11. DATA PROCESSING

## 11.5.5. Generalization of the procedure for averaging reflection intensities

Once the scale factors of all frames are determined, they need to be applied to the reflection intensities and error estimates. The reflection intensities with the same reduced Miller indices can then be averaged.

When method 2 is used for averaging, the determination of $\langle I_h \rangle$ is more complicated as there are as many estimates of the full intensity $I_{hi}$ as there are partial reflections $h_{im}$. Therefore, intensity averaging of reflection $h$ has to be done in two steps. First, for every reflection $h_i$, the intensity estimates from all partial observations will be the weighted mean, where the weights are based on the estimated standard deviations of each intensity measurement. In the second step, the average is taken over the $i$ different scaled intensities for the observed reflections.

The selection of reflections useful for averaging is the same as for scaling (Table 11.5.3.1), except that it is no longer necessary to reject reflections that have insignificant intensities. Applying a $\sigma$ cutoff while averaging the scaled intensities will lead to a statistical bias of the weaker reflection intensities.

For samples of three or more equivalent reflections, it is necessary to consider the absolute values of the differences between individual intensities and the median of the sample: $|I_{hi} - I_{\text{median}}|$. The outliers can be detected by several statistical tests and, once detected, can be either down-weighted or rejected. When the sample consists of only two reflections, they can be considered a 'discordant pair' if the difference between their intensities is not warranted by the estimated errors and, hence, both reflections can be rejected (Blessing, 1997).

Averaging intensities estimated according to method 2 has an advantage over method 1 as outliers and discordant pairs can be 'screened' at two levels: firstly, when the estimates of the full reflection intensity $I_{hi}$, calculated by expression (11.5.2.2) from different parts of the same reflection, are considered, and secondly when the mean intensities $\langle I_{hi} \rangle$ from different reflections are considered.

## 11.5.6. Estimating the quality of data scaling and averaging

A commonly used estimate of the quality of scaled and averaged Bragg reflection intensities is $R_{\text{merge}}$. Useful definitions of $R$ factors are:

$$R_{\text{merge}} = R_1 = \left[ \left( \sum_h \sum_i |I_{hi} - \langle I_h \rangle| \right) \Big/ \sum_h \sum_i |I_{hi}| \right] \times 100\%,$$
(11.5.6.1)

$$R_2 = \left\{ \left[ \sum_h \sum_i (I_{hi} - \langle I_h \rangle)^2 \right] \Big/ \sum_h \sum_i I_{hi}^2 \right\} \times 100\%$$
(11.5.6.2)

$$\text{and} \quad R_w = \left\{ \left[ \sum_h \sum_i W_{hi}(I_{hi} - \langle I_h \rangle)^2 \right] \Big/ \sum_h \sum_i W_{hi} I_{hi}^2 \right\} \times 100\%.$$
(11.5.6.3)

The linear ($R_1$), square ($R_2$) and weighted ($R_w$) $R$ factors can be subdivided into resolution ranges, intensity ranges, reflection classes, frame number and regions of the detector surface. When method 1 is used, reflections $h_i$ can be grouped in terms of the sums of partialities of contributing partial reflections $h_{im}$.

The $R$-factor variation depends on the properties of the detector with respect to intensities. Generally the $R$ factor decreases as intensity increases. Thus, the $R$ factor generally increases with resolution. Any deviation from this behaviour might indicate a problem in the data collection due to nonlinearity of the detector response, ice diffuse diffraction, or any other stray effects superimposed on the crystal diffraction.

A useful indicator of the quality of the intensity estimates of partial reflections is the mean ratio of calculated partiality to observed partiality:

$$r_p = \langle p_{him}^{\text{calc}} / p_{him}^{\text{obs}} \rangle = \langle p_{him}^{\text{calc}} \langle I_h \rangle / I_{him} \rangle.$$
(11.5.6.4)

The deviation of this ratio from unity can be examined as a function of the reflection intensity, resolution and calculated partiality.

The comparison of $R$ factors for centric and noncentric reflections can be used to determine the significance of an anomalous-scattering effect. The quality of the anomalous-dispersion signal can be assessed by calculation of the scatter, $\sigma_{Ih}$, where

$$\sigma_{Ih} = \left\{ [1/(n-1)] \sum_n (\langle I_h \rangle - I_{hn})^2 \right\}^{1/2}$$
(11.5.6.5)

and $\langle I_h \rangle$ is the average of the $n$ measurements of the full reflection intensities $I_{hn}$. The $\sigma_{Ih}$ values for noncentric reflections can be compared to the scatter, $\sigma_{Ih}^+$ or $\sigma_{Ih}^-$, of reflections differing only in absorption while excluding Bijvoet opposites. The mean scatter is calculated from all $\sigma_{Ih}$ values,

$$\langle \sigma_{Ih} \rangle = (1/h) \sum_h \left\{ [1/(n-1)] \sum_n (\langle I_h \rangle - I_{hn})^2 \right\}^{1/2}.$$
(11.5.6.6)

The ratios $\langle \sigma_{Ih} \rangle / \langle \sigma_{Ih}^+ \rangle$ and $\langle \sigma_{Ih} \rangle / \langle \sigma_{Ih}^- \rangle$ should be larger than unity for significant anomalous-dispersion data.

## 11.5.7. Experimental results

11.5.7.1. *Variation of scale factors versus frame number*

If scale factors are to make physical sense, their behaviour with respect to the frame number has to be in accordance with the known changes in the beam intensity, crystal condition and detector response.

The scaling of a $\varphi$X174 procapsid data set (Dokland *et al.*, 1997) was performed using methods 1 and 2 as described here and using *SCALEPACK* (Otwinowski & Minor, 1997) (Fig. 11.5.7.1). Graphs (*a*) and (*b*) in Fig. 11.5.7.1 have four segments corresponding to four synchrotron beam 'fills'. All three methods give scale factors within 5% of each other (Figs. 11.5.7.1*c* and *d*). However, for the first and last frame of each 'fill' the results can differ by as much as 15%. Both method 1 and *SCALEPACK* produce physically wrong results in that the scale factors of these frames look like outliers compared to the scale factors of the neighbouring frames. By contrast, method 2 provides consistent scale factors for these frames. Although the algorithm used by *SCALEPACK* for scaling frames with partial reflections has never been disclosed, the similar results obtained by method 1 and *SCALEPACK* suggest that *SCALEPACK* might be using an algorithm similar to that of method 1 (Fig. 11.5.7.1*d*).

Attempts at scaling a data set of a frozen crystal of HRV14 (Rossmann *et al.*, 1985, 1997) failed with method 1 as a result of gaps in the rotation range for the first 20 frames, causing singularity of the normal equations matrix. When frames without useful neighbours were excluded, the cubic symmetry of the crystal was sufficient for successful scaling. In contrast, method 2 did not have any problems with the whole data set, and the results obtained with method 2 showed greater consistency than those obtained with method 1 or *SCALEPACK* (Fig. 11.5.7.2).
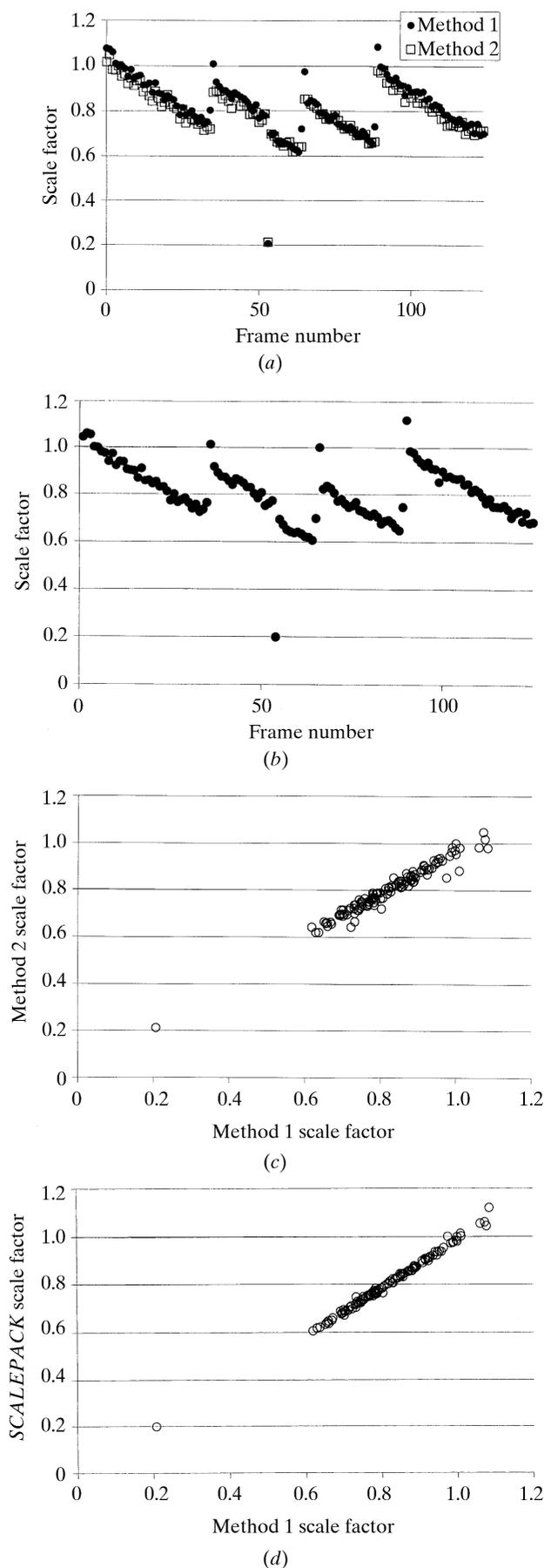
238

(a)



(b)



(c)



(d)

Fig. 11.5.7.1. Linear scale factors as a function of frame number for a $\varphi$X174 data set (Dokland *et al.*, 1997). Results from (a) method 1 and method 2, (b) *SCALEPACK*. Comparison of (c) method 2 *versus* method 1, and (d) *SCALEPACK versus* method 1.
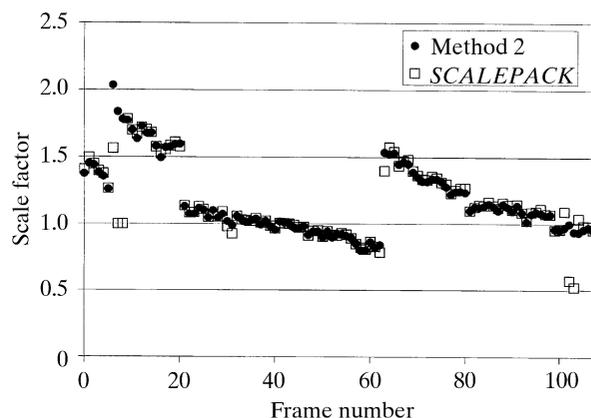


Fig. 11.5.7.2. Linear scale factor as a function of frame number for an HRV14 data set (Rossmann *et al.*, 1985, 1997).

The accuracy and robustness of method 2 is also demonstrated by the scaling results for a Sindbis virus capsid protein (SCP), residues 114–264 (Choi *et al.*, 1991, 1996). The behaviour of the scale factor with respect to the frame number reflects the anisotropy of the thin plate-shaped crystal (Fig. 11.5.7.3). For the first 40 frames (frame numbers 0 to 39), even-numbered frames have higher scale factors than odd-numbered frames. Data collection was stopped after frame number 39 and restarted. After frame number 39, odd-numbered frames have higher scale factors than even-numbered frames. This effect presumably relates to the use of the two alternating image plates with slightly different sensitivities in the *R*-axis camera used in the data collection.

### 11.5.7.2. *R factor as a function of 'sum-of-partialities' (method 1)*

In order to determine the limits of tolerance that can be permitted when method 1 is used, the *R* factor was examined as a function of the sum-of-partialities for the $\varphi$X174 procapsid data (Fig. 11.5.7.4). Reflections with sum-of-partialities of $1 \pm 0.3$ were used. The *R* factor changes sharply when the sum-of-partialities is outside $1 \pm 0.15$. Hence, $\pm 0.15$ were acceptable limits of tolerance for this data set.
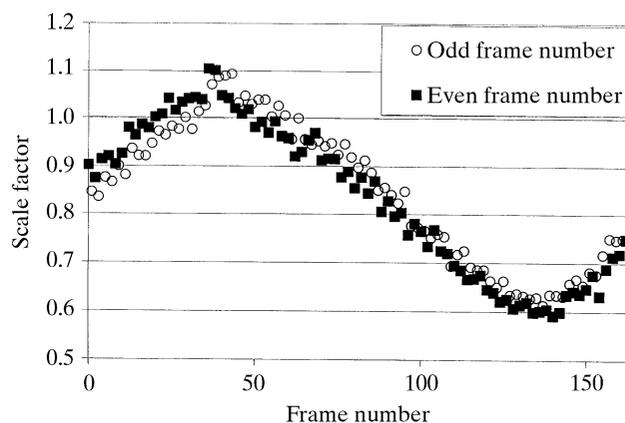


Fig. 11.5.7.3. Linear scale factor determined by method 2 as a function of frame number for an SCP(114–264) data set (Choi *et al.*, 1991, 1996). The sine-like pattern reflects the anisotropy of a thin plate-shaped crystal.
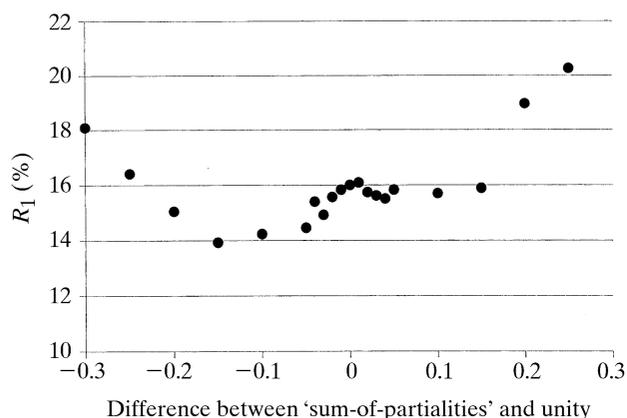
239

Fig. 11.5.7.4. *R* factor as a function of the difference of calculated 'sum-of-partialities' and unity for the estimates of full reflections when method 1 is used for the scaling and averaging of a $\varphi$X174 data set (Dokland *et al.*, 1997).



Fig. 11.5.7.6. The observed partialities plotted against calculated partialities for a $\varphi$X174 data set (Dokland *et al.*, 1997) processed by method 2. The observed partialities for individual partial reflections were averaged in bins of calculated partialities. The broken line represents the ideal relationship $p_{obs} = p_{calc}$.

### 11.5.7.3. *Statistics for rejecting reflections and data quality as a function of frame number*

The behaviour of the *R* factor *versus* frame number (Fig. 11.5.7.5) is more monotonic when method 1 is used compared to method 2. In method 1, the data-quality estimates for neighbouring frames are strongly correlated because the full reflections used in the statistics are obtained by summing partials from consecutive frames. By contrast, in method 2 every frame produces estimates of full reflection intensities independently of the neighbouring frames. Therefore, the *R* factors per frame calculated after scaling with method 2 truly represent the data quality for individual frames.

### 11.5.7.4. *Observed versus calculated partiality*

The relationship between observed and calculated partialities (Fig. 11.5.7.6) deviates from the ideal line $p_{obs} = p_{calc}$, especially for the smaller calculated partialities where $p_{obs} > p_{calc}$. This suggests errors in the measurements of $p_{obs}$ or the calculations of $p_{calc}$. The latter may be improved by a post refinement of the orientation matrix and crystal mosaicity (Rossmann *et al.*, 1979).

### 11.5.7.5. *Anisotropic mosaicity*

Refinement of the effective mosaicity can show both the anisotropic nature of the crystal (Fig. 11.5.7.7) as well as the impact of radiation damage. The effective mosaicity is the convolution of
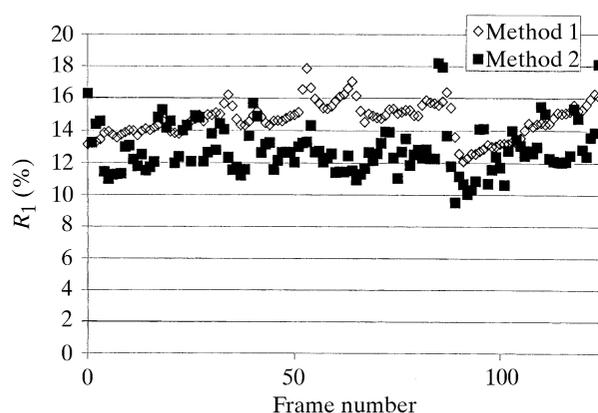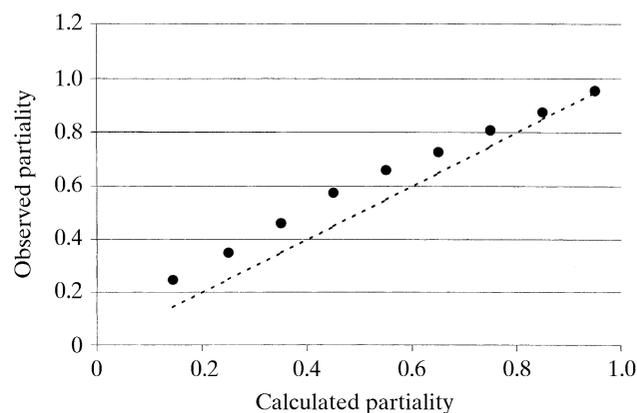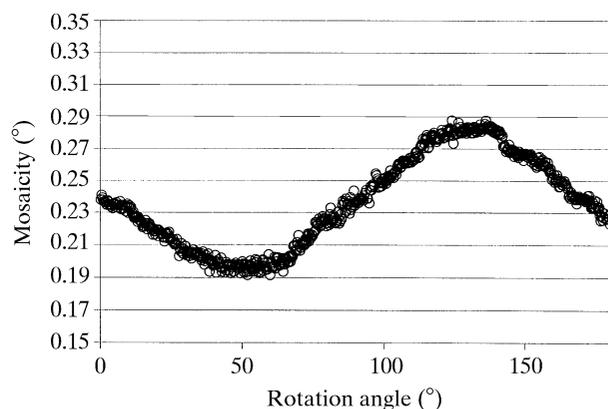


Fig. 11.5.7.7. Variation of (unconstrained) mosaicity for a monoclinic crystal of the bacterial virus alpha3 (Bernal *et al.*, 1998) showing the crystal anisotropy.

the mosaic spread of the crystal, the beam divergence and the wavelength divergence of the incident X-ray beam. Hence, X-ray diffraction data collected at a synchrotron-radiation source necessitate the differentiation of the effective mosaicity in the horizontal and vertical planes. A more general approach is the introduction of six parameters reflecting the anisotropic effective mosaicity.
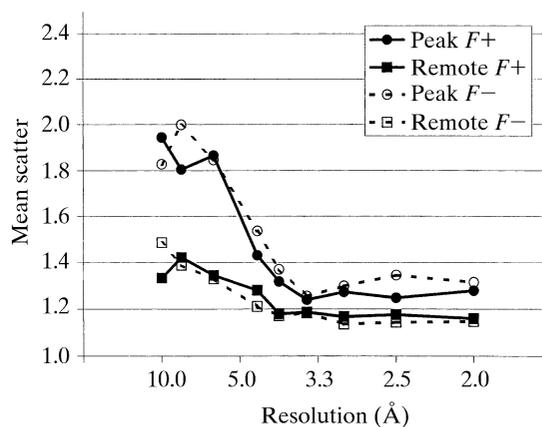


Fig. 11.5.7.5. *R* factor per frame as a function of frame number for a $\varphi$X174 data set (Dokland *et al.*, 1997).



Fig. 11.5.7.8. Quality of anomalous-dispersion data for an SeMet derivative of dioxygenase Rieske ferredoxin (Colbert & Bolin, 1999).

240

### 11.5.7.6. *Anomalous dispersion*

The quality of anomalous-dispersion data can be assessed by calculation of the average scatter, expression (11.5.6.6). The ratios $\langle\sigma_{Ih}\rangle/\langle\sigma_{Ih}^+\rangle$ and $\langle\sigma_{Ih}\rangle/\langle\sigma_{Ih}^-\rangle$ should be larger than unity for significant anomalous data (Fig. 11.5.7.8). Note the much larger ratios for the scatter among measurements of $I_h$ for data measured at the absorption edge of Se, as opposed to measurements remote from the edge. The decreasing values of the ratios with resolution are due to the decrease of $I_h$ value, thus causing the error in the measurement of $I_h$ to approach the difference in intensity of Bijvoet opposites.

### 11.5.8. Conclusions

The generalized HRS method allows scaling and averaging of X-ray diffraction data collected with an oscillation camera while simultaneously using full and partial reflections. The procedure is as useful for thin slices of reciprocal space as it is for thicker slices.

The results of data processing with the two different algorithms indicate that method 1, based on adding partial reflections, may fail to scale data sets with gaps in the rotation range or with low redundancy. The values of the scale factors obtained with both methods are similar, except for cases where there are gaps in the rotation range or dramatic changes in the true scale factors between consecutive frames. In these cases, method 1 produces a physically wrong result. The algorithm used by method 1 is probably similar to that used by *SCALEPACK* (Otwinowski & Minor, 1997).

Method 2 is more stable and versatile than method 1, and allows the scaling of data sets with incompletely measured reflections and low redundancy. The major drawback of method 2 is that errors in the crystal orientation matrix and mosaicity, as well as inadequacies of the theoretical model for reflection partiality, contribute to errors in the scaled intensities. Therefore, post refinement is needed for method 2 to perform at its best.

### Appendix 11.5.1.

### Partiality model (Rossmann, 1979; Rossmann *et al.*, 1979)

Small differences in the orientation of domains within the crystal, as well as the cross fire of the incident X-ray beam, will give rise to a series of possible Ewald spheres. Their extreme positions will subtend an angle $2m$ at the origin of the reciprocal space, and their centres lie on a cusp of limiting radius $\delta = m/\lambda$, where $m$ is the half-angle effective mosaic spread. As the reciprocal lattice is rotated around the axis $(Oy)$ perpendicular to the mean direction of the incident radiation $(Oz)$, a point $P$ will gradually penetrate the effective thickness of the reflection sphere (Fig. A11.5.1.1). Initially, only a few domain blocks will satisfy Bragg's law, but upon further rotation the number of blocks that are in a reflecting condition will increase. The maximum will be reached when the point $P$ has penetrated halfway through the sphere's effective thickness, after which there will be a decline of the crystal volume able to diffract.

Let $q$ be a measure of the fraction of the path travelled by $P$ between the extreme reflecting positions $P_A$ and $P_B$, and let $p$ be the fraction of the energy already diffracted. Then the relation between $p$ and $q$ must have the general form shown in Fig. A11.5.1.2. It is physically reasonable to assume that the curve for $p$ is tangential to $q = 0$ at $p = 0$ and to $q = 1$ at $p = 1$.

A reasonable approximation to the above conditions can be obtained by considering the fraction of the volume of a sphere removed by a plane a distance $q$ from its surface (Fig. A11.5.1.2). It is easily shown that if $p$ is the volume, then
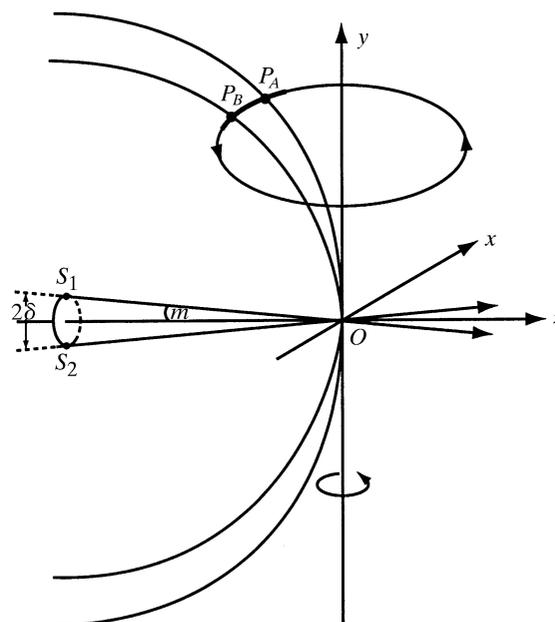


Fig. A11.5.1.1. Penetration of a reciprocal-lattice point $P$ into the sphere of reflection by rotation around $Oy$. The extremes of reflecting conditions at $P_A$ and $P_B$ are equivalent to X-rays passing along the lines $S_1O$ and $S_2O$ with centres of the Ewald spheres at $S_1$ and $S_2$ and subtending an angle of $2m$ at $O$. Hence, in three dimensions, the extreme reflecting spheres will lie with their centres on a circle of radius $\delta = m/\lambda$ at $z = -1/\lambda$.

$$p = 3q^2 - 2q^3. \qquad (A11.5.1.1)$$

This curve is shown in Fig. A11.5.1.2 and corresponds to assuming that the reciprocal-lattice point is a sphere of finite volume cutting an infinitely thin Ewald sphere. Also shown in Fig. A11.5.1.2 is the line $p = q$ which would result if the reciprocal-lattice point were a rectangular block whose surfaces were parallel and perpendicular to the Ewald sphere at the point of penetration.

Assuming a right-handed coordinate system $(x, y, z)$ in reciprocal space fixed to the camera, it is easily shown (Wonacott, 1977) that the condition for reflection is

$$d^{*2} + (2z/\lambda) = 0, \qquad (A11.5.1.2)$$

where $d^*$ is the distance of a reciprocal-lattice point $P(x, y, z)$ from the origin, $O$, of reciprocal space. Similarly, it can be shown that at the ends of the path of the reciprocal-lattice point through the finite thickness of the sphere,

$$d^{*2} + \delta^2 + (2z/\lambda) - 2\delta(x_A^2 + y_A^2)^{1/2} = 0 \quad \text{and}$$
$$d^{*2} + \delta^2 + (2z/\lambda) - 2\delta(x_B^2 + y_B^2)^{1/2} = 0. \qquad (A11.5.1.3)$$

Therefore,

$$z_A = (\lambda/2)\left[-d^{*2} - \delta^2 + 2\delta(x_A^2 + y_A^2)^{1/2}\right],$$
$$z_B = (\lambda/2)\left[-d^{*2} - \delta^2 + 2\delta(x_B^2 + y_B^2)^{1/2}\right]. \qquad (A11.5.1.4)$$

Since $\delta$ is small, it can be assumed that $2\delta(x^2 + y^2)^{1/2}$ is independent of the position of the reciprocal-lattice point $P$ between the extreme positions $P_A$ and $P_B$ (Fig. A11.5.1.1). Hence, the length of the path through the finite thickness of the sphere is proportional to

$$z_A - z_B = 2\lambda\delta(x_P^2 - y_P^2)^{1/2}. \qquad (A11.5.1.5)$$

references