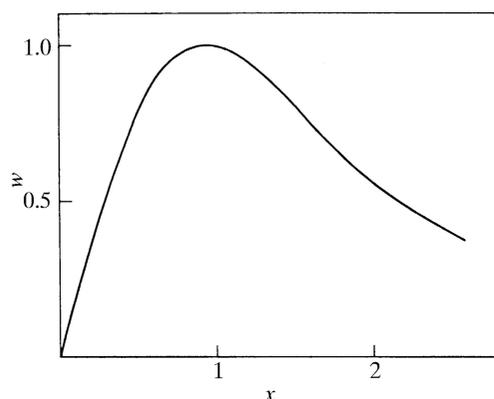


2. RECIPROCAL SPACE IN CRYSTAL-STRUCTURE DETERMINATION


 Fig. 2.2.7.1. The form of w as given by (2.2.7.2).

usually available in direct procedures) are considered as additional *a priori* information so that (2.2.7.1) may be replaced by

$$\tan \varphi_{\mathbf{h}} \simeq \frac{\sum_j \beta_j \sin(\varphi_{\mathbf{k}_j} + \varphi_{\mathbf{h}-\mathbf{k}_j})}{\sum_j \beta_j \cos(\varphi_{\mathbf{k}_j} + \varphi_{\mathbf{h}-\mathbf{k}_j})}, \quad (2.2.7.3)$$

where β_j is the solution of the equation

$$D_1(\beta_j) = D_1(G_j)D_1(\alpha_{\mathbf{k}_j})D_1(\alpha_{\mathbf{h}-\mathbf{k}_j}). \quad (2.2.7.4)$$

In (2.2.7.4),

$$G_j = 2|E_{\mathbf{h}}E_{\mathbf{k}_j}E_{\mathbf{h}-\mathbf{k}_j}|\sqrt{N}$$

or the corresponding second representation parameter, and $D_1(x) = I_1(x)/I_0(x)$ is the ratio of two modified Bessel functions.

In order to promote (in accordance with the aims of Hull and Irwin) the agreement between α and $\langle \alpha \rangle$, the distribution of α may be used (Cascarano, Giacovazzo, Burla *et al.*, 1984; Burla *et al.*, 1987); in particular, the first two moments of the distribution: accordingly,

$$w = \left\{ \exp \left[\frac{-(\alpha - \langle \alpha \rangle)^2}{2\sigma_\alpha^2} \right] \right\}^{1/3}$$

may be used, where σ_α^2 is the estimated variance of α .

Stage 7: Figures of merit. The correct solution is found among several by means of figures of merit (FOMs) which are expected to be extreme for the correct solution. Largely used are (Germain *et al.*, 1970)

$$(a) \quad \text{ABSFOM} = \frac{\sum_{\mathbf{h}} \alpha_{\mathbf{h}}}{\sum_{\mathbf{h}} \langle \alpha_{\mathbf{h}} \rangle},$$

which is expected to be unity for the correct solution.

$$(b) \quad \text{PSIO} = \frac{\sum_{\mathbf{h}} |\sum_{\mathbf{k}} E_{\mathbf{k}} E_{\mathbf{h}-\mathbf{k}}|}{\sum_{\mathbf{h}} (\sum_{\mathbf{k}} |E_{\mathbf{k}} E_{\mathbf{h}-\mathbf{k}}|^2)^{1/2}}.$$

The summation over \mathbf{k} includes (Cochran & Douglas, 1957) the strong $|E|$'s for which phases have been determined, and indices \mathbf{h} correspond to very small $|E_{\mathbf{h}}|$. Minimal values of PSIO (≤ 1.20) are expected to be associated with the correct solution.

$$(c) \quad R_\alpha = \frac{\sum_{\mathbf{h}} |\alpha_{\mathbf{h}} - \langle \alpha_{\mathbf{h}} \rangle|}{\sum_{\mathbf{h}} \langle \alpha_{\mathbf{h}} \rangle}.$$

That is, the Karle & Karle (1966) residual between the actual and the estimated α 's. After scaling of $\alpha_{\mathbf{h}}$ on $\langle \alpha_{\mathbf{h}} \rangle$ the correct solution should be characterized by the smallest R_α values.

$$(d) \quad \text{NQEST} = \sum_j G_j \cos \Phi_j,$$

where G is defined by (2.2.5.21) and

$$\Phi = \varphi_{\mathbf{h}} - \varphi_{\mathbf{k}} - \varphi_{\mathbf{l}} - \varphi_{\mathbf{h}-\mathbf{k}-\mathbf{l}}$$

are quartet invariants characterized by large basis magnitudes and small cross magnitudes (De Titta *et al.*, 1975; Giacovazzo, 1976). Since G is expected to be negative as well as $\cos \Phi$, the value of NQEST is expected to be positive and a maximum for the correct solution.

Figures of merit are then combined as

$$\begin{aligned} \text{CFOM} = & w_1 \frac{\text{ABSFOM} - \text{ABSFOM}_{\min}}{\text{ABSFOM}_{\max} - \text{ABSFOM}_{\min}} \\ & + w_2 \frac{\text{PSIO}_{\max} - \text{PSIO}}{\text{PSIO}_{\max} - \text{PSIO}_{\min}} \\ & + w_3 \frac{R_{\alpha_{\max}} - R_\alpha}{R_{\alpha_{\max}} - R_{\alpha_{\min}}} \\ & + w_4 \frac{\text{NQEST} - \text{NQEST}_{\min}}{\text{NQEST}_{\max} - \text{NQEST}_{\min}}, \end{aligned}$$

where w_i are empirical weights proportional to the confidence of the user in the various FOMs.

Different FOMs are often used by some authors in combination with those described above: for example, enantiomorph triplets and quartets are supplementary FOMs (Van der Putten & Schenk, 1977; Cascarano, Giacovazzo & Viterbo, 1987).

Different schemes of calculating and combining FOMs are also used: one scheme (Cascarano, Giacovazzo & Viterbo, 1987) uses

$$(a1) \quad \text{CPHASE} = \frac{\sum w_j G_j \cos(\Phi_j - \theta_j) + w_j G_j \cos \Phi_j}{\sum_{\text{s.i.}+\text{s.s.}} w_j G_j D_1(G_j)},$$

where the first summation in the numerator extends over symmetry-restricted one-phase and two-phase s.s.'s (see Sections 2.2.5.9 and 2.2.5.10), and the second summation in the numerator extends over negative triplets estimated *via* the second representation formula [equation (2.2.5.13)] and over negative quartets. The value of CPHASE is expected to be close to unity for the correct solution.

(a2) $\alpha_{\mathbf{h}}$ for strong triplets and $E_{\mathbf{k}}E_{\mathbf{h}-\mathbf{k}}$ contributions for PSIO triplets may be considered random variables: the agreements between their actual and their expected distributions are considered as criteria for identifying the correct solution.

(a3) correlation among some FOMs is taken into account.

According to this scheme, each FOM (as well as the CFOM) is expected to be unity for the correct solution. Thus one or more figures are available which constitute a sort of criterion (on an absolute scale) concerning the correctness of the various solutions: FOMs (and CFOM) $\simeq 1$ probably denote correct solutions, CFOMs $\ll 1$ should indicate incorrect solutions.

Stage 8: Interpretation of E maps. This is carried out in up to four stages (Koch, 1974; Main & Hull, 1978; Declercq *et al.*, 1973):

- (a) peak search;
- (b) separation of peaks into potentially bonded clusters;
- (c) application of stereochemical criteria to identify possible molecular fragments;
- (d) comparison of the fragments with the expected molecular structure.

2.2.8. Other multisolution methods applied to small molecules

In very complex structures a large initial set of known phases seems to be a basic requirement for a structure to be determined.

2.2. DIRECT METHODS

Table 2.2.8.1. Magic-integer sequences for small numbers of phases (n) together with the number of sets produced and the root-mean-square error in the phases

n	Sequence								No. of sets	R.m.s. error ($^\circ$)
1	1								4	26
2	2	3							12	29
3	3	4	5						20	37
4	5	7	8	9					32	42
5	8	11	13	14	15				50	45
6	13	18	21	23	24	25			80	47
7	21	29	34	37	39	40	41		128	48
8	34	47	55	60	63	65	66	67	206	49

This aim can be achieved, for example, by introducing a large number of permutable phases into the initial set. However, the introduction of every new symbol implies a fourfold increase in computing time, which, even in fast computers, quickly leads to computing-time limitations. On the other hand, a relatively large starting set is not in itself enough to ensure a successful structure determination. This is the case, for example, when the triplet invariants used in the initial steps differ significantly from zero. New strategies have therefore been devised to solve more complex structures.

(1) Magic-integer methods

In the classical procedure described in Section 2.2.7, the unknown phases in the starting set are assigned all combinations of the values $\pm\pi/4$, $\pm 3\pi/4$. For n unknown phases in the starting set, 4^n sets of phases arise by quadrant permutation; this is a number that increases very rapidly with n . According to White & Woolfson (1975), phases can be represented for a sequence of n integers by the equations

$$\varphi_i = m_i x \pmod{2\pi}, \quad i = 1, \dots, n. \quad (2.2.8.1)$$

The set of equations can be regarded as the parametric equation of a straight line in n -dimensional phase space. The nature and size of errors connected with magic-integer representations have been investigated by Main (1977) who also gave a recipe for deriving magic-integer sequences which minimize the r.m.s. errors in the represented phases (see Table 2.2.8.1). To assign a phase value, the variable x in equation (2.2.8.1) is given a series of values at equal intervals in the range $0 < x < 2\pi$. The enantiomorph is defined by exploring only the appropriate half of the n -dimensional space.

A different way of using the magic-integer method (Declercq *et al.*, 1975) is the *primary-secondary P-S method* which may be described schematically in the following way:

(a) Origin- and enantiomorph-fixing phases are chosen and some one-phase s.s.'s are estimated.

(b) Nine phases [this is only an example: very long magic-integer sequences may be used to represent primary phases (Hull *et al.*, 1981; Debaerdemaeker & Woolfson, 1983)] are represented with the approximated relationships:

$$\begin{cases} \varphi_{i_1} = 3x \\ \varphi_{i_2} = 4x \\ \varphi_{i_3} = 5x \end{cases} \quad \begin{cases} \varphi_{j_1} = 3y \\ \varphi_{j_2} = 4y \\ \varphi_{j_3} = 5y \end{cases} \quad \begin{cases} \varphi_{p_1} = 3z \\ \varphi_{p_2} = 4z \\ \varphi_{p_3} = 5z. \end{cases}$$

Phases in (a) and (b) constitute the *primary set*.

(c) The phases in the *secondary set* are those defined through \sum_2 relationships involving pairs of phases from the primary set: they, too, can be expressed in magic-integer form.

(d) All the triplets that link together the phases in the combined primary and secondary set are now found, other than triplets used to obtain secondary reflections from the primary ones. The general algebraic form of these triplets will be

$$m_1x + m_2y + m_3z + b \equiv 0 \pmod{1},$$

where b is a phase constant which arises from symmetry translation. It may be expected that the 'best' value of the unknown x, y, z corresponds to a maximum of the function

$$\psi(x, y, z) = \sum |E_1 E_2 E_3| \cos 2\pi(m_1x + m_2y + m_3z + b),$$

with $0 \leq x, y, z < 1$. It should be noticed that ψ is a Fourier summation which can easily be evaluated. In fact, ψ is essentially a figure of merit for a large number of phases evaluated in terms of a small number of magic-integer variables and gives a measure of the internal consistency of \sum_2 relationships. The ψ map generally presents several peaks and therefore can provide several solutions for the variables.

(2) The random-start method

These are procedures which try to solve crystal structures by starting from random initial phases (Baggio *et al.*, 1978; Yao, 1981). They may be so described:

(a) A number of reflections (say NUM ~ 100 or larger) at the bottom of the CONVERGE map are selected. These, and the relationships which link them, form the system for which trial phases will be found.

(b) A pseudo-random number generator is used to generate M sets of NUM random phases. Each of the M sets is refined and extended by the tangent formula or similar methods.

(3) *Accurate calculation of s.i.'s and s.s.'s with 1, 2, 3, 4, ..., n phases*

Having a large set of good phase relationships allows one to overcome difficulties in the early stages and in the refinement process of the phasing procedure. Accurate estimates of s.i.'s and s.s.'s may be achieved by the application of techniques such as the representation method or the neighbourhood principle (Hauptman, 1975; Giacovazzo, 1977a, 1980b). So far, second-representation formulae are available for triplets and one-phase seminvariants; in particular, reliably estimated negative triplets can be recognized, which is of great help in the phasing process (Casarano, Giacovazzo, Camalli *et al.*, 1984). Estimation of higher-order s.s.'s with upper representations or upper neighbourhoods is rather difficult, both because the procedures are time consuming and because the efficiency of the present joint probability distribution techniques deteriorates with complexity. However, further progress can be expected in the field.

(4) *Modified tangent formulae and least-squares determination and refinement of phases*

The problem of deriving the individual phase angles from triplet relationships is greatly overdetermined: indeed the number of triplets, in fact, greatly exceeds the number of phases so that any φ_h may be determined by a least-squares approach (Hauptman *et al.*, 1969). The function to be minimized may be

$$M = \frac{\sum_k w_k [\cos(\varphi_h - \varphi_k - \varphi_{h-k}) - C_k]^2}{\sum w_k},$$

where C_k is the estimate of the cosine obtained by probabilistic or other methods.

Effective least-squares procedures based on linear equations (Debaerdemaeker & Woolfson, 1983; Woolfson, 1977) can also be used. A triplet relationship is usually represented by

$$(\varphi_p \pm \varphi_q \pm \varphi_r + b) \approx 0 \pmod{2\pi}, \quad (2.2.8.2)$$

where b is a factor arising from translational symmetry. If (2.2.8.2) is expressed in cycles and suitably weighted, then it may be written as

$$w(\varphi_p \pm \varphi_q \pm \varphi_r + b) = wn,$$

where n is some integer. If the integers were known then the equation would appear (in matrix notation) as

2. RECIPROCAL SPACE IN CRYSTAL-STRUCTURE DETERMINATION

$$\mathbf{A}\Phi = \mathbf{C}, \quad (2.2.8.3)$$

giving the least-squares solution

$$\Phi = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{C}. \quad (2.2.8.4)$$

When approximate phases are available, the nearest integers may be found and equations (2.2.8.3) and (2.2.8.4) constitute the basis for further refinement.

Modified tangent procedures are also used, such as (Sint & Schenk, 1975; Busetta, 1976)

$$\tan \varphi_{\mathbf{h}} \simeq \frac{\sum_j G_{\mathbf{h}, \mathbf{k}_j} \sin(\varphi_{\mathbf{k}_j} + \varphi_{\mathbf{h}-\mathbf{k}_j} - \Delta_j)}{\sum_j G_{\mathbf{h}, \mathbf{k}_j} \cos(\varphi_{\mathbf{k}_j} + \varphi_{\mathbf{h}-\mathbf{k}_j} - \Delta_j)},$$

where Δ_j is an estimate for the triplet phase sum ($\varphi_{\mathbf{h}} - \varphi_{\mathbf{k}_j} - \varphi_{\mathbf{h}-\mathbf{k}_j}$).

(5) *Techniques based on the positivity of Karle–Hauptman determinants*

(The main formulae have been briefly described in Section 2.2.5.7.) The maximum determinant rule has been applied to solve small structures (de Rango, 1969; Vermin & de Graaff, 1978) *via* determinants of small order. It has, however, been found that their use (Taylor *et al.*, 1978) is not of sufficient power to justify the larger amount of computing time required by the technique as compared to that required by the tangent formula.

(6) *Tangent techniques using simultaneously triplets, quartets, . . .*

The availability of a large number of phase relationships, in particular during the first stages of a direct procedure, makes the phasing process easier. However, quartets are sums of two triplets with a common reflection. If the phase of this reflection (and/or of the other cross terms) is known then the quartet probability formulae described in Section 2.2.5.5 cannot hold. Similar considerations may be made for quintet relationships. Thus triplet, quartet and quintet formulae described in the preceding paragraphs, if used without modifications, will certainly introduce systematic errors in the tangent refinement process.

A method which takes into account correlation between triplets and quartets has been described (Giacovazzo, 1980c) [see also Freer & Gilmore (1980) for a first application], according to which

$$\tan \varphi_{\mathbf{h}} \simeq \frac{\sum_{\mathbf{k}} G \sin(\varphi_{\mathbf{k}} + \varphi_{\mathbf{h}-\mathbf{k}}) - \sum_{\mathbf{k}, \mathbf{l}} G' \sin(\varphi_{\mathbf{k}} + \varphi_{\mathbf{l}} + \varphi_{\mathbf{h}-\mathbf{k}-\mathbf{l}})}{\sum_{\mathbf{k}} G \cos(\varphi_{\mathbf{k}} + \varphi_{\mathbf{h}-\mathbf{k}}) - \sum_{\mathbf{k}, \mathbf{l}} G' \cos(\varphi_{\mathbf{k}} + \varphi_{\mathbf{l}} + \varphi_{\mathbf{h}-\mathbf{k}-\mathbf{l}})},$$

where G' takes into account both the magnitudes of the cross terms of the quartet and the fact that their phases may be known.

(7) *Integration of Patterson techniques and direct methods (Egert & Sheldrick, 1985) [see also Egert (1983, and references therein)]*

A fragment of known geometry is oriented in the unit cell by real-space Patterson rotation search (see Chapter 2.3) and its position is found by application of a translation function (see Section 2.2.5.4 and Chapter 2.3) or by maximizing the weighted sum of the cosines of a small number of strong translation-sensitive triple phase invariants, starting from random positions. Suitable FOMs rank the most reliable solutions.

(8) *Maximum entropy methods*

A common starting point for all direct methods is a stochastic process according to which crystal structures are thought of as being generated by randomly placing atoms in the asymmetric unit of the unit cell according to some *a priori* distribution. A non-uniform prior distribution of atoms $p(\mathbf{r})$ gives rise to a source of random atomic positions with entropy (Jaynes, 1957)

$$H(p) = - \int_V p(\mathbf{r}) \log p(\mathbf{r}) \, d\mathbf{r}.$$

The maximum value $H_{\max} = \log V$ is reached for a uniform prior $p(\mathbf{r}) = 1/V$.

The strength of the restrictions introduced by $p(\mathbf{r})$ is not measured by $H(p)$ but by $H(p) - H_{\max}$, given by

$$H(p) - H_{\max} = - \int_V p(\mathbf{r}) \log [p(\mathbf{r})/m(\mathbf{r})] \, d\mathbf{r},$$

where $m(\mathbf{r}) = 1/V$. Accordingly, if a prior prejudice $m(\mathbf{r})$ exists, which maximizes H , the revised relative entropy is

$$S(p) = - \int_V p(\mathbf{r}) \log [p(\mathbf{r})/m(\mathbf{r})] \, d\mathbf{r}.$$

The maximization problem was solved by Jaynes (1957). If $G_j(p)$ are linear constraint functionals defined by given constraint functions $C_j(\mathbf{r})$ and constraint values c_j , *i.e.*

$$G_j(p) = \int_V p(\mathbf{r}) C_j(\mathbf{r}) \, d\mathbf{r} = c_j,$$

the most unbiased probability density $p(\mathbf{r})$ under prior prejudice $m(\mathbf{r})$ is obtained by maximizing the entropy of $p(\mathbf{r})$ relative to $m(\mathbf{r})$. A standard variational technique suggests that the constrained maximization is equivalent to the unconstrained maximization of the functional

$$S(p) + \sum_j \lambda_j G_j(p),$$

where the λ_j 's are Lagrange multipliers whose values can be determined from the constraints.

Such a technique has been applied to the problem of finding good electron-density maps in different ways by various authors (Wilkins *et al.*, 1983; Bricogne, 1984; Navaza, 1985; Navaza *et al.*, 1983).

Maximum entropy methods are strictly connected with traditional direct methods: in particular it has been shown that:

(a) the maximum determinant rule (see Section 2.2.5.7) is strictly connected (Britten & Collins, 1982; Piro, 1983; Narayan & Nityananda, 1982; Bricogne, 1984);

(b) the construction of conditional probability distributions of structure factors amounts precisely to a reciprocal-space evaluation of the entropy functional $S(p)$ (Bricogne, 1984).

Maximum entropy methods are under strong development: important contributions can be expected in the near future even if a multipurpose robust program has not yet been written.

2.2.9. Some references to direct-methods packages: the small-molecule case

Some references for direct-methods packages are given below. Other useful packages using symbolic addition or multisolution procedures do exist but are not well documented.

CRUNCH: Gelder, R. de, de Graaff, R. A. G. & Schenk, H. (1993). *Automatic determination of crystal structures using Karle–Hauptman matrices*. *Acta Cryst.* **A49**, 287–293.

DIRDIF: Beurskens, P. T., Beurskens G., de Gelder, R., Garcia-Granda, S., Gould, R. O., Israel, R. & Smits, J. M. M. (1999). *The DIRDIF-99 program system*. Crystallography Laboratory, University of Nijmegen, The Netherlands.

MITHRIL: Gilmore, C. J. (1984). *MITHRIL. An integrated direct-methods computer program*. *J. Appl. Cryst.* **17**, 42–46.

MULTAN88: Main, P., Fiske, S. J., Germain, G., Hull, S. E., Declercq, J.-P., Lessinger, L. & Woolfson, M. M. (1999). *Crystallographic software: teXsan for Windows*. <http://www.rigaku.com/downloads/journal/Vol15.1.1998/texsan.pdf>.

PATSEE: Egert, E. & Sheldrick, G. M. (1985). *Search for a fragment of known geometry by integrated Patterson and direct methods*. *Acta Cryst.* **A41**, 262–268.