

2. RECIPROCAL SPACE IN CRYSTAL-STRUCTURE DETERMINATION

transform (FFT) algorithm if only the xy translation is sought (2D FFT), or if only the rotation angle is needed (1D FFT).

2D alignment methods can be divided into three classes: (1) those that employ exhaustive searches in order to find three orientation parameters; (2) those that perform exhaustive searches by using either simplifications (separate searches for translation and rotation parameters) (Penczek *et al.*, 1992) or by taking advantage of invariant image representations (Schatz & van Heel, 1990; Frank *et al.*, 1992 and the following discussion; Schatz & van Heel, 1992; Marabini & Carazo, 1996); or finally (3) those that are aimed at improvement of previously determined parameters and employ local searches.

In practice, as the windowed particles are approximately centred, the search for translation parameters can be restricted to relatively small values. A very efficient algorithm that takes advantage of the geometry is based on resampling to polar coordinates of the area of the image that roughly corresponds to the particle size. The resampling is done around centres placed on pixels located within a distance from the image centre that corresponds to a preset maximum translation (Joyeux & Penczek, 2002) (Fig. 2.5.7.2). For each translation, a 1D rotational cross-correlation function in polar coordinates is calculated. Overall, the alignment method based on resampling to polar coordinates comprises the following steps: (1) the image is resampled to polar coordinates; (2) 1D FFTs of various lengths are calculated, appropriately weighted and padded with zeros to equalize their lengths; (3) complex multiplications with 1D Fourier transforms of the similarly processed referenced image are calculated; (4) the inverse 1D FFT is calculated and the position of the maximum is found. The last step yields the rotation angle. Steps (1)–(4) are repeated with the image that is being aligned shifted to account for translations. In addition, the rotation angle for one of the images being mirrored is efficiently calculated in parallel with step (3) by repeating the multiplication with the 1D Fourier transforms of the reference image complex conjugated. This additional check is a necessity in the analysis of single-particle data sets, as usually one can expect on average half of the images to be mirrored versions of the other half in the data set. Overall, the method is very accurate, because only data under the circular mask enter the calculation.

For a set of N images containing the same object in various orientations and corrupted by an additive noise, the problem of alignment would be relatively simple. For proteins that have strong preferred orientation and particularly when a staining technique is used for grid preparation, this is certainly the case. In the procedure called *reference-based alignment*, one of the images that appears ‘typical’ is selected and used as a reference to align the remaining images. After all available images are aligned their average is calculated and used as a reference in a repeated alignment of all images. The process is iterated until the orientations of the images stabilize (Frank *et al.*, 1982).

More formally, Frank *et al.* (1988) proposed the definition of a set of N images f_k , $k = 1, \dots, N$, aligned if a set of transformations \mathbf{T}_k , $k = 1, \dots, N$, (rotation angles and translations) is found such that all pairs of images are mutually brought into register, so the expression

$$\begin{aligned} L_1(\{f\}, \{\mathbf{T}\}) &= \sum_{k=1}^{N-1} \sum_{l=k+1}^N \|f_k(\mathbf{T}_k \mathbf{x}) - f_l(\mathbf{T}_l \mathbf{x})\|^2 \\ &= \sum_{k=1}^{N-1} \sum_{l=k+1}^N \left(\|f_k(\mathbf{T}_k \mathbf{x})\|^2 + \|f_l(\mathbf{T}_l \mathbf{x})\|^2 - 2f_k(\mathbf{T}_k \mathbf{x})f_l(\mathbf{T}_l \mathbf{x}) \right) \end{aligned} \quad (2.5.7.9)$$

is minimized. Although there is no simple way to minimize L_1 , the interesting observation is that there is no requirement of the images to represent the same particle, not even a similar one. This

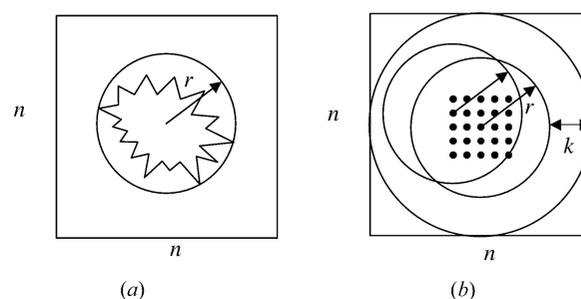


Fig. 2.5.7.2. The geometrical constraints of the 2D alignment problem. (a) The reference 2D particle is placed within a square image frame $n \times n$ pixels and its size is such that it can be bounded by a circle with a radius r no larger than $0.9n$. (b) The particle projection, the size of which is bounded by the same radius as the reference view, can be located within a circle centred on discrete locations within the image frame, such that the maximum translation is $k = (n/2) - r$. The number of possible translations is $(2k + 1)^2$. Reprinted from Joyeux & Penczek (2002) with permission from Elsevier.

leads to the conclusion that if the minimum of L_1 could be found, a set of diverse images could be aligned; moreover, upon alignment similar images would have similar orientation and subsequent classification of such an aligned data set would reveal subsets of similar images.

A practical method of minimizing, called a *reference-free alignment*, was proposed by Penczek *et al.* (1992) by showing that minimization of L_1 is equivalent to maximization of

$$L_2(\{f\}, \{\mathbf{T}\}) = \sum_{k=1}^{N-1} \|f_k(\mathbf{T}_k \mathbf{x}) - \langle f \rangle_k\|^2, \quad (2.5.7.10)$$

where

$$\langle f \rangle_k = \frac{1}{N-1} \sum_{l=1, l \neq k}^N f_l(\mathbf{T}_l \mathbf{x}) \quad (2.5.7.11)$$

is the partial average of the set of images calculated with the exclusion of the k th image. The method is based on the observation that given a set of approximately aligned images, it should be possible to minimize L_2 by sequentially correcting alignments of individual images using the cross-correlation function between each image and the average of the remaining ones. On each step, depending whether the orientation of the image changes or not, (2.5.7.10) will decrease or remain constant.

The outcome of the reference-free alignment algorithm is an aligned set of N images, so all particles that have similar shapes will have similar orientations. Thus, it is natural (and because of the alignment possible) to divide the data set into classes of images that have similar shapes and orientations, *i.e.*, to cluster them. A number of well known clustering algorithms have been adopted for EM applications (Frank, 1990). The general purpose of clustering is to organize objects (in the case of EM, images) into classes whose members are similar to each other, while dissimilar to objects from other classes.

Reference-free alignment with subsequent clustering works well as long as all particles share the same overall shape (*i.e.*, the very low frequency component), as is the case for ribosomes. However, some molecules yield projections that have quite different shapes, as for example is the case for barrel-like proteins GroEL (Roseman *et al.*, 1996) with rectangular views and circular end views or flat and rectangular hemocyanin (Boisset *et al.*, 1995). In this case, the reference-free alignment tends to be unstable, as (2.5.7.10) has multiple local minima, which in practice means that the global average of the whole data set can vary significantly depending on the initiation of the procedure. In general, reference-free alignment is an ‘alignment first, classifi-