with a radius of 6 Å. Local r.m.s. values are then calculated. The contrast  $c$  is now defined as the standard deviation  $\sigma$  of the local r.m.s. density values multiplied by a factor to normalize for the fraction of solvent  $sc$  in the crystal:

$$c = [(1 - sc)/sc]^{1/2} \sigma. \quad (2.2.6.1)$$

**Skewness of electron density,  $S$ .** A high value for the skewness  $S$  of the electron density in an electron-density map indicates the presence of local electron-density maxima with high positive density values. The skewness is defined as the third moment of the electron density:

$$S = \langle \rho^3 \rangle / \langle \rho^2 \rangle^{3/2}. \quad (2.2.6.2)$$

In order to compute the mean values of  $\rho^3$  and  $\rho^2$ , all density grid points in the asymmetric unit of the electron-density map are taken into account (Terwilliger *et al.*, 2009).

**Overlap of NCS-related density,  $O_{\text{NCS}}$ .** The presence of correlated electron density at noncrystallographic symmetry (NCS)-related regions in a map can be used as an indicator for the quality of the electron-density map (Cowtan & Main, 1998; Vellieux *et al.*, 1995; Terwilliger *et al.*, 2009). The overlap ( $O_{\text{NCS}}$ ) between density values at NCS-related locations is also often

used to evaluate the presence of local symmetry or non-space-group symmetry:

$$O_{\text{NCS}} = \langle \rho_i \rho_j \rangle. \quad (2.2.6.3)$$

$\rho_i$  and  $\rho_j$  are the normalized electron-density values in the NCS-related regions in the asymmetric unit. The average is calculated over the whole region where NCS is present. This region may be defined as the region where overlap values are 0.3 or greater, or by a mask. If there are more than two NCS groups, the average is taken over all NCS pairs.

**$R$  factor,  $R_{\text{DENMOD}}$ , and phase correlation,  $m_{\text{DENMOD}}$ , from statistical density modification.** The amplitudes and phases of structure factors calculated using statistical density modification can be compared with the observed amplitudes and experimental phases (Cowtan & Main, 1996; Terwilliger, 2001; Terwilliger *et al.*, 2009). These comparisons yield an  $R$  value ( $R_{\text{DENMOD}}$ ) for the amplitudes and a mean cosine of the phase difference ( $m_{\text{DENMOD}}$ ) for the phases.

$$R_{\text{DENMOD}} = \sum_{hkl} \left| |F_{\text{obs}}| - |F_{\text{DENMOD}}| \right| / \sum_{hkl} |F_{\text{obs}}|, \quad (2.2.6.4)$$

$$m_{\text{DENMOD}} = (1/N) \sum_{hkl} \cos(\alpha_{\text{obs}} - \alpha_{\text{DENMOD}}). \quad (2.2.6.5)$$

### Correlation coefficient CC of chain trace against native data.

The quality of density modification in *SHELXE* (Sheldrick, 2002, 2010) can be assessed by computing a Pearson linear correlation coefficient [see equation (2.2.2.13)] of the calculated structure-factor amplitudes for a chain trace against the native structure-factor amplitudes. If the poly-Ala trace yields a CC value higher than 0.25 and if the mean chain length of the trace is >10, the solution is almost always correct. This criterion is currently implemented in the program *ARCIMBOLDO* (Rodríguez *et al.*, 2009).

### 2.2.7. Quality indicators for molecular replacement

In a case where a known structure is assumed to be similar to that of a target molecule (structural similarity is typically inferred by the degree of sequence similarity), the known structure, also termed the search model, can be used to determine the structure of the target molecule. The approach, termed molecular replacement, was first described in 1962 (Rossmann & Blow, 1962). Nowadays, about two thirds of all newly determined structures are determined by molecular replacement (Long *et al.*, 2008).

**Rotation function, RF.** The rotation function RF is a measure of the overlap or the agreement of the stationary Patterson function  $P2$  calculated from the observed data and the rotated Patterson function  $P1$  from the search model.

$$\text{RF} = \int_r P2 \underline{R} P1 dr. \quad (2.2.7.1)$$

In the equation for RF,  $\underline{R}$  is the rotation operator. The integration is performed between a minimum value and a maximum value for the radius  $r$ . These values are chosen according to the size of the search model, with the aim of including as many intramolecular Patterson peaks (self vectors) as possible and to exclude as many intermolecular Patterson peaks (cross vectors) as possible. The ratio of the height of a peak in the RF to the background level is used as an indicator of how likely it is that this peak describes the orientation of a molecule in the target structure.

**Translation function, TF.** There are numerous ways of defining a translation function TF, making it impractical to discuss quality

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

## 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.

## 2. BASIC CRYSTALLOGRAPHY

- Diederichs, K. & Karplus, P. A. (1997a). Improved R-factors for diffraction data analysis in macromolecular crystallography. *Nat. Struct. Biol.* **4**, 269–275.
- Diederichs, K. & Karplus, P. A. (1997b). Improved R-factors for diffraction data analysis in macromolecular crystallography. Erratum. *Nat. Struct. Biol.* **4**, 592.
- Engh, R. A. & Huber, R. (1991). Accurate bond and angle parameters for X-ray protein structure refinement. *Acta Cryst.* **A47**, 392–400.
- Evans, P. (2006). Scaling and assessment of data quality. *Acta Cryst.* **D62**, 72–82.
- Fisher, S. J., Helliwell, J. R., Khurshid, S., Govada, L., Redwood, C., Squire, J. M. & Chayen, N. E. (2008). An investigation into the protonation states of the C1 domain of cardiac myosin-binding protein C. *Acta Cryst.* **D64**, 658–664.
- Grosse-Kunstleve, R. W. & Adams, P. D. (2003). Substructure search procedures for macromolecular structures. *Acta Cryst.* **D59**, 1966–1973.
- Hauptman, H. & Karle, J. (1953). Solution of the phase problem. I. The centrosymmetric crystal. Am. Crystallogr. Assoc. Monograph No. 3. Dayton, Ohio: Polycrystal Book Service.
- Hirshfeld, F. L. & Rabinovich, D. (1973). Treating weak reflexions in least-squares calculations. *Acta Cryst.* **A29**, 510–513.
- Hoof, R. W. W., Vriend, G., Sander, C. & Abola, E. E. (1996). Errors in protein structures. *Nature (London)*, **381**, 272.
- Howell, P. L. & Smith, G. D. (1992). Identification of heavy-atom derivatives by normal probability methods. *J. Appl. Cryst.* **25**, 81–86.
- James, R. W. (1948). False detail in three-dimensional Fourier representations of crystal structures. *Acta Cryst.* **1**, 132–134.
- Jones, T. A., Zou, J.-Y., Cowan, S. W. & Kjeldgaard, M. (1991). Improved methods for building protein models in electron density maps and the location of errors in these models. *Acta Cryst.* **A47**, 110–119.
- Kabsch, W. (1988). Evaluation of single-crystal X-ray diffraction data from a position-sensitive detector. *J. Appl. Cryst.* **21**, 916–924.
- Kabsch, W. (1993). Automatic processing of rotation diffraction data from crystals of initially unknown symmetry and cell constants. *J. Appl. Cryst.* **26**, 795–800.
- Kabsch, W. (2010). XDS. *Acta Cryst.* **D66**, 125–132.
- Kleywegt, G. & Brünger, A. T. (1996). Checking your imagination: applications of the free R value. *Structure*, **4**, 897–904.
- Long, F., Vagin, A. A., Young, P. & Murshudov, G. N. (2008). BALBES: a molecular-replacement pipeline. *Acta Cryst.* **D64**, 125–132.
- Luzzati, V. (1952). Traitement statistique des erreurs dans la détermination des structures cristallines. *Acta Cryst.* **5**, 802–810.
- McCoy, A. J., Grosse-Kunstleve, R. W., Adams, P. D., Winn, M. D., Storoni, L. C. & Read, R. J. (2007). Phaser crystallographic software. *J. Appl. Cryst.* **40**, 658–674.
- Matthews, B. W. (1968). Solvent content of protein crystals. *J. Mol. Biol.* **33**, 491–497.
- Otwinowski, Z. & Minor, W. (1997). Processing of X-ray diffraction data collected in oscillation mode. *Methods Enzymol.* **276**, 307–326.
- Panjikar, S. & Tucker, P. A. (2002). Phasing possibilities using different wavelengths with a xenon derivative. *J. Appl. Cryst.* **35**, 261–266.
- Parkinson, G., Vojtechovsky, J., Clowney, L., Brünger, A. T. & Berman, H. M. (1996). New parameters for the refinement of nucleic acid-containing structures. *Acta Cryst.* **D52**, 57–64.
- Pflugrath, J. W. (1999). The finer things in X-ray diffraction data collection. *Acta Cryst.* **D55**, 1718–1725.
- Read, R. J. (1986). Improved Fourier coefficients for maps using phases from partial structures with errors. *Acta Cryst.* **A42**, 140–149.
- Read, R. J. (2001). Pushing the boundaries of molecular replacement with maximum likelihood. *Acta Cryst.* **D57**, 1373–1382.
- Rodgers, J. L. & Nicewander, W. A. (1988). Thirteen ways to look at the correlation coefficient. *Am. Stat.* **42**, 59–66.
- Rodriguez, D. D., Grosse, C., Himmel, S., Gonzalez, C., de Ilarduya, I. M., Becker, S., Sheldrick, G. M. & Uson, I. (2009). Crystallographic *ab initio* protein structure solution below atomic resolution. *Nat. Methods*, **6**, 651–653.
- Rossmann, M. G. & Blow, D. M. (1962). The detection of sub-units within the crystallographic asymmetric unit. *Acta Cryst.* **15**, 24–31.
- Schneider, T. R. & Sheldrick, G. M. (2002). Substructure solution with SHELXD. *Acta Cryst.* **D58**, 1772–1779.
- Sheldrick, G. M. (2002). Macromolecular phasing with SHELXE. *Z. Kristallogr.* **217**, 644–650.
- Sheldrick, G. M. (2008). A short history of SHELX. *Acta Cryst.* **A64**, 112–122.
- Sheldrick, G. M. (2010). Experimental phasing with SHELXC/D/E: combining chain tracing with density modification. *Acta Cryst.* **D66**, 479–485.
- Stenkamp, R. E. & Jensen, L. H. (1984). Resolution revisited: limit of detail in electron density maps. *Acta Cryst.* **A40**, 251–254.
- Terwilliger, T. C. (2001). Map-likelihood phasing. *Acta Cryst.* **D57**, 1763–1775.
- Terwilliger, T. C., Adams, P. D., Read, R. J., McCoy, A. J., Moriarty, N. W., Grosse-Kunstleve, R. W., Afonine, P. V., Zwart, P. H. & Hung, L.-W. (2009). Decision-making in structure solution using Bayesian estimates of map quality: the PHENIX AutoSol wizard. *Acta Cryst.* **D65**, 582–601.
- Terwilliger, T. C. & Berendzen, J. (1999). Discrimination of solvent from protein regions in native Fouriers as a means of evaluating heavy-atom solutions in the MIR and MAD methods. *Acta Cryst.* **D55**, 501–505.
- Vaguine, A. A., Richelle, J. & Wodak, S. J. (1999). SFCHECK: a unified set of procedures for evaluating the quality of macromolecular structure-factor data and their agreement with the atomic model. *Acta Cryst.* **D55**, 191–205.
- Vellieux, F. M. D. A. P., Hunt, J. F., Roy, S. & Read, R. J. (1995). DEMON/ANGEL: a suite of programs to carry out density modification. *J. Appl. Cryst.* **28**, 347–351.
- Weeks, C. M., De Titta, G. T., Miller, R. & Hauptman, H. A. (1993). Application of the minimal principle to peptide structures. *Acta Cryst.* **D49**, 179–181.
- Weeks, C. M., DeTitta, G. T., Hauptman, H. A., Thuman, P. & Miller, R. (1994). Structure solution by minimal-function phase refinement and Fourier filtering. II. Implementation and applications. *Acta Cryst.* **A50**, 210–220.
- Weiss, M. S. (2001). Global indicators of X-ray data quality. *J. Appl. Cryst.* **34**, 130–135.
- Weiss, M. S. & Hilgenfeld, R. (1997). On the use of the merging R factor as a quality indicator for X-ray data. *J. Appl. Cryst.* **30**, 203–205.
- Wilson, A. J. C. (1949). The probability distribution of X-ray intensities. *Acta Cryst.* **2**, 318–321.