

## 2. RECIPROCAL SPACE IN CRYSTAL-STRUCTURE DETERMINATION

exceeding atomic resolution, but the difficulties in overcoming the very low signal-to-noise ratio (SNR) and low contrast in the data, combined with the adverse effects of the contrast transfer function (CTF) of the microscope, hamper progress in fulfilling the potential of the technique. However, in recent years, cryo-EM has proven its power in the structure determination of large macromolecular assemblies and machines which are too large and complex for the more traditional techniques of structural biology, *i.e.*, X-ray crystallography and NMR spectroscopy.

Single-particle reconstruction is based on the assumption that a protein exists in solution in multiple copies of the same basic structure. Unlike in crystallography, no ordering of the structure within a crystal grid is required; the enhancement of the SNR is achieved by bringing projection images of different (but structurally identical) proteins into register and averaging them. This is why the technique is sometimes called ‘crystallography without crystals’.

Within the linear weak-phase-object approximation of the image formation process in the microscope [see equation (2.5.2.43) in Section 2.5.2], 2D projections represent line integrals of the Coulomb potential of the particle under examination convoluted with the point-spread function of the microscope,  $s$ , as introduced in Section 2.5.1. In addition, we have to consider the translation  $\mathbf{t}$  of the projection in the plane of micrograph, suppression of high-frequency information by the envelope function  $E$  of the microscope, and two additive noises  $m^B$  and  $m^S$ . The first one is a coloured background noise, while the second is attributed to the residual scattering by the solvent or the supporting thin layer of carbon, if used, assumed to be white and affected by the transfer function of the microscope in the same way as the imaged protein. In order to have the image formation model correspond more closely to the physical reality of data collection, we write equation (2.5.6.4) from Section 2.5.6 such that the projection operation is always realized in the  $z$  direction of the coordinate system (corresponding to the direction of propagation of the electron beam), while the molecule is rotated arbitrarily by three Eulerian angles:

$$d_n(\mathbf{x}) = s_n(\mathbf{x}) * e_n(\mathbf{x}) * \left[ \int f(\mathbf{T}_n \mathbf{r}) dz + m_n^S(\mathbf{x}) \right] + m_n^B(\mathbf{x}),$$

$$n = 1, \dots, N. \quad (2.5.7.1)$$

Here  $f \in R^{n^3}$  represents the three-dimensional (3D) electron density of the imaged macromolecule and  $d \in R^{n^2}$  is the  $n$ th observed two-dimensional (2D) projection image. The total number of projection images  $N$  depends on the structure determination project, and can vary from a few hundred to hundreds of thousands. Further,  $e$  is the inverse Fourier transform of the envelope function,  $\mathbf{x} = [x \ y]^T$  is a vector of coordinates in the plane of projections,  $\mathbf{r} = [r_x \ r_y \ r_z \ 1]^T$  is a vector of coordinates associated with  $n$ th macromolecule,  $\mathbf{T}$  is the  $4 \times 4$  transformation matrix given by

$$\mathbf{T}(\mathbf{R}, \mathbf{t}) = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 1 \end{bmatrix}, \quad \begin{bmatrix} \mathbf{x} \\ z \\ 1 \end{bmatrix} = \mathbf{T} \mathbf{r}, \quad (2.5.7.2)$$

with  $\mathbf{t} = [t_x \ t_y \ 1]^T$  being the shift vector of translation of the object (and its projection) in the  $xy$  plane (translation in  $z$  is irrelevant due to the projection operation) and  $\mathbf{R}(\psi, \theta, \varphi)$  is the  $3 \times 3$  rotation matrix specified by three Eulerian angles. As in Section 2.5.6, two of the angles define the direction of projection  $\tau(\theta, \varphi)$ , while the third angle  $\psi$  results in rotation of the projection image in the plane of the formed image  $xy$ ; changing this angle does not provide any additional information about the structure  $f$ . Both types of noise are assumed to be mutually uncorrelated and independent between projection images (*i.e.*,

$\langle m_i^k m_j^l \rangle_{i \neq j} = 0$ ;  $k, l = S, B$ ) and also uncorrelated with the signal ( $\langle d_i m_i^k \rangle = 0$ ;  $k = S, B$ ). Model (2.5.7.1) is semi-empirical in that, unlike in the standard model, we have two contributions to the noise. Although in principle amorphous ice should not be affected by the CTF, so the term  $m^S$  should be absorbed into  $m^B$ , in practice the buffer in which the protein is purified is not pure water and it is possible to observe CTF effects by imaging frozen buffer alone. Moreover, if a thin support carbon is used, it will be a source of very strong CTF-affected noise also included in  $m^B$ .

In Fourier space, (2.5.7.1) is written by taking advantage of the central section theorem [equation (2.5.6.8) of Section 2.5.6]: the Fourier transform of a projection is extracted as a Fourier plane  $uv$  of a rotated Fourier transform of a 3D object:

$$D_n(\mathbf{u}) = \text{CTF}(\mathbf{u}; \Delta f_n, q) E_n(\mathbf{u}) \left\{ [F(\mathbf{T} \mathbf{v})]_{u_z=0} + M_n^S(\mathbf{u}) \right\} + M_n^B(\mathbf{u}). \quad (2.5.7.3)$$

The capital letters denote Fourier transforms of objects appearing in (2.5.7.1) while CTF (a Fourier transform of  $s$ ) depends, among other parameters that are set very accurately (such as the accelerating voltage of the microscope), on the defocus setting  $\Delta f_n$  and the amplitude contrast ratio  $0 \leq q < 1$  that reflects the presence of the amplitude contrast that is due to the removal of widely scattered electrons [the real term in (2.5.5.14)]. For the range of frequency considered,  $q$  is assumed to be constant and the CTF is written in terms of the phase perturbation function  $\chi$  [given by equation (2.5.2.33)] as

$$\begin{aligned} \text{CTF}(\mathbf{u}; \Delta f) &= [1 + 2q(q-1)]^{-1/2} \left\{ (1-q) \sin[\chi(|\mathbf{u}|; \Delta f)] \right. \\ &\quad \left. - q \cos[\chi(|\mathbf{u}|; \Delta f)] \right\} \\ &= \sin \left\{ \chi(|\mathbf{u}|; \Delta f) - \arctan[q/(1-q)] \right\}, \end{aligned} \quad (2.5.7.4)$$

where for simplicity we assumed no astigmatism. Finally, the rotationally averaged power spectrum of the observed image, calculated as the expectation value of its squared Fourier intensities (2.5.7.3), is given by

$$P_d(u) = \text{CTF}^2(u) E^2(u) [P_f(u) + P_S(u)] + P_B(u), \quad (2.5.7.5)$$

where  $u = |\mathbf{u}|$  is the modulus of spatial frequency.

## 2.5.7.2. Structure determination in single-particle reconstruction

The goal of single-particle reconstruction is to determine the 3D electron-density map  $f$  of a biological macromolecule such that its projections agree in a least-squares sense with a large number of collected 2D electron-microscopy projection images,  $d_n \in R^{n^2}$  ( $n = 1, 2, \dots, N$ ), of isolated (single) particles with random and unknown orientations. Thus, we seek a least-squares solution to the problem stated by (2.5.7.1) [or, equivalently, in Fourier space, to (2.5.7.3)]. This is formally written as a nonlinear optimization problem (Yang *et al.*, 2005),

$$\begin{aligned} \min_{\psi_n, \theta_n, \varphi_n, t_x, t_y, f, \Delta f_n, q, \dots} & L(\psi_n, \theta_n, \varphi_n, t_x, t_y, f, \Delta f_n, q, \dots) \\ & \equiv \frac{1}{2} \sum_{n=1}^N \| s_n(\mathbf{x}) * e_n(\mathbf{x}) * \int f(\mathbf{T}_n \mathbf{r}) dz - d_n(\mathbf{x}) \|^2. \end{aligned} \quad (2.5.7.6)$$

The factor of  $\frac{1}{2}$  is included merely for convenience. The objective function in (2.5.7.6) is clearly nonlinear due to the coupling

## 2.5. ELECTRON DIFFRACTION AND ELECTRON MICROSCOPY IN STRUCTURE DETERMINATION

EM data collection					
Film			CCD		
Analysis of power spectra					
Estimation of astigmatism, defocus, envelope function, background noise, signal-to-noise ratio of the data					
Particle picking					
Manual		Semi-automated		Automated	
2D alignment					
Using invariants		Reference-free		Multireference	
2D classification					
K-means		Hierarchical		Other	
Initial model					
Random	Guessed	Homology	Ab initio using class averages	Experimental	
				Random conical tilt	Tomography
3D refinement					
3D projection matching		Unified	Fourier space refinement	Multireference alignment	
Analysis of the 3D map					
Surface representation	Docking	Segmentation	Detection of secondary structure elements	Real-space variance	Conformational modes

Fig. 2.5.7.1. Typical steps performed in a single-particle cryo-EM structure determination project.

between the orientation parameters  $\psi_n, \theta_n, \varphi_n, t_{x_n}, t_{y_n}$  ( $n = 1, 2, \dots, N$ ) and the 3D density  $f$ .

The parameters in (2.5.7.6) to be determined can be separated into two groups. (1) The orientation parameters  $\psi_n, \theta_n, \varphi_n, t_{x_n}, t_{y_n}$  that have to be determined entirely by solving (2.5.7.6) and for which there are no initial guesses, and the structure  $f$  itself, for which we may or may not have an initial guess. The number of parameters in this group is very large:  $n^3 + 5m$ . Note that in single-particle reconstruction, the number of projection data  $m$  is far greater than the linear size of the data in pixels, *i.e.*,  $m \gg n$ . (2) Various parameters which we will broadly call the parameters of the image formation model (2.5.7.1)–(2.5.7.4): the defocus settings of the microscope  $\Delta f_n$ , the amplitude contrast ratio  $q$  and, if analytical forms of the envelope function  $E$ , the power spectrum of the background noise  $M$ , or the structure  $F$  are adopted, the parameters of these equations. Some of the parameters in the second group are usually known very accurately or can be estimated from micrograph data before one attempts to solve (2.5.7.6) (see Section 2.5.7.4), but they can also be refined during the structure determination process [for the method for correcting the defocus settings, see Mouche *et al.* (2001)].

Owing to the very large number of parameters in (2.5.7.6) and the nonlinearities present, one almost never attempts to solve the problem directly. Instead, structure determination using the single-particle technique involves several steps. (i) The macromolecular complex is prepared with a purity of at least 90%. (ii) The sample is flash-frozen in liquid ethane. Alternatively, cryo-negative stain techniques or traditional negative stain methods can be used. (iii) Pictures of the macromolecular complexes are taken. (iv) Exhaustive analysis of 2D particle images aimed at increasing the SNR of the data and evaluation of the homogeneity of the sample is performed. (v) An initial low-resolution model of the structure is established using either experimental techniques or computational methods. (vi) The initial structure is refined in order to increase the resolution using an enlarged data set. Only in this step does one attempt to minimize (2.5.7.6) more

or less directly. (vii) Visualization and interpretation of the resulting 3D electron-density map is the last step; it often involves docking of X-ray structures of molecules into EM density maps in order to reveal the arrangement of known molecules within the EM envelope (Fig. 2.5.7.1). As within the weak-phase-object approximation of the image formation in EM the relation between densities in collected images and the 3D electron density of the imaged macromolecule is linear [(2.5.7.1)], all data-processing methods employed in the structure determination project should be linear, so the densities in the cryo-EM 3D model can be interpreted in terms of the electron density of the protein.

In the actual single-particle project not all the steps have to be executed in the order outlined above. The technique has proved to be particularly useful in studies of functional complexes of proteins whose base state is known to a certain resolution or even of functional complexes whose atomic (X-ray crystallographic) structure is known. In these cases, steps (iv) and (v) can be omitted and the structure of the functional complex (for examples with ligands bound to it) can be relatively easily determined using the native structure as a starting point for step (vi).

In addition to difficulties with obtaining good cryo-EM data, the technique is computationally intensive. The reason is that in order to obtain a sufficient SNR in the 3D structure, processing of hundreds of thousands of EM projection images of the molecule might be necessary. For each, five orientation parameters have to be determined, and this is in addition to determination of the image-formation parameters required for the optimization of correlation searches. In effect, it is not unusual for single-particle projects to consume weeks of the computer time of multi-processing clusters. This also explains why the knowledge of the base structure simplifies the work to a large degree: when it is known, initial values of the orientation parameters can be easily established, reducing not only the computational time, but also possibilities of errors in the structure-determination process.

### 2.5.7.3. Electron microscopy and data digitization

The electron microscope is a phase imaging system; *i.e.*, in order to create contrast in images, they have to be underfocused. Owing to the particular form of the CTF of the microscope [(2.5.7.4)], not only the amplitudes of the image in Fourier space are modified, but information in some ranges of spatial frequencies is set to zero and some phases have reversed sign. Therefore, in order to obtain possibly uniform coverage of Fourier space, the standard practice is to take pictures using different defocus settings and merge them computationally in order to fill gaps in Fourier space. The problem is compounded by the relation between underfocus and the envelope function of the microscope. Far-from-focus images have high contrast, but the envelope function has a relatively steep fall-off limiting the range of useful spatial frequencies. Conversely, close-to-focus images have little contrast, but the envelope function is decreasing, slowly extending useful information to high spatial frequencies. In effect, it is easier to process computationally far-from-focus data and to obtain accurate alignment of particles, but the results have severely limited resolution. Processing of close-to-focus data is challenging and results tend to be less accurate, but there is the potential to obtain high-resolution information.

The experimental techniques of initial structure determination (random conical tilt, tomography) require collection of tilt data. This is facilitated by dedicated microscope stages that can be rotated inside the microscope column yielding additional views of the same field. However, collection of high-quality tilt images is difficult. The quality of tilted images tends to be adversely affected by charging and drift effects. Moreover, as the stage is tilted the effective ice thickness increases (inversely proportionally to the cosine of the tilt angle, so at 60° the factor is two) and the contrast of the images decreases correspondingly. Finally,