## 11. DATA PROCESSING

from the expected distribution $\omega(\varepsilon_1, \varepsilon_2, \varepsilon_3)$ defined in Section 11.3.2.3 and removed from the final background value.

### 11.3.3.3. *Standard profiles*

Reflection profiles are represented on the Ewald sphere within a domain $D_0$ comprising $2n_1 + 1, 2n_2 + 1, 2n_3 + 1$ equidistant grid-points along $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$, respectively. The sampling distances between adjacent grid points are then $\Delta_1 = \delta_D/(2n_1 + 1)$, $\Delta_2 = \delta_D/(2n_2 + 1)$, $\Delta_3 = \delta_M/(2n_3 + 1)$. Thus, grid coordinate $\nu_3$ ($\nu_3 = -n_3, \ldots, n_3$) covers the set of rotation angles

$$\Gamma_{\nu_3} = \{\varphi'|(\nu_3 - 1/2)\Delta_3 \le (\varphi' - \varphi) \cdot \zeta \le (\nu_3 + 1/2)\Delta_3\}.$$

Contributions to the spot intensity come from one or several adjacent data images ($j = j_1, \ldots, j_2$), each covering the set of rotation angles

$$\Gamma_j = \{\varphi'|\varphi_0 + (j-1)\Delta_\varphi \le \varphi' \le \varphi_0 + j\Delta_\varphi\}.$$

Assuming Gaussian profiles along $\mathbf{e}_3$ for all reflections (see Section 11.3.2.3), the fraction of counts (after subtraction of the background) contributed by data frame $j$ to grid coordinate $\nu_3$ is

$$f_{\nu_3 j} \approx \int\limits_{\Gamma_j \cap \Gamma_{\nu_3}} \exp[-(\varphi' - \varphi)^2/2\sigma^2]\, \mathrm{d}\varphi'$$

$$\times \left\{\int\limits_{\Gamma_j} \exp[-(\varphi' - \varphi)^2/2\sigma^2]\, \mathrm{d}\varphi'\right\}^{-1},$$

where $\sigma = \sigma_M/|\zeta|$. The integrals can be expressed in terms of the error function, for which efficient numerical approximations are available (Abramowitz & Stegun, 1972). Finally, each pixel on data image $j$ belonging to the reflection is subdivided into $5 \times 5$ areas of equal size, and $f_{\nu_3 j}/25$ of the pixel signal is added to the profile value at grid coordinates $\nu_1, \nu_2, \nu_3$ corresponding to each subdivision.

This complicated procedure leads to more uniform intensity profiles for all reflections than using their untransformed shape. This simplifies the task of modelling the expected intensity distribution needed for integration by profile fitting. As implemented in *XDS*, reference profiles are learnt every 5° of crystal rotation at nine positions on the detector, each covering an equal area of the detector face. In the learning phase, profile boxes of the strong reflections are normalized and added to their nearest reference profile boxes. The contributions are weighted according to the distance from the location of the reference profile. Each grid point within the average profile boxes is classified as signal if it is above 2% of the peak maximum. Finally, each profile is scaled such that the sum of its signal pixels normalizes to one. The analytic expression $\omega(\varepsilon_1, \varepsilon_2, \varepsilon_3)$ defined in Section 11.3.2.3 for the expected intensity distribution is only a rough initial approximation which is now replaced by the empirical reference profiles.

### 11.3.3.4. *Intensity estimation*

If an expected intensity distribution $\{p_i|i \in D_0\}$ of the observed profile is given in a domain $D_0$, the reflection intensity $I$ can be estimated as

$$I = \sum_{i \in D}(c_i - b_i)p_i/v_i \Big/ \sum_{i \in D}p_i^2/v_i,$$

which minimizes the function

$$\psi(I) = \sum_{i \in D}(c_i - I \cdot p_i - b_i)^2/v_i, \qquad \sum_{i \in D_0}p_i = 1.$$

$b_i, c_i, v_i$ ($i \in D$) are background, contents and variance of pixels observed in a subdomain $D \subseteq D_0$ of the expected distribution. The

background $b_i$ underneath a diffraction spot is often assumed to be a constant which is estimated from the neighbourhood around the reflection. Determination of reflection intensities by profile fitting has a long tradition (Diamond, 1969; Ford, 1974; Kabsch, 1988*b*; Otwinowski, 1993). Implementations of the method differ mainly in their assumptions about the variances $v_i$. Ford uses constant variances, which works well for films, which have a high intrinsic background. In *XDS*, which was originally designed for a multiwire detector, $v_i \propto p_i$ was assumed, which results in a straight summation of background-subtracted counts within the expected profile region, $I = \sum_{i \in D}(c_i - b_i)/\sum_{i \in D}p_i$. This particular simple formula is very satisfactory for the low background typical of these detectors. For the general case, however, better results can be obtained by using $v_i = b_i + Ip_i$ for the pixel variances as shown by Otwinowski and implemented in *DENZO* and in the later version of *XDS*. Starting with $v_i = b_i$, the intensity is now found by an iterative process which is terminated if the new intensity estimate becomes negative or does not change within a small tolerance, which is usually reached after three cycles. It can be shown that the solution thus obtained is unique.

### 11.3.4. Scaling

Usually, many statistically independent observations of symmetry-related reflections are recorded in the rotation images taken from one or several similar crystals of the same compound. The squared structure-factor amplitudes of equivalent reflections should be equal and the idea of scaling is to exploit this *a priori* knowledge to determine a correction factor for each observed intensity. These correction factors compensate to some extent for effects such as radiation damage, absorption, and variations in detector sensitivity and exposure times, as well as variations in size and disorder between different crystals.

The usual methods of scaling split the data into batches of roughly the same size, each covering one or more adjacent rotation images, and then determine a single scaling factor for all reflections in each batch. Neighbouring reflections may then receive quite different corrections if they are assigned to different batches. Since the selection of batch boundaries is to some extent arbitrary, a more continuous correction function would be preferable. This function could be modelled analytically (for example by using spherical harmonics) or empirically, as implemented in *XSCALE* and described below.

For each reflection, observational equations are defined as

$$\psi_{hl\alpha} = (I_{hl} - g_\alpha I_h)/\sigma_{hl}.$$

The subscript $h$ represents the unique reflection indices and $l$ enumerates all symmetry-related reflections to $h$. By definition, the unique reflection indices have the largest $h$, then $k$, then $l$ value occuring in the set of all indices related by symmetry to the original indices, including Friedel mates. Thus, two reflections are symmetry-related if and only if their unique indices are identical. $I_h$ is the unknown 'true' intensity and $I_{hl}, \sigma_{hl}$ are symmetry-related observed intensities and their standard deviations, respectively. The subscript $\alpha$ denotes the coordinates at which the scaling function $g_\alpha$ should be evaluated. As implemented in *XDS* and *XSCALE*, $\alpha = 1, \ldots, 9$ denotes nine positions uniformly distributed in the detector plane at the beginning of data collection, $\alpha = 10, \ldots, 18$ the same positions on the detector but after the crystal has been rotated by, say, 5°, and so on. The scaling factors $g_\alpha$ and the estimated intensities $I_h$ are found at the minimum of the function

$$\Psi = \sum_{hl\alpha} w_{hl\alpha} \Psi_{hl\alpha}^2.$$

222

The main difference from the method of Fox & Holmes (1966) is the introduction of the weights $w_{hl\alpha}$. These weights depend upon the distance between each reflection $hl$ and the positions $\alpha$. They are monotonically decreasing functions of this distance, implemented as Gaussians in *XDS* and *XSCALE*. This results in a smoothing of the scaling factors since each reflection contributes to the observational equations in proportion to the weights $w_{hl\alpha}$.

Minimization of $\Psi$ is done iteratively. After each step, the $g_\alpha$ are replaced by $\Delta g_\alpha + g_\alpha$ and rescaled to a mean value of 1. The corrections $\Delta g_\alpha$ are determined from the normal equations

$$\sum_\beta A_{\alpha\beta}\Delta g_\beta = b_\alpha,$$

where

$$A_{\alpha\beta} = \sum_h [\delta_{\alpha\beta}I_h^2 u_{h\alpha} + (r_{h\alpha}v_{h\beta} + v_{h\alpha}r_{h\beta} - v_{h\alpha}v_{h\beta})/u_h]$$

$$b_\alpha = \sum_h I_h r_{h\alpha}$$

$$I_h = v_h/u_h$$

$$r_{h\alpha} = v_{h\alpha} - g_\alpha u_{h\alpha} I_h$$

$$u_{h\alpha} = \sum_l w_{hl\alpha}/\sigma_{hl}^2$$

$$v_{h\alpha} = \sum_l w_{hl\alpha}I_{hl}/\sigma_{hl}^2$$

$$u_h = \sum_\alpha g_\alpha^2 u_{h\alpha}$$

$$v_h = \sum_\alpha g_\alpha v_{h\alpha}.$$

In case a 'true' intensity $I_h$ is available from a reference data set, the non-diagonal elements are omitted from the sum over $h$ in the normal matrix $A_{\alpha\beta}$. The corrections $\Delta g_\alpha$ are expanded in terms of the eigenvectors of the normal matrix, thereby avoiding shifts along eigenvectors with very small eigenvalues (Diamond, 1966). This filtering method is essential since the normal matrix has zero determinant if no reference data set is available.

### 11.3.5. Post refinement

The number of fully recorded reflections on each single image rapidly declines for small oscillation ranges and the complete intensities of the partially recorded reflections have to be estimated. This presented a serious obstacle in early structural work on virus crystals, as the crystal had to be replaced after each exposure on account of radiation damage. A solution of this problem, the 'post refinement' technique, was found by Schutt, Winkler and Harrison, and variants of this powerful method have been incorporated into most data-reduction programs [for a detailed discussion, see Harrison *et al.* (1985); Rossmann (1985)]. The method derives complete intensities of reflections only partially recorded on an image from accurate estimates for the fractions of observed intensity, the 'partiality'. The partiality of each reflection can always be calculated as a function of orientation, unit-cell metric, mosaic spread of the crystal and model intensity distributions. Obviously, the accuracy of the estimated full reflection intensity then strongly depends on a precise knowledge of the parameters describing the diffraction experiment. Usually, for many of the partial reflections, symmetry-related fully recorded ones can be found, and the list of such pairs of intensity observations can be used to refine the required parameters by a least-squares procedure. Clearly, this refinement is carried out after all images have been processed, which explains why the procedure is called 'post refinement'.

Adjustments of the diffraction parameters $s_\mu$ $(\mu = 1, \ldots, k)$ are determined by minimization of the function $E$, which is defined as the weighted sum of squared residuals between calculated and observed partial intensites.

$$E = \sum_{hj} w_{hj}(\Delta_{hj})^2$$

$$\Delta_{hj} = R_j(\varphi_{hj})g_j I_h - I_{hj}$$

$$w_{hj} = 1/\{\sigma^2(I_{hj}) + [R_j(\varphi_{hj})g_j]^2\sigma^2(I_h)\}.$$

Here, $I_{hj}$ is the intensity recorded on image $j$ of a partial reflection with indices summarized as $hj$, $I_h$ is the mean of the observed intensities of all fully recorded reflections symmetry-equivalent to $hj$, $g_j$ is the inverse scaling factor of image $j$, $\varphi_{hj}$ is the calculated spindle angle of reflection $hj$ at diffraction and $R_j$ is the computed fraction of total intensity recorded on image $j$.

Expansion of the residuals $\Delta_{hj}$ to first order in the parameter changes $\delta s_\mu$ and minimization of $E(\delta s_\mu)$ leads to the $k$ normal equations

$$\sum_{\mu'=1}^{k}\left(\sum_{hj}w_{hj}\frac{\partial\Delta_{hj}}{\partial s_\mu}\frac{\partial\Delta_{hj}}{\partial s_{\mu'}}\right)\delta s_{\mu'} = -\sum_{hj}w_{hj}\Delta_{hj}\frac{\partial\Delta_{hj}}{\partial s_\mu}.$$

Often, the normal matrix is ill-conditioned, since changes in some unit-cell parameters or small rotations of the crystal about the incident X-ray beam do not significantly affect the calculated partiality $R_j$. To take care of these difficulties, the system of equations is rescaled to yield unit diagonal elements for the normal matrix and the correction vector $\delta s_\mu$ is filtered by projection into a subspace defined by the eigenvectors of the normal matrix with sufficiently large eigenvalues (Diamond, 1966).

The parameters are corrected by the filtered $\delta s_\mu$ and a new cycle of refinement is started until a minimum of $E$ is reached. The weights, residuals and their gradients are calculated using the current values for $s_\mu$ and $g_j$ at the beginning of each cycle. The derivatives

$$\frac{\partial\Delta_{hj}}{\partial s_\mu} = g_j I_h\left(\frac{\partial R_j}{\partial\varphi_{hj}}\frac{\partial\varphi_{hj}}{\partial s_\mu} + \frac{\partial R_j}{\partial\sigma_M}\frac{\partial\sigma_M}{\partial s_\mu} + \frac{\partial R_j}{\partial|\zeta_{hj}|}\frac{\partial|\zeta_{hj}|}{\partial s_\mu}\right)$$

appearing in the normal equations can be worked out from the definitions given in Sections 11.3.2.2 and 11.3.2.4 (to simplify the following equations, the subscript $hj$ is omitted). The fraction $R_j$ of the total intensity can be expressed in terms of the error function (see Section 11.3.2.4) as

$$R_j = [\mathrm{erf}(z_1) - \mathrm{erf}(z_2)]/2$$

$$z_1 = |\zeta|(\varphi_0 + j\Delta_\varphi - \varphi)/(2)^{1/2}\sigma_M$$

$$z_2 = |\zeta|[\varphi_0 + (j-1)\Delta_\varphi - \varphi]/(2)^{1/2}\sigma_M.$$

Using the relation $\mathrm{d}\,\mathrm{erf}(z)/\mathrm{d}z = [2/(\pi)^{1/2}]\exp(-z^2)$, the derivatives of $R_j$ are

$$\partial R_j/\partial\varphi = [\exp(-z_2^2) - \exp(-z_1^2)]|\zeta|/[\sigma_M(2\pi)^{1/2}]$$

$$\partial R_j/\partial\sigma_M = [z_2\exp(-z_2^2) - z_1\exp(-z_1^2)]/[\sigma_M(\pi)^{1/2}]$$

$$\partial R_j/\partial|\zeta| = [z_1\exp(-z_1^2) - z_2\exp(-z_2^2)]/[|\zeta|(\pi)^{1/2}].$$

It remains to work out the derivatives $\partial\varphi/\partial s_\mu$, $\partial\sigma_M/\partial s_\mu$ and $\partial|\zeta|/\partial s_\mu$ (not shown here). As discussed in detail by Greenhough & Helliwell (1982), spectral dispersion and asymmetric beam cross fire lead to some variation of $\sigma_M$, which makes it necessary to include additional parameters in the list $s_\mu$. The effect of these parameters on the partiality is dealt with easily by the derivatives $\partial\sigma_M/\partial s_\mu$.