

11. DATA PROCESSING

for lune separation is distinct from the requirement for spot separation. If lunes overlap, spots may have more than one index consistent with a particular position on the detector.

The autoindexing procedure described above is not dependent on prior knowledge of the crystal unit cell; however, for efficiency reasons, the search is restricted to a reasonable range of unit-cell dimensions, obtained, for example, from the requirement of spot separation. In *DENZO*, this default can be overridden by the keyword 'longest vector', but the need to use this keyword is a sign of a problem that should be fixed. Either the defined spot size should be decreased or data should be recollected with the detector further away from the crystal.

11.4.3.4. *Misindexing*

Autoindexing is sensitive to inaccuracy in the description of the detector geometry. The specified position of the beam on the detector should correspond to the origin of the Bragg-peaks lattice (Miller index 000). Autoindexing will shift the origin of the lattice to the nearest Bragg lattice point. An incorrect beam position will result in the nearest Bragg lattice point not having the index 000. In such a situation, all reflections will have incorrectly determined indices. Such misindexing can be totally self-consistent until the intensities of symmetry-related reflections are compared. This dependence of the indexing correctness on the assumed beam position is the main source of difficulties in indexing (Gewirth, 1996; Otwinowski & Minor, 1997). The beam position has to be precise, as the largest acceptable error is one half of the shortest distance between spots.

Indexing limited to determining h , k , l triplets is not very sensitive to other detector parameters. Errors by a degree or two in rotation or by 10% in distance are unlikely to produce wrong values of h , k and l . Sometimes even a very large error, such as the distance being too large by a factor of 5, will still produce the correct h , k , l triplets. The detector position error will be compensated by an error in the lattice determined by autoindexing. For this reason, the accuracy of the lattice is not a function of the autoindexing procedure, but depends mainly on the accuracy of the detector description. By the same token, the distortion of the lattice also depends on the accuracy of the detector parameters.

11.4.3.5. *Twins*

Special care has to be taken if more than one crystal contributes to the diffraction image. When there is a large disproportion between volumes (e.g. the presence of a satellite crystal), autoindexing may work without any modifications. In the case of similar volumes, the manual editing of weaker reflections and resolution cuts can make the proportion of reflections from one crystal in the peak-search list large enough for the autoindexing method to succeed. If the crystals have a similar orientation, using only very low resolution data may be the right method. In the case of twinned crystals, autoindexing sometimes finds a superlattice that results in integer indices simultaneously for both crystals. In such a case, *DENZO* solves the problem of finding the best three-dimensional lattice that incorporates all of the observed peaks. Unfortunately, for a twinned crystal, this is a mathematically correct solution to an incorrectly posed problem.

11.4.4. Coordinate systems

There are four natural coordinate systems used to describe a diffraction experiment, defined by the order in which the data are stored, the beam and gravity, or the beam and the goniostat axes

(spindle or 2θ). These coordinate systems will be called, respectively, *data*, *beam-gravity*, *beam-spindle* and *beam- 2θ* .

11.4.4.1. *Beam-gravity*

To visualize a diffraction pattern, beam-gravity is the coordinate system clearly preferred by human physiology. The universal preference to relate to the gravity direction is revealed by the observation that people generally perceive an image in a mirror as inverted left-right rather than top-down. Hence *XdisplayF* uses the beam-gravity coordinate system, except when diffraction data cannot be related to gravity.*

11.4.4.2. *Data*

The first (1983) *DENZO* implementation used the data coordinate system to describe the beam position on the detector and to define the integration box. This is still the case in order to keep backward compatibility.

11.4.4.3. *Beam-spindle*

Until 1998, *DENZO* supported only a single-axis goniostat and used a beam-spindle coordinate system to define crystal and detector orientation and polarization. Initially, the goniostat spindle axis was assumed to be horizontal, so the direction perpendicular to the beam and spindle was described by the keyword *vertical*, which in reality may not relate to the gravity direction for some goniostats. The keyword *rotx* relates to rotation around the spindle axis, *roty* around the *vertical axis* and *rotz* around the beam axis. The definition of the orientation matrix in the communication file between *DENZO* and *SCALEPACK* uses an unintuitive convention: the letter y in *roty* relates to the first element of the vector, x in *rotx* relates to the second and z in *rotz* to the third. However, the matrix always has a positive determinant, so this convention has no impact on the handedness of the coordinate system. This unfortunate choice of convention, preserved for backward compatibility reasons, appears only in the communication file and has no significance for anybody who does not inspect the matrix.

11.4.4.4. *Beam- 2θ*

The recent addition of a general goniostat introduced a conceptual change in the *DENZO* coordinate system. The data-collection axis can be oriented in any direction, so in principle *rotx*, *roty* and *rotz* no longer need to be defined relative to the data-collection axis. However, to keep the useful correlation between refinable parameters (*crystal rotz* and *detector rotz* being close to 100% correlated), one real and two virtual goniostats are used simultaneously in *DENZO*. Refinable crystal parameters (*crystal rotx*, *roty*, *rotz*) are still defined, as in the past, by the data-collection axis and the beam. This means that the directions of rotations defined by *fit crystal rotx*, *roty* and *rotz* do not rotate around the data-collection axis as the program advances from one image to another. This coordinate system changes with the change in direction of the data-collection axis. Crystal orientation is defined by three constant, perpendicular axes, which, in the current version, no longer have to be aligned with the physical crystal goniostat. However, the so-called *2 theta* rotation has a fixed axis, and, if it exists, it defines the *DENZO* coordinate system together with the beam axis. Thus the current coordinate system in *DENZO* should be called beam- 2θ . Fortunately for the user, the conversions between different coordinate systems are handled transparently. For example, the refined change in the crystal orientation is converted from the refined goniostat to the crystal-orientation goniostat. The

* There may occasionally be an exception to this when the experimental system is not known.

movements of the physical goniostat are converted into appropriate changes in the diffraction pattern. The physical goniostat appears only to describe the data collection and, optionally, to calculate the physical goniostat angles that achieve particular crystal alignments.

The *DENZO* coordinate system (Gewirth, 1996) is used in the definition of crystal goniostats, 2θ goniostat, Weissenberg coupling and polarization.

This discussion of the coordinate systems shows that the conceptual complexity of the program description does not result in complexity of the actual use of the program. The success of data analysis does not require a full understanding of the relations between internal *DENZO* goniostats and the coordinate systems. The reason for this complexity was to create a simple pattern of correlations between crystal and detector parameters in *DENZO* refinement. This in turn allows for simple and easy-to-understand control of the refinement process and simplifies problem diagnostics. For example: the definition of refined *crystal rotx* as rotation around the data-collection axis makes hardware problems when driving the spindle and shutter result only in fluctuations of *crystal rotx*. The constant nonzero value of the refined shifts between frames of *crystal roty* and *rotz* is a sign of misalignment of the data-collection axis. Although the program compensates for this misalignment with changes in crystal orientation, this introduces a small error in the Lorentz factor. The nature of these problems is such that they do not result in a complete failure of the experiment, but they do have an impact on the quality of the result. It is up to the experimenter and the instrument manager to assess the significance of these indications.

11.4.5. Experimental assumptions

To achieve the main target of a diffraction experiment – the estimation of structure factors – three components need to be determined, with maximum possible precision:

(1) the crystal response function (the relationship between the crystal structure factor and the number of diffracted X-ray photons, which depends also on the X-ray source characteristics);

(2) the detector response function; and

(3) the geometrical description of the detector relative to the directions of the X-ray beam and crystal goniostat axes.

The main difficulty of data analysis in protein crystallography is the complexity of the process that determines these components. *HKL* can determine all three directly from the data produced by the analogue-to-digital converter (ADC). The only extra program needed is one that sends the raw ADC signal to the computer disk. For charge-coupled-device (CCD) detectors, spatial detector distortion and sensitivity per pixel functions need to be established in a separate experiment. Usually it is worthwhile to establish a geometrical description of the detector in a separate diffraction experiment. A precise determination requires a well diffracting, high symmetry, non-slipping crystal and a special data-collection procedure.

11.4.5.1. Crystal diffraction

The crystal response function consists of two types of factors included in the analysis: additive factors, which are represented by the background, and a number of multiplicative factors, such as exposed crystal volume, overall and resolution-dependent decay, Lorentz factor, flux variation, polarization, *etc.* Other factors, like extinction and non-decay radiation damage (radiation damage can result not only in decay, but also in a change in the crystal lattice, often a main source of error in an experiment), are ignored by *HKL*, except for their contribution to error estimates.

11.4.5.2. Data model

The detector response function is the main component for the data model. *HKL* supports

(1) data stored in 8 or 16 bit fields;

(2) overflow table;

(3) linear, bilinear, polynomial and exponential response, with the error model represented by an arbitrary scale;

(4) saturation limit;

(5) value representing lack of data;

(6) constant offsets per read-out channel;

(7) pattern noise;

(8) lossless compression;

(9) flood-field response; and

(10) sensitivity response.

HKL supports most data formats, which represent particular combinations of the above features. The formats define the coordinate system, the pixel size, the detector size, the active area and the fundamental shape (cylindrical, spherical, flat rectangular or circular, single or multi-module) of the detector.

The main complexity of the data-analysis program and the difficulties in using it are not in application of the data model but rather in the determination of the unknown data-model parameters. The refinement of the data-model parameters is an order of magnitude more complex (in terms of the computer code) than the integration of the Bragg peaks when the parameters are known.

The data model is a compromise between an attempt to describe the measurement process precisely and the ability to find parameters describing this process. For example, the overlap between the Bragg peaks is typically ignored due to the complexity of spot-shape determination when reflections overlap. The issue is not only to implement the parameterization, but also to do it with acceptable speed and stability of the numerical algorithms. A more complex data model can be more precise (realistic) under specific circumstances, but can result in a less stable refinement and produce less precise final results in most cases. An apparently more realistic (complex) data model may end up being inferior to a simpler and more robust approach. The complexity of model-quality analysis is due to the fact that some types of errors may be much less significant than others. In particular, an error that changes the intensities of all reflections by the same factor only changes the overall scale factor between the data and the atomic model. Truncation of the integration area results in a systematic reduction of calculated reflection intensities. A variable integration area may result in a different fraction of a reflection being omitted for different reflections. The goal of an integration method is to minimize the variation in the omitted fraction, rather than its magnitude. Similarly, if there is an error in predicting reflection-profile shape, this constant error has a smaller impact than a variable error of the same magnitude.

The magnitude and types of errors are very different in different experiments. The compensation of errors also differs between experiments, making it hard to generalize about an optimal approach to data analysis when the data do not fully satisfy the assumptions of the data model. For intense reflections, when counting statistics are not a limiting factor, none of the current data models accounts for all reproducible errors in experiments. This issue is critical in measuring small differences originating from dispersive effects.

11.4.5.3. Data-model refinement

The parameters of the data model can be classified into four groups:

(1) Those refinable from self-consistency of the data by a (nonlinear) least-squares method.