

The Use of Discrete Hilbert Transforms in Phase Extension and Improvement

GIUSEPPE ZANOTTI,^{a*} FIORENZA FOGALE^a AND GUIDO CAPITANI^b

^aDipartimento di Chimica Organica e Centro Studi Biopolimeri, Via Marzolo 1, 35131 Padova, Italy, and

^bBiozentrum der Universitaet Basel, Abteilung Strukturbiologie, Klingelbergstrasse 70, CH-4056 Basel, Switzerland. E-mail: zanotti@pdchor.unipd.it

(Received 22 February 1996; accepted 13 May 1996)

Abstract

A new totally model independent procedure for phase extension and improvement in X-ray crystallography, based on the discrete Hilbert transforms, is presented. The method has been tested using simulated diffraction data of a small molecule and simulated and experimental data of a protein crystal. Starting from a randomly incomplete set of correct phases, it allows calculation of the unknown phases. Moreover, a set of phases affected by a mean phase error of $\pm 90^\circ$ can be improved to a mean error of $\pm 25^\circ$ if the correct figures of merit for the reflections are known. The performance and the limitations of the technique, as well as the perspectives for further developments, are discussed.

1. Introduction

The problem of phase extension, *i.e.* that of calculating phases associated with all reflections if only a limited number of phases are known, is a classical one in crystallography. Phase improvement is a related problem. Both are already routinely solved for small-molecule structures, using probabilistic relationships (see *e.g.* Giacovazzo, 1980). However, in macromolecular crystallography, where these problems are particularly relevant, this is not quite the case. In protein structure solution with the multiple-isomorphous-replacement (MIR) technique, very often the isomorphism of the heavy-atom derivatives does not extend beyond 3 Å resolution or so. Usually, a set of phases obtained with the isomorphous-replacement method must be improved further, in order to obtain an interpretable electron-density map. In the single-isomorphous-replacement (SIR) approach, this is even a *conditio sine qua non*. Several techniques based on density modification have been proposed in the past to extend phases and to improve protein electron-density maps (*e.g.* Davies & Rollet, 1976; Schevitz, Podjarny, Zwick, Hughes & Sigler, 1981; Bath & Blow, 1982; Cannillo, Oberti & Ungaretti, 1983; Bryan & Banner, 1987; Wang, 1985; Shiono & Woolfson, 1992; Refaat & Woolfson, 1993). The most commonly used among them is the procedure currently known as *solvent flattening*, as implemented in slightly different ways in

several programs. It consists in detecting, either manually or *via* an automated procedure, the solvent regions in a protein electron-density map and in creating a mask that defines solvent regions. The latter are appropriately flattened, upon which map inversion generates new phases for all reflections, including those that were eventually left out of the ordinary phasing procedure. Solvent flattening is a powerful technique in improving an existing set of phases if a good mask can be devised. In practice, this is only possible if the starting phase set is reasonably reliable.

Another procedure for phase extension and improvement is based on entropy-maximization techniques (Bricogne, 1988). Recently, a more sophisticated approach has been proposed, based on entropy maximization constrained by solvent flattening (Xiang, Carter, Bricogne & Gilmore, 1993).

Here we describe a method of phase extension and improvement based on discrete Hilbert transforms. This procedure is based on a completely different principle from those previously described and has the advantage of being absolutely model independent.

2. Background of the method

Hilbert transforms, also called dispersion relations or Kronig–Kramers relations (Toll, 1956), despite having been known to the crystallographic community for a long time (Ramachandran, 1969), have never found practical applications in the field. In contrast, they are widely used in many other areas, including optics and spectroscopy (VanderNoot, 1992; Williams & Marshall, 1992; Joo & Albrecht, 1993).

For the sake of simplicity, we shall first illustrate the use of Hilbert transforms in the one-dimensional case. The extension to three dimensions will be performed at the end of this section. Assume we have a complex function $F(X) = A(X) + iB(X)$. The dispersion relationships relate the real and the imaginary parts of this function in a simple way:

$$A(X) = (1/\pi)P \int_{-\infty}^{\infty} [B(X')/(X' - X)] dX', \quad (1a)$$

$$B(X) = -(1/\pi)P \int_{-\infty}^{\infty} [A(X')/(X' - X)] dX', \quad (1b)$$

where P stands for the principal part of the Cauchy integral. Relations (1a) and (1b) form a pair of so-called 'Hilbert transforms' (HT). It has been demonstrated by Toll (1956) that (1) is valid in any scattering experiment if $F(X)$ is the Fourier transform of a function $f(t)$ with the property $f(t) = 0$ for $t < 0$ (condition of strict causality). The causality applies of course in the case of time series: when the function f does not depend on time but on space, causality cannot be used in the original sense defined by Toll (1956). Nevertheless, all the formalisms can be applied simply by assuming $f(x) = 0$ for $x < 0$ (see also Burge, Fiddy, Greenaway & Ross, 1976). An important aspect is that no other *a priori* hypothesis on the properties of $f(x)$, except the previous one, is necessary in order to ensure the validity of (1).

The major obstacle in the application of dispersion relations to single-crystal experiments has been the fact that, in this particular case, the function $F(X)$ vanishes in practice for all values of X , except for the lattice points. Consequently, the integrals (1) diverge (Kaufmann, 1985). However, they are valid in the case of scattering by amorphous samples and can be used to constrain the phase values (Makowski, 1981). Only recently, the problem was partially overcome by Mishnev (1993), who, applying Shannon's sampling theorem (Shannon, 1949), derived the following expression for the structure factors:

$$F(h/2) = -(1/j) \sum_{k=-\infty}^{\infty} F(k/2)[1 - \cos \pi(h - k)]/\pi(h - k), \quad (2)$$

where $j = (-1)^{1/2}$, h and k are positive or negative integers and $h \neq k$. If the real and imaginary parts are separated, (2) can be written

$$A(h/2) = -(1/\pi) \sum_{k=-\infty}^{\infty} B(k/2)[1 - (-1)^{h-k}]/(h - k), \quad (3a)$$

$$B(h/2) = (1/\pi) \sum_{k=-\infty}^{\infty} A(k/2)[1 - (-1)^{h-k}]/(h - k). \quad (3b)$$

Equations (3a) and (3b) constitute the crystallographic counterpart of (1a) and (1b) and they relate the imaginary to the real part of the structure factor and *vice versa*. Note that half-integer indices are introduced in order to overcome the problem of discrete functions. The physical reason for this is to be found in the causal transform condition: the Shannon sampling theorem would allow the reconstruction of $F(X)$ with integer sampling points only (crystallographic sampling) by defining $f(x)$ in the interval $-a/2, +a/2$ (where a is the

repetition period of the crystalline array). However, in this case, the condition $f(x) = 0$ for $x < 0$, *i.e.* the condition of causal Fourier transform, would not be fulfilled. The use of twice the crystallographic sampling rate becomes a necessary condition for the application of the Hilbert transforms to structure factors [more complete discussion of the validity conditions of (3) can be found in Mishnev (1993)].

In practice, the real and imaginary components of structure factors with half-integer indices are given by the sum of the imaginary and real components, respectively, of structure factors with integer indices and *vice versa*. Consequently, the real part of a structure factor with integer index \mathbf{h} can only be calculated if we know the imaginary components of all the reflections with half-integer indices ($\mathbf{h}/2$), which is not the case. For this reason, the relationships (3) seem at first glance of little use because we cannot measure reflections with non-integer indices.

In three dimensions, relationship (2) can be written*

$$F(\mathbf{h}/2) = -j \sum_{k_1} \sum_{k_2} \sum_{k_3} F(\mathbf{k}/2) \times \prod_{i=1}^3 [1 - \cos \pi(h_i - k_i)]/\pi(h_i - k_i), \quad (4)$$

where $\mathbf{h} \equiv (h_1, h_2, h_3)$ and $\mathbf{k} \equiv (k_1, k_2, k_3)$.

3. The phase-extension and -improvement procedure

To understand the meaning of coefficients with half-integer indices, let us look at the behaviour of the real and imaginary components of the structure factor in one dimension, using as an example a one-dimensional structure of ten atoms. The values of $A(h)$ and $B(h)$, for $0 < h < 22$, are listed in a graphical form in Figs. 1(a) and 2(a) for integer and half-integer reflections, respectively. Whilst structure factors were calculated *via* the classical formula using the atomic coordinates, coefficients for half-integer reflections were calculated from the former using (3a) and (3b). Of course, if we now recalculate the structure factors, using the same relationship, we will obtain the starting values again, except for a small approximation due to the truncation error. A relevant point is illustrated in Figs. 1(b) and 2(b): if some reflections are deleted from the starting set, the half-integer coefficients calculated with the reduced set keep the same general behaviour. The reason for this is apparent from (3a) and (3b): in the calculation of $A(h/2)$, each term in the sum is divided by $(h - k)$ and the relevant terms will presumably be those with h not too different from k and/or those with a large value of $B(k)$. All others will be negligible in practice.

* This is the final equation in Mishnev (1993), except for a constant factor π^2 .

This is more evident in Figs. 3(a) and (b), where the real and imaginary parts of the structure factors were calculated using only the real or imaginary part of the highest terms of coefficients with half-indices. In practice, only the imaginary parts of the coefficients that have modulus greater than 5 were used in the sum (3a) and the real parts in the sum (3b). Comparison with

the correct values, illustrated in Fig. 3, shows a surprisingly good agreement.

These results suggest that, if the phases of an appropriate number of reflections are known, it could be possible to predict the remaining ones. Besides, it is possible to modify or select the half-index coefficients to be used in calculations, introducing in this way a sort of filter in the transforms. In the following, we describe how these principles were used in order to extend or improve phases.

The considerations discussed above apply to the three-dimensional case in a similar way. Using (4) and separating the real and imaginary parts, one can write

$$A(\mathbf{h}/2) = -(1/\pi^3) \sum_{k_1} \sum_{k_2} \sum_{k_3} B(\mathbf{k}/2) \times \prod_{i=1}^3 [1 - (-1)^{h_i - k_i} / (h_i - k_i)] \quad (5a)$$

$$B(\mathbf{h}/2) = (1/\pi^3) \sum_{k_1} \sum_{k_2} \sum_{k_3} A(\mathbf{k}/2) \times \prod_{i=1}^3 [1 - (-1)^{h_i - k_i} / (h_i - k_i)]. \quad (5b)$$

For clarity, let us use \mathbf{h} to indicate a triplet of even integers and \mathbf{k} one of odd integers: we can now rewrite relationships (5) in a more convenient way [note that, under these conditions, $h_i - k_i$ for (6a), (6b) and $k_i - h_i$ for (7a), (7b) are always odd]:

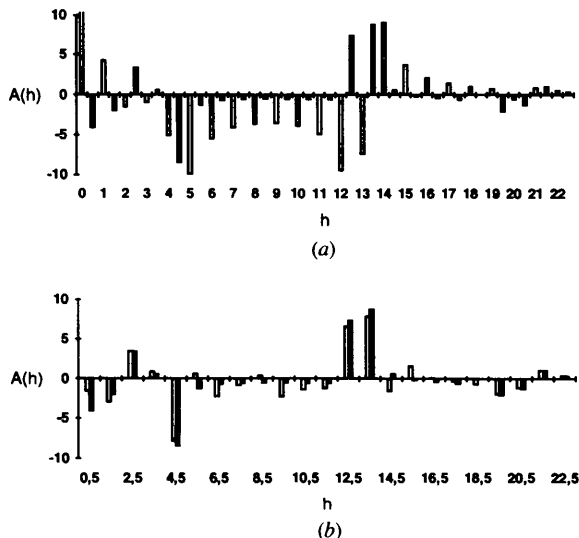


Fig. 1. Real components of the structure factors for a one-dimensional arrangement of atoms. (a) Calculated values of structure factors (i.e. those with integer indices) are represented by white bars. Black bars represent coefficients with half-integer indices, calculated from the previous ones using formula (3a). The value of $A(0)$ is truncated. (b) Comparison of the real part of the coefficients with half-integer indices: those calculated from a complete and correct starting set are in white, those calculated from an incomplete starting set, where reflections 2, 5, 7, 15 were deleted, in black.

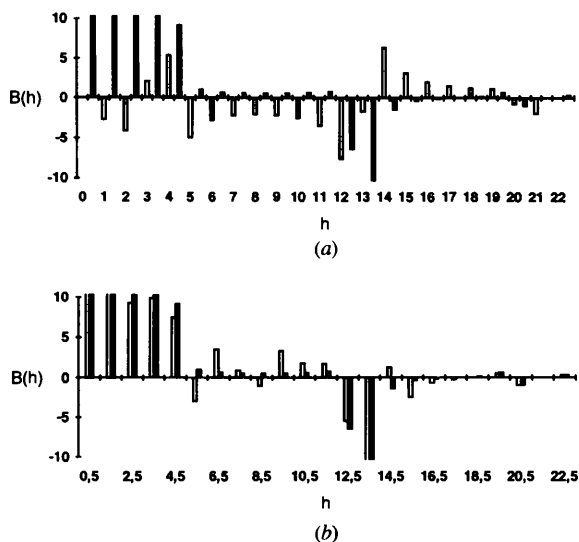


Fig. 2. Same as Fig. 1 for the imaginary components of the structure factors.

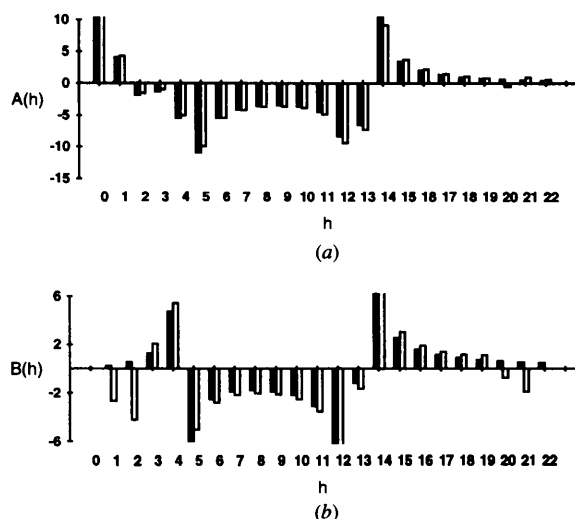


Fig. 3. (a) Real and (b) imaginary components of the structure factors. While bars represent the components calculated with a complete set of half-integer indices coefficients using relationships (3). Black bars show the same coefficients calculated using only a portion of those with half-integer indices, i.e. those with modulus greater than 5. It appears that the exclusion of the smaller terms in the sum does not significantly affect the values of $A(h)$ and $B(h)$; only for reflections 1, 2 and 20, 21 will the phases be affected by a relevant error.

$$A(\mathbf{h}/2) = -(8/\pi^3) \sum_{k_1} \sum_{k_2} \sum_{k_3} B(\mathbf{k}/2) \prod_{i=1}^3 1/(h_i - k_i) \quad (6a)$$

$$B(\mathbf{h}/2) = (8/\pi^3) \sum_{k_1} \sum_{k_2} \sum_{k_3} A(\mathbf{k}/2) \prod_{i=1}^3 1/(h_i - k_i) \quad (6b)$$

$$A(\mathbf{k}/2) = -(8/\pi^3) \sum_{h_1} \sum_{h_2} \sum_{h_3} B(\mathbf{h}/2) \prod_{i=1}^3 1/(k_i - h_i) \quad (7a)$$

$$B(\mathbf{k}/2) = (8/\pi^3) \sum_{h_1} \sum_{h_2} \sum_{h_3} A(\mathbf{h}/2) \prod_{i=1}^3 1/(k_i - h_i). \quad (7b)$$

If we assume that the $A(\mathbf{h})$ and $B(\mathbf{h})$ components of the structure factors with integer indices are known, or at least a fraction of them, we obviously can use (7a) and (7b) to calculate $A(\mathbf{k})$ and $B(\mathbf{k})$ of hypothetical reflections with half-indices. From these, in turn, it is possible, using (6a) and (6b), to recalculate the original $A(\mathbf{h})$ and $B(\mathbf{h})$. Sums in (6) and (7) extend from $-\infty$ to $+\infty$, but the high-resolution terms will presumably only make a small contribution to the sum since they are in general smaller than the low-resolution terms.

3.1. The phase-extension process

Let us divide the set of reflections into two subsets: $\{\Gamma_h\}$ will be the subset of reflections with known phases and $\{\Omega_h\}$ the fraction for which phases are unknown. Via (7), we can use subset $\{\Gamma_h\}$ to obtain approximate values of $A(\mathbf{k}/2)$ and $B(\mathbf{k}/2)$ for the entire set, $\{\Gamma_k + \Omega_k\}$; from these, we can subsequently, by using (6), recalculate the entire set of reflections $\{\Gamma_h + \Omega_h\}$. It is important to note that we obtain new phases for all reflections but since we know the moduli we can substitute them with the newly obtained ones and start a new phase calculation cycle. We can summarize the steps in the process of phase extension as follows:

(i) Using the set of reflections with known phases $\{\Gamma_h\}$ (it can be a more or less complete set at limited resolution or simply a set made up of a limited number of strong phased reflections, as happens with direct methods), calculate $A(\mathbf{k}/2)$ and $B(\mathbf{k}/2)$ via (7) for all reflections to a resolution R_1 . Consider that for this calculation the value of $F(000)$ is needed, which, having an intensity much larger than that of all other reflections, plays an important role in determining the phases of the very low resolution reflections.

(ii) From $A(\mathbf{k}/2)$ and $B(\mathbf{k}/2)$, calculate $A(\mathbf{h}/2)$ and $B(\mathbf{h}/2)$, i.e. the set $\{\Gamma_h + \Omega_h\}$, up to a resolution R_2 . Empirical tests have shown that $R_1 > R_2$ helped in reducing the truncation effects in the sums.

(iii) Renormalize the values of the newly calculated reflections based on their known moduli $|F(\mathbf{h})|$, using $|F(\mathbf{h})|^2 = A(\mathbf{h})^2 + B(\mathbf{h})^2$. Now repeat the process from step (i).

Modifications can be introduced in step (ii) and eventually in step (i). For example, one can limit the calculation to the stronger reflections or to classes of reflections that are believed to be more accurate. It is also possible to introduce weights in the sum. These modifications, along with the fact that the experimental information contained in structure-factor moduli is preserved through the macrocycles, are the basis for the convergence of the procedure. These constraints, however, are not necessarily sufficient to ensure convergence in all cases: some examples of their possible use in some successful situations are discussed in §4.

In the following, these three steps together will be called one macrocycle.

3.2. The phase-improvement process

Let us consider the realistic situation where a complete set of phases is known to a given resolution but these phases are affected by experimental errors. In the case of a protein structure determined by the multiple isomorphous replacement method, to each phase is assigned a figure of merit (m) as a weighting factor, which indicates the reliability of the phase itself. Phases with a figure of merit m of 1.0 are in principle expected to be correct, while those with a figure of merit close to 0 completely wrong. To introduce this information in the calculations, half-integer coefficients are calculated by a modified version of (7):

$$A(\mathbf{k}/2) = -(8/\pi^3) \sum_{h_1} \sum_{h_2} \sum_{h_3} mB(\mathbf{h}/2) \prod_{i=1}^3 1/(k_i - h_i) \quad (8a)$$

$$B(\mathbf{k}/2) = (8/\pi^3) \sum_{h_1} \sum_{h_2} \sum_{h_3} mA(\mathbf{h}/2) \prod_{i=1}^3 1/(k_i - h_i). \quad (8b)$$

In this way, each coefficient in the sum is weighted according to its reliability. In the following cycles, new phases are obtained. Of course, we want to modify only the wrong phases and keep the correct ones. The last can be distinguished by their figure of merit, which has given them a strong weight in the sum: the calculated phases for reflections with a high m should therefore be similar to the previous ones. To take these facts into account, in every macrocycle the new phase, φ_N , was obtained as the weighted mean of the starting phase, φ_s , and the newly calculated one, φ_c :

$$\varphi_N = m\varphi_s + (1 - m)\varphi_c. \quad (9)$$

The procedure is then repeated from the beginning, calculating half-integer coefficients from (8) and structure factors using (6).

Table 1. Fractional atomic coordinates of the small molecule (a phenylalanyl residue) used as a test case

An isotropic B factor of 7 \AA^2 was used for all atoms. The molecule was positioned in a $P1$ crystal cell, with parameters $a = 12$, $b = 6$, $c = 4 \text{ \AA}$, $\alpha = \beta = \gamma = 90^\circ$.

	x	y	z
N	0.25742	0.51183	0.73150
$C\alpha$	0.37992	0.51133	0.73900
$C\beta$	0.42550	0.39117	0.43050
$C\gamma$	0.55292	0.38400	0.42100
$C\delta_1$	0.61483	0.48167	0.67900
$C\epsilon_1$	0.73142	0.47500	0.67025
$C\zeta$	0.78608	0.37083	0.40350
$C\epsilon_2$	0.72417	0.27317	0.14550
$C\delta_2$	0.60758	0.27983	0.15425
C	0.42375	0.75067	0.74150
O	0.52525	0.78967	0.74775

4. Results and discussion

4.1. Test case I. A small-molecule crystal

Structure factors were calculated for a hypothetical molecule of ten atoms (a phenylalanyl residue) in a $P1$ crystal cell, as reported in Table 1. Two different kinds of test were performed: in the first, the data were considered to be error free but only a limited set of phases was assumed to be known and phases were extended to the remaining reflections; in the second, all the data were used in the calculations but random errors were introduced in the starting set and the procedure aimed at improving the existing phases. In both cases, only reflections within a resolution limit of 0.95 \AA were used.

(a) *Phase extension.* In this test, the n th reflection was deleted from the starting set and all statistics refer only to the deleted fraction.* Only the results for $n = 2$ and $n = 5$ will be discussed here. For $n = 5$, the subset of reflections $\{\Omega_h\}$ was obtained by deleting 20% of the reflections in a random way. After application of relationships (7) and (6), the mean overall phase error† for subset $\{\Omega_h\}$ was 54° . At this point, only reflections of subset $\{\Omega_h\}$ whose modulus was greater than a fixed threshold were used to calculate the half-integer coefficients, according to (7), along with all the reflections of subset $\{\Gamma_h\}$. A similar criterion was used to recalculate integer coefficients: only half-integer terms greater than a given threshold were used in (6). This was suggested by an examination of the statistics after every cycle, which showed a strong negative correlation between the phase error and the value of the modulus. The threshold cut-off values were decreased

* Deleting a reflection means that its value is not considered in the sum in the first macrocycle but the value of its modulus is kept and used from the second macrocycle onwards.

† The mean phase error is defined, with respect to the 'true' phase φ_T , as $MPE = \sum_i (PD)_i / N$, where $PD = |\varphi_i - \varphi_T|$ and, if PD is greater than 180° , it is set equal to $360^\circ - PD$.

at every cycle so that the entire set was used in the last cycles. The phase refinement converges after 12 cycles and the mean overall phase error for extended reflections is 30.5° . The distribution of errors as a function of moduli (Fig. 4a) shows that for reflections with $|F(\mathbf{h})| > 1.0$ the uncertainty is about 23° .*

For $n = 2$, the results of the phase prediction after the first cycle are given in Fig. 4(b). The error is 69° but the distribution of the phase error as a function of the modulus $|F(\mathbf{h})|$ shows that it is smaller for the highest reflections [for $|F(\mathbf{h})| > 3.1$, the mean error is about 39°]. This is also true for half-index coefficients. Therefore, only coefficients with a modulus higher than a given threshold were used to calculate the half-integer values of step (i), in order to speed up the convergence of the procedure, and the same was done for the half-integer coefficients used in step (ii). The limiting threshold was decreased at each cycle so that all reflections were used in the calculations of cycle 8. The final result is illustrated in Fig. 4(b), where the mean phase errors are reported as a function of moduli of structure factors. Despite the relatively large overall error, the highest reflections [$|F(\mathbf{h})| > 3.1$] have a mean error of only about 23° .

* Reflections were calculated on an absolute scale and $F(000) = 69.2$.

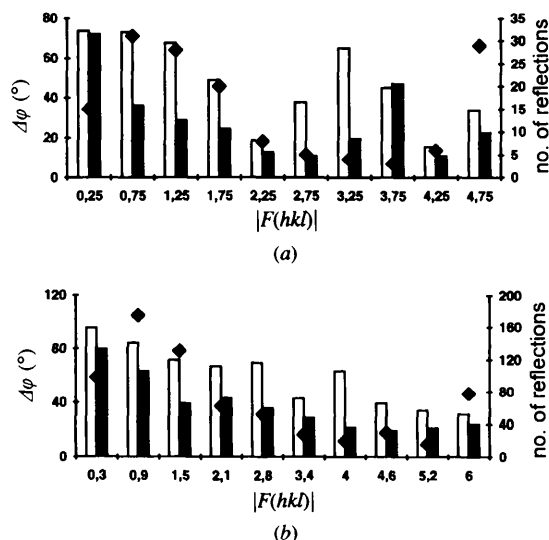


Fig. 4. Mean phase error, as a function of modulus, in the case of the phase-extension process for a small test molecule. The mean phase error is reported on the left, the number of reflections on the right. The abscissa gives the mean value of the modulus within the interval. Black diamonds represent the number of reflections in the interval. (a) 20% of reflections were selected in a random way and their phases were deleted from the starting set. Phase errors were calculated only for the 149 reflections with 'unknown' phase. White bar: mean phase error after the first cycle; grey bar: mean phase error after 12 cycles. (b) Same as (a), except that now the phase of every second reflection was deleted from the starting set in a systematic way. White bar: mean error after the first cycle; grey bar: mean phase error after 50 cycles.

(b) *Phase improvement.* Two different tests were performed. In the first, a mean error $\Delta\varphi$ of $\pm 30^\circ$ on the phases of reflections was introduced, in the second $\pm 70^\circ$. Initial errors were introduced by using a random-number generator to produce numbers between -0.5 and 0.5 and by adding the resulting values (labelled rand), multiplied by $4\Delta\varphi$, to each phase. A figure of merit was assigned to every reflection according to the error introduced for the relative phase, *i.e.* $m = 1 - 2|\text{rand}|$. Half-integer coefficients were then calculated using (8). In this way, the terms with the lowest phase error obtained a higher weight and produced a set of coefficients that gave, using (6), structure factors with an improved phase. In this case, the prediction of the moduli was relatively poor but it is not important because they are known from the experimental work.

In both cases, after the first cycle, the phase error decreased consistently for all reflections, even for those with a small modulus, and the improvement continued in successive cycles (from 30 to 19° in the first case). In Fig. 5, the distribution of the phase error as a function of the moduli is presented; the starting distribution is

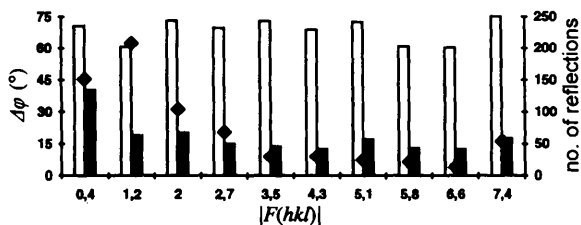


Fig. 5. Test of phase improvement for a small-molecule crystal. A mean phase error of $\pm 70^\circ$ was introduced in a random way on the original set of reflections. The mean phase error in the starting set as a function of the reflection moduli is shown by white bars, the mean phase error after 10 cycles by grey bars. Black diamonds represent the number of reflections in the interval (scale on the right).

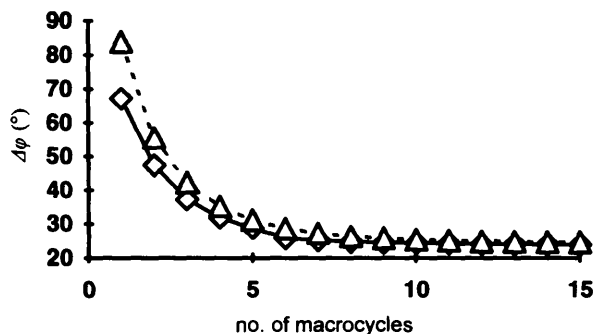
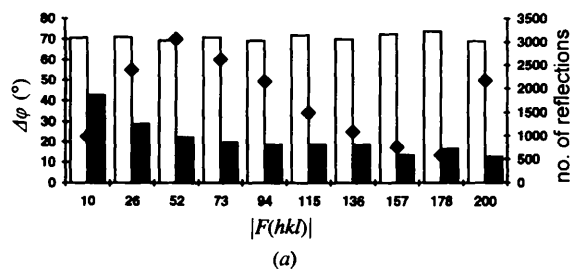


Fig. 6. Mean overall phase error as a function of the number of macrocycles for phase improvement in a hypothetical protein crystal. Two tests are shown: in one, a random error of $\pm 90^\circ$ was introduced (triangles), in the other a random error of $\pm 70^\circ$ (diamonds). In both cases, the errors converge to similar values, around $\pm 24^\circ$, and convergence is reached in about ten cycles.

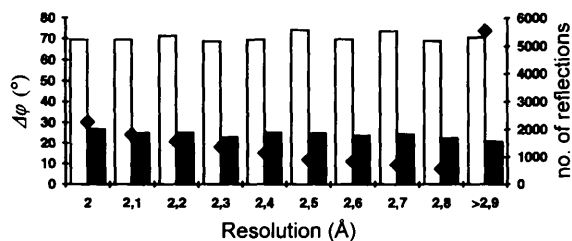
compared with that obtained after ten cycles of calculation and the improvement is clearly visible from the figure: in the second test case, the mean phase error for the starting set of reflections is $\pm 68^\circ$, the final one $\pm 22^\circ$, with an improvement of nearly 46° .

4.2. Test case II. A small protein

A hypothetical protein crystal was built up by using the atomic coordinates of human muscle fatty-acid-binding protein, MFABP, a single-chain protein of 132 amino acids, with 1027 atoms (Zanotti, Scapin, Spadon, Veerkamp & Sacchettini, 1992). The molecular model was positioned in a triclinic $P1$ cell, with $a = b = c = 40 \text{ \AA}$ and $\alpha = \beta = \gamma = 90^\circ$, and data were calculated at 2 \AA resolution. As in the case of the small-molecule crystal, a random error on the phases was introduced and the figure of merit was evaluated accordingly. From a starting mean overall phase error of ± 70 and $\pm 90^\circ$, mean errors of ± 23.8 and $\pm 24.5^\circ$, respectively, were obtained after 15 macrocycles. The mean phase error as a function of the number of macrocycles is shown in Fig. 6: the procedure approaches convergence in a few cycles with an exponential behaviour. Interestingly, the method can reduce the phase error to an acceptable level (in the perspective of Fourier-map interpretation) even when the starting data set is a very low quality one. The phase-error distribution as a function of structure-factor moduli and resolution is presented in Fig. 7. It is quite



(a)



(b)

Fig. 7. Test of phase improvement for a simulated protein crystal. After introducing a mean phase error of $\pm 70^\circ$ in a random way, phases were improved by 20 macrocycles of refinement, as described in the text. The final mean error is $\pm 23.5^\circ$. The behaviour of the mean phase error as a function of modulus and of resolution is reported in (a) and (b), respectively. White bars show the starting mean phase errors, grey bars the final ones, black diamonds the number of reflections in the interval.

similar to that obtained for the small molecule. The mean error increases for smaller moduli, as expected. The effect of resolution is surprisingly low: the error distribution is relatively flat, indicating that truncation errors are negligible.

In order to show the effect of the process on map quality, in Fig. 8 three maps are compared: (a) the 'correct' one, based on calculated phases; (b) the map with starting phases, *i.e.* $\pm 70^\circ$ of mean phase random error; (c) the map calculated from the phase set of (b) after 20 macrocycles of HT refinement (mean phase error $\pm 23.5^\circ$). In map (b), it is not even possible to discern the shape of the molecule, while map (c)

shows nearly all the structural information contained in (a).

4.3. Test case III. Experimental data of a protein crystal

The real crystal structure of MFABP belongs to space group $P2_12_12_1$ ($a = 35.4$, $b = 56.7$, $c = 72.7$ Å) and has been solved and refined to a resolution of 2 Å. Data for the native protein were collected on a Siemens X1000 area detector system coupled to a Rigaku RU-200 rotating-anode X-ray generator. 35 568 reflections up to 2.1 Å resolution were measured and

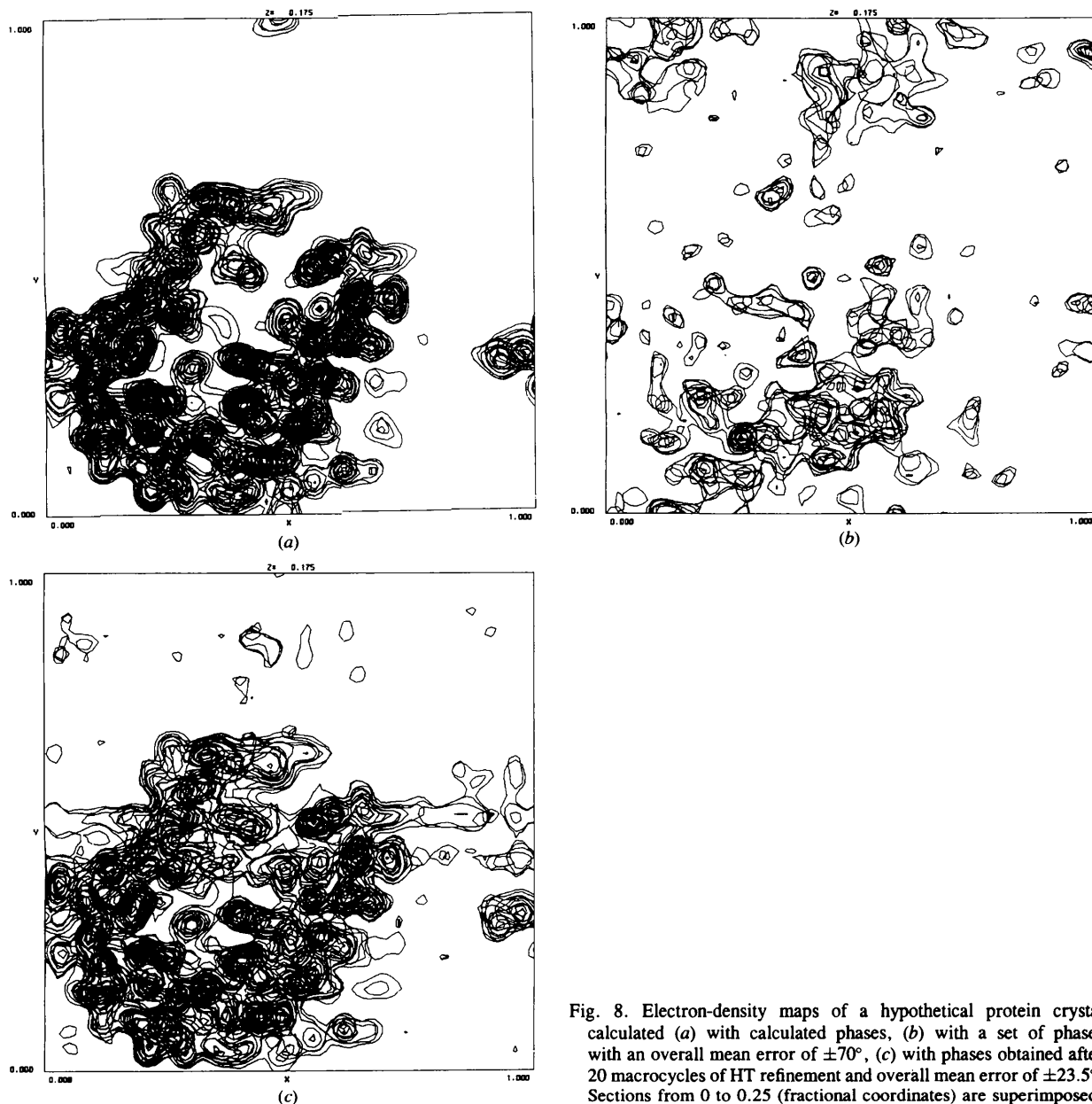


Fig. 8. Electron-density maps of a hypothetical protein crystal calculated (a) with calculated phases, (b) with a set of phases with an overall mean error of $\pm 70^\circ$, (c) with phases obtained after 20 macrocycles of HT refinement and overall mean error of $\pm 23.5^\circ$. Sections from 0 to 0.25 (fractional coordinates) are superimposed.

reduced to 7576 independent reflections with an overall R_{merge} on intensity of 0.052 (Zanotti *et al.*, 1992). Phases were calculated from the refined model, which includes 1027 protein atoms, 56 water molecules and 16 C atoms of a fatty acid chain. As in the previous cases, a random error on phases was introduced and the figure of merit was evaluated accordingly. No attempt was made to put the data on an absolute scale and no σ cut-off was applied. Reflections absent from the data set (the completeness to 2.1 Å is about 90%) were simply ignored and not considered in the statistics. From a starting mean overall phase error of $\pm 70^\circ$, a mean error of $\pm 33.3^\circ$ was obtained after ten macrocycles. The mean overall error is about 9° greater than that obtained with the simulated data of the previous example, starting from the same initial error. This behaviour can be understood by looking at the phase-error distribution as a function of moduli (Fig. 9a), which is quite different from that obtained in the previous cases: the mean phase error, which is around $\pm 19^\circ$ for $|F_{\text{obs}}|$ greater than 219, drastically increases to $\pm 72^\circ$ for $|F_{\text{obs}}|$ less than 10. Evidently, experimental errors on intensities play a role, which is particularly relevant for very weak reflections. The effect on resolution is also significant (Fig. 9b), but this is probably attributable to the previous fact rather than to the effect of truncation. Despite this, the phase improvement is highly significant and the inaccuracy in measurement of the data does not affect the procedure overall.

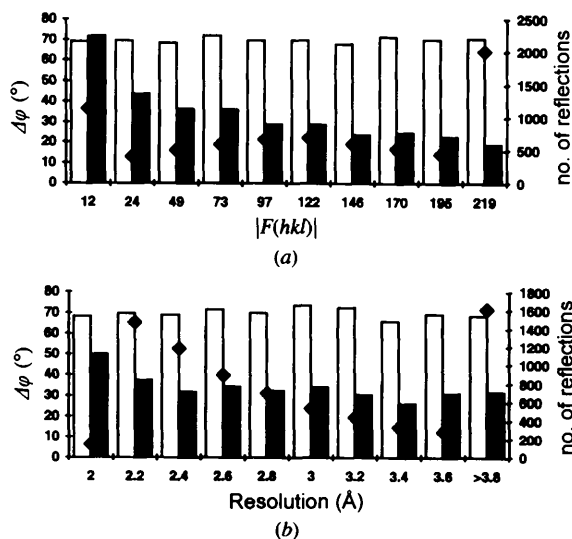


Fig. 9. Test of phase improvement for a real protein crystal (space group $P2_12_12_1$, $Z = 4$). With experimental moduli and phases calculated from the molecular model, a mean phase error of $\pm 70^\circ$ on phases was introduced in a random way. Phases were improved by 10 macrocycles of refinement, as described in the text. The final mean error is $\pm 33.3^\circ$. The behaviour of mean phase error as a function of modulus and of resolution is reported in (a) and (b), respectively. White bars show the starting mean phase errors, grey bars the final ones, black diamonds the number of independent reflections in the interval.

5. Conclusions

Two different tasks were considered in this paper: (a) starting from an incomplete set of correct phases within a resolution sphere, to calculate the values of the remaining ones; (b) to improve a set of phases affected by a given error, assuming that correct estimates of errors are available. The degree of success in the first case strongly depends on the completeness of the starting set of phases: if most phases are known, the remaining ones can be predicted with a small error (which increases with the incompleteness of the starting set). In any case, at least about 50% of the possible phases must be known, in order to obtain a reliable estimate of the remaining set. Another limitation of the method is presented by the distribution of the known phases: in fact, the phase of reflection h mostly depends, in the present approach, on the reflections with indices close to it: the procedure seems therefore more suitable for the extension of randomly unphased reflections than of an entire shell of resolution.

For the phase-improvement task, the critical point seems to be the weighting in the sums: the procedure can distinguish between 'good' and 'bad' terms (and use the former to improve the latter) only if a reliable figure of merit can be calculated. This was not a problem in our tests with calculated phases but it will play a fundamental role in the extension of a MIR or SIR data set. This task will rely on the estimates of the protein phase probability, which are calculated from the derivative data: the related figures of merit may suffer from systematic errors. Our present activity is therefore aimed at assessing how well the described HT procedure behaves in the common situation when experimental phases are available, but they need to be improved in order to attain an interpretable electron-density map.

Finally, the use of discrete Hilbert transforms is based on the so-called 'causality' condition and no assumptions need to be made on electron-density characteristics, as for atomicity or positivity in classical direct methods; nor on the molecular model, like in some density-modification procedures. For that reason, the present approach appears to be complementary to the previous ones and could eventually be combined with them to help in solving the phase problem.

We thank Professor S. Ciccariello and Dr A. Cervellino for valuable discussions. We are grateful to Professor J. N. Jansonius for generously allowing GC to devote part of his research time to this project and for reading the manuscript.

References

- Bath, N. T. & Blow, D. M. (1982). *Acta Cryst.* **A38**, 21–29.
- Bricogne, G. (1988). *Acta Cryst.* **A44**, 517–545.

- Bryan, R. K. & Banner, D. W. (1987). *Acta Cryst.* **A43**, 556–564.
- Burge, R. E., Fiddy, M. A., Greenaway, A. H. & Ross, G. (1976). *Proc. R. Soc. London Ser. A*, **350**, 191–212.
- Cannillo, E., Oberti, R. & Ungaretti, L. (1983). *Acta Cryst.* **A39**, 68–74.
- Davies, A. R. & Rollet, J. S. (1976). *Acta Cryst.* **32**, 17–23.
- Giacovazzo, C. (1980). *Direct Methods in Crystallography*. London: Academic Press.
- Joo, T. & Albrecht, A. C. (1993). *J. Phys. Chem.* **97**, 1262–1264.
- Kaufmann, B. (1985). *Acta Cryst.* **A41**, 152–155.
- Makowski, L. (1981). *J. Appl. Cryst.* **14**, 160–168.
- Mishnev, A. F. (1993). *Acta Cryst.* **A49**, 159–161.
- Ramachandran, G. N. (1969). *Mater. Res. Bull.* **4**, 525–534.
- Refaat, L. S. & Woolfson, M. M. (1993). *Acta Cryst.* **D49**, 367–371.
- Schevitz, R. W., Podjarny, A. D., Zwick, M., Hughes, J. J. & Sigler, P. B. (1981). *Acta Cryst.* **A37**, 669–677.
- Shannon, C. E. (1949). *Proc. Inst. Radio Eng. NY*, **37**, 10–41.
- Shiono, M. & Woolfson, M. M. (1992). *Acta Cryst.* **A48**, 451–456.
- Toll, J. S. (1956). *Phys. Rev.* **104**, 1760–1770.
- VanderNoot, T. J. (1992). *J. Electroanal. Chem.* **332**, 9–24.
- Wang, B. C. (1985). *Methods Enzymol.* **115**, 90–112.
- Williams, C. P. & Marshall, A. (1992). *Anal. Chem.* **64**, 916–923.
- Xiang, S., Carter, C. W., Bricogne, G. & Gilmore, C. J. (1993). *Acta Cryst.* **D49**, 193–212.
- Zanotti, G., Scapin, G., Spadon, P., Veerkamp, J. H. & Sacchettini, J. C. (1992). *J. Biol. Chem.* **267**, 18541–18550.