# Simulated Annealing for Phasing using Spatial Constraints

BY P. BÉRAN* AND A. SZÖKE

*University of California, Lawrence Livermore National Laboratory, PO Box 808, Livermore, CA 94550, USA*

## Abstract

Phase refinement using partial structural information is an important step in the treatment of crystallographic data. An optimization procedure is presented that, in contrast to conventional Fourier methods, takes into account the spatial dependence of the accuracy of the partial information for electron density. Annealing studies for a test case show that this optimization procedure is tractable for system sizes relevant to the treatment of protein diffraction data and requires substantially less structural information for accurate phase refinement than Fourier methods.

## 1. Introduction

The Fourier method has been known since the early days of X-ray crystallography (Cochran, 1951) and is ubiquitous in techniques such as solvent flattening (Wang, 1985; Cura, Krishnaswamy & Podjarny, 1992) or molecular replacement (Rossmann & Blow, 1962; Fitzgerald, 1991) used nowadays to phase protein diffraction data. In this method, a partial knowledge of the electron density is used to estimate the phase of the structure factors of the whole structure. However, it often occurs in practice that the partial knowledge of the electron density is insufficient to lead to accurate estimates for the phase of the structure factors and the resulting electron-density map is difficult or impossible to interpret. Probabilistic approaches (Sim, 1959; Srinavasan, 1966; Bricogne, 1976; Read, 1986) have been suggested to improve the efficiency of the Fourier method in the case of a limited partial knowledge of the electron density.

It has recently been suggested (Szöke, 1993) that an efficient use of the partial knowledge of the electron density could be achieved in cases where one has an estimate of the electron density that is accurate in a restricted region of the unit cell. In such cases, the dependence of the accuracy of the estimate on the position within the unit cell can be considered as an additional piece of information that may help

to improve phase recovery. The Fourier method unfortunately does not take this type of information into account. This type of information can be included by using an optimization procedure (Szöke, 1993) in which the electron density is simultaneously constrained to be consistent both with the knowledge of the density in the restricted region and with the diffraction data. The idea of using this type of information is interesting because the identification of a protein domain or of a region of a disordered solvent in an imperfect density map often leads to accurate knowledge of the electron density in a restricted region.

We propose here an optimization procedure for phasing using the knowledge of the electron density in a restricted region of the unit cell as a constraint. The amplitudes of the structure factors of the electron density are fixed to their true values, while the phases are used as variational parameters. We define the cost function to be optimized as the average, *restricted to the region of known density*, of the squared difference between the electron density corresponding to the variational phases and the density estimate available in the restricted region. This cost function is then minimized using simulated annealing (Kirkpatrick, Gelatt & Vecchi, 1983).

We study the performance of this procedure for accurate phase refinement on a test case given by the model of a known protein containing 149 residues. In order to characterize the accuracy of phase retrieval, we evaluate the average error of the recovered phases, as well as their sensitivity to noise in data for the density estimate in the restricted region and to noise in data for structure-factor amplitudes.

We find that this procedure requires substantially less structural information for accurate phase refinement than Fourier methods. For example, assuming the electron density to be known with a given uncertainty in a restricted domain occupying 61% of the unit cell and consisting mostly of disordered solvent, this procedure allows an estimation of the electron density in the rest of the unit cell with an error smaller than twice the uncertainty of the density estimate in the restricted domain. In contrast, using the same partial knowledge for the electron density, we find that the Fourier estimates for phases

---

* Present address: Université Paul Sabatier, Laboratoire de Physique Quantique, URA 505, IRSAMC, 118 route de Narbonne, 31062 Toulouse CEDEX, France.

differ from the true values by more than 90° on average. This optimization procedure is practical: it takes *ca* 2 h on an IBM 6000 computer to phase the structure up to 2 Å resolution.

In addition, studies of the relaxation time for small systems suggest that the optimization problem presented here is of moderate complexity when data for the density estimate in the restricted region and data for structure-factor amplitudes are sufficient to lead to unambiguous determination of the whole electron density. When the latter condition is not satisifed, we observe a slowing down in relaxation characteristic of spin glass systems (Mackenzie & Young, 1982).

This procedure treats electron-density data for protein structure and disordered solvent on the same footing. It can thus make an efficient use of knowledge for the position of parts of disordered solvent, in contrast to Fourier methods which require the knowledge of the whole protein envelope.

In the case where the amount of structural information is insufficient to lead to accurate phase refinement, it frequently occurs that the density map obtained by Fourier methods still allows the identification of additional parts of the structure, leading to the complete determination of the protein structure in several steps. The question of the performance of the optimization procedure in such a case deserves further study. In particular, it would be interesting to know whether or not this procedure is able to lead to more interpretable density maps than Fourier methods in the case of a limited partial knowledge of the electron density.

## 2. Optimization procedure

### 2.1. *Cost function*

Let us state the problem for the simple case of an orthorhombic crystal. The structure factors $F_k$ and the electron density $\rho(\mathbf{r})$ are related by

$$F_k = \int_{\mathrm{u.c.}} d^3r\rho(\mathbf{r})\exp(2\pi i\{[(k_1r_1)/a_1] + [(k_2r_2)/a_2]$$

$$+ [(k_3r_3)/a_3]\})\qquad(2.1)$$

and by

$$\rho(\mathbf{r}) = V^{-1}\sum_k F_k\exp(-2\pi i\{[(k_1r_1)/a_1] + [(k_2r_2)/a_2]$$

$$+ [(k_3r_3)/a_3]\}),\qquad(2.2)$$

where the integral in (2.1) is performed over the unit cell, $V$ denotes the volume of the unit cell and $a_1$, $a_2$ and $a_3$ denote the unit-cell dimensions. The amplitudes of structure factors $F_k$ can be inferred from diffraction data but their phases are unknown. Let us suppose that we are given an estimate $\rho^{\mathrm{est}}(\mathbf{r})$ of the electron density, which is accurate only in a restricted region $\mathscr{D}$ of the unit cell. We wish to reconstruct the electron density in the whole unit cell by simul-

taneously requiring $\rho(\mathbf{r})$ to be consistent with diffraction data for $|F_k|$ and to match $\rho^{\mathrm{est}}(\mathbf{r})$ in a restricted region $\mathscr{D}$.

Let us first suppose that $\rho(r)$ and $\rho^{\mathrm{est}}(\mathbf{r})$ are bandwidth limited, *i.e.* their structure factors $F_k$ and $F_k^{\mathrm{est}}$ vanish for $\lambda_k/2 < R$, where $R$ is the resolution limit and where

$$\lambda_k = [(k_1^2/a_1^2) + (k_2^2/a_2^2) + (k_3^2/a_3^2)]^{-1/2}\qquad(2.3)$$

is the wavelength associated with wave vector $\mathbf{k}$.

The knowledge of the amplitude of structure factors can be incorporated by expressing the electron density as

$$\rho_{\varphi_k}(\mathbf{r}) = V^{-1}\sum_k{}'A_k\exp(i\varphi_k)\exp(-2\pi i\{[(k_1r_1)/a_1]$$

$$+ [(k_2r_2)/a_2] + [(k_3r_3)/a_3]\}),\qquad(2.4)$$

where $A_k$ and $\varphi_k$, respectively, denote the amplitude and phase angle of structure factor $F_k$ and where the primed sum is performed over index $k$ satisfying $\lambda_k/2 \geq R$. Subscript $\varphi_k$ on the left-hand side of (2.4) indicates that the density is parametrized by the finite set of phase angles $\varphi_k$ satisfying $\lambda_k/2 \geq R$. We wish to fit the phase parameters $\varphi_k$ so that $\rho_{\varphi_k}(\mathbf{r})$ matches $\rho^{\mathrm{est}}(\mathbf{r})$ in domain $\mathscr{D}$. For this, we minimize the cost function $E(\varphi_k)$ given by

$$E(\varphi_k) = (N_1N_2N_3)^{-1}\sum_{n_\alpha = 0}^{N_\alpha - 1}\chi_n|\rho_{\varphi_k}(\mathbf{r}_n) - \rho^{\mathrm{est}}(\mathbf{r}_n)|^2,$$

$$(2.5)$$

where $N_1$, $N_2$ and $N_3$ are positive integers satisfying $a_\alpha/N_\alpha < R$ for $\alpha = 1$, 2 and 3; the sum $n_\alpha = 0$, ..., $N_\alpha - 1$ is performed over all three components of $\mathbf{n} = (n_1, n_2, n_3)$; $\mathbf{r}_n$ defines a grid spanning the unit cell in the manner

$$\mathbf{r}_n = \left(\frac{n_1}{N_1}a_1, \frac{n_2}{N_2}a_2, \frac{n_3}{N_3}a_3\right);\qquad(2.6)$$

$\chi_n$ is the characteristic function of the domain $\mathscr{D}$ on lattice $\mathbf{r}_n$ defined by

$$\chi_n = \begin{cases} 1 & \text{if } \mathbf{r}_n \in \mathscr{D} \\ 0 & \text{otherwise.} \end{cases}\qquad(2.7)$$

The cost function $E(\varphi_k)$ can be explicitly written as a function of phases $\varphi_k$ as

$$E(\varphi_k) = \sum_{k,q}{}'\hat{\chi}_{q-k}A_kA_q\exp[i(\varphi_k - \varphi_q)]$$

$$- 2\mathrm{Re}\left[\sum_{k,q}{}'\hat{\chi}_{q-k}F_q^{\mathrm{est}*}A_k\exp(i\varphi_k)\right]$$

$$+ \sum_{k,q}{}'\hat{\chi}_{q-k}F_k^{\mathrm{est}}F_q^{\mathrm{est}*},\qquad(2.8)$$

where $\hat{\chi}_k$ is the discrete Fourier transform of $\chi_n$

given by

$$\hat{\chi}_k = \frac{V^{-2}}{N_1 N_2 N_3} \sum_{n_\alpha = 0}^{N_\alpha - 1} \chi_n \exp(2\pi i\{[(k_1 n_1)/N_1]$$

$$+ [(k_2 n_2)/N_2] + [(k_3 n_3)/N_3]\}). \qquad (2.9)$$

We now discuss the relation between the Fourier method and the procedure which consists of minimizing the cost function given by (2.5). The phases $\varphi_k$ obtained by conventional Fourier methods on the basis of the estimated density $\rho^{est}(\mathbf{r})$ can be expressed as the phase configuration which minimizes the cost function $E(\underline{\varphi_k})$ when the domain $\mathcal{D}$ occupies the whole unit cell. Indeed, in this case, (2.8) can be rewritten as

$$E(\underline{\varphi_k}) = V^{-2} \sum_k{}' |A_k \exp(i\varphi_k) - F_k^{est}|^2. \qquad (2.10)$$

This quantity is minimal when $\varphi_k$ is equal to the phase angle of $F_k^{est}$. The phases obtained by Fourier methods thus depend on the value of the estimate $\rho^{est}(\mathbf{r})$ at every position $\mathbf{r}$ in the unit cell, whether or not $\rho^{est}(\mathbf{r})$ is accurate at position $\mathbf{r}$. In the case where we have an estimate of the electron density that is accurate only in a restricted region $\mathcal{D}$, the definition of $\rho^{est}(\mathbf{r})$ outside $\mathcal{D}$ becomes arbitrary, leading to inaccurate phases. In such cases, it is more correct to evaluate phases by minimizing the cost function of the type in (2.5), which only takes into account the information contained in the estimate $\rho^{est}(\mathbf{r})$, restricted to domain $\mathcal{D}$.

We now consider the case in which the densities $\rho(\mathbf{r})$ and $\rho^{est}(\mathbf{r})$ are not limited in resolution. In this case, we first need to low-pass filter $\rho(\mathbf{r})$ and $\rho^{est}(\mathbf{r})$ so as to reduce the number of phase variables corresponding to nonzero structure factors to a finite number. This can be done by using the substitution

$$A_k \rightarrow \begin{cases} \exp[-(\pi W/\lambda_k)^2] A_k & \text{if } \lambda_k/2 \ge R \\ 0 & \text{otherwise} \end{cases} \qquad (2.11)$$

and

$$F_k^{est} \rightarrow \begin{cases} \exp[-(\pi W/\lambda_k)^2] F_k^{est} & \text{if } \lambda_k/2 \ge R \\ 0 & \text{otherwise,} \end{cases} \qquad (2.12)$$

where $W$ is a parameter. In the limit $W \gg R$, this substitution corresponds to a convolution in real space using a Gaussian function of the form $\exp[-(r/W)^2]$. In the limit $W \ll R$, this substitution corresponds to a convolution using a long-range function of the type $[2R \sin(r/2R) - r \cos(r/2R)]/r^3$. In the latter limit, the value of the convoluted density $\rho^{est}(\mathbf{r})$ at the postion $\mathbf{r} \in \mathcal{D}$ may contain significant contributions from the unconvoluted density $\rho^{est}(\mathbf{r}')$ at distant positions $\mathbf{r}'$ out of domain $\mathcal{D}$, where $\rho^{est}(\mathbf{r}')$ does not describe accurately the actual electron density. In such circumstances, $\rho^{est}(\mathbf{r})$ may not

describe accurately $\rho(\mathbf{r})$ in domain $\mathcal{D}$ after convolution. It is, therefore, appropriate to use a value of $W$ such that $W \ge R$, which eliminates the long-range tail of the convolution kernel and thus preserves the quality of the estimate $\rho^{est}(\mathbf{r})$ in domain $\mathcal{D}$ at a distance $W$ from its boundary.

Taking into account the crystallographic symmetry and the reality of the electron density leads to linear relations among phases $\varphi_k$. The cost function $E$ can thus be expressed as a function of a restricted set of independent phase variables that we denote $\{\overline{\varphi}_1, ..., \overline{\varphi}_M\}$. Each independent phase variable $\overline{\varphi}_j$ corresponds to a unique reflection. These independent phase variables can be further classified into non-centrosymmetric phases, which can taken continuous values between 0 and 360°, and centrosymmetric phases, whose set of allowed values consists either of the doublet (0,180°) or the doublet (90, 270°).

## 2.2. Simulated annealing

Our implementation of simulated annealing is directly inspired by the work by Kirkpatrick, Gelatt & Vecchi (1983). The cost function $E(\overline{\varphi}_1, ..., \overline{\varphi}_M)$ defines the 'energy' of the system at the position $(\overline{\varphi}_1, ..., \overline{\varphi}_M)$ in the configuration space of the independent phases. The standard Metropolis method is used to generate a stochastic walk in the configuration space in a thermal equilibrium at a 'temperature' $T$. At each step, a phase $\overline{\varphi}_j$ is randomly chosen and given a random change $\Delta\overline{\varphi}_j$. The energy difference $\Delta E$ between the new and the previous configuration is calculated. The new configuration is accepted with probability $P = \min[\exp(-\Delta E/T), 1]$. In case of acceptance, the new configuration is used as the starting point of the next step. Otherwise, the previous configuration is used again as a starting point

This stochastic process enables the sampling of the configuration space according to the Boltzmann probability distribution given by

$$\mathcal{P}(\overline{\varphi}_1, ..., \overline{\varphi}_M) \simeq \exp\{-[E(\overline{\varphi}_1, ..., \overline{\varphi}_M)]/T\}. \qquad (2.13)$$

Both cost function $E$ and 'temperature' $T$ are in units of $\text{Å}^{-6}$ and have the physical dimension of a density squared.

For the annealing schedule, we choose an initial temperature $T_0$ twice as large as the largest value of the energy obtained from a few configurations taken at random. These configurations are obtained by assigning to each phase $\overline{\varphi}_j$ a random value consistent with its type (centrosymmetric or noncentrosymmetric). We then cool exponentially using $T_n = (1 - \varepsilon)^n T_0$, with $n = 1, ..., N_{cycle}$. The initial configuration is taken at random and a number $N_{step}$ of Metropolis iterations is performed at each temperature $T_n$. For relatively small $\varepsilon$ and for values of $N_{step}$ and $N_{cycle}$

large enough, the stochastic process converges to the global minimum of the cost function.

For the sampling of configurations, we use $\Delta\varphi_j = 180°$ in the case of centrosymmetric phases and choose $\Delta\overline{\varphi}_j$ at random in an interval $[-\beta_j, \beta_j]$ in the case of noncentrosymmetric phases. The value $\beta_j$ is initially set equal to 180°. In order to optimize the search in configuration space (Vanderbilt & Louie, 1983), $\beta_j$ is then progressively reduced during the annealing process so that the acceptance rate for Metropolis steps affecting the value of the phase $\overline{\varphi}_j$ is nearly equal to 1/2. In order to do this, we recalculate the acceptance rate for moves of phase $\overline{\varphi}_j$ after every $100M$ Metropolis iterations, where $M$ is the number of independent phase variables. If this acceptance rate is smaller than 0.4, the $\beta_j$ value is reduced by a factor equal to twice the value of the acceptance rate.

### 3. Numerical test

In order to test our optimization procedure, we use the electron density of the model of the protein *Staphylococcus nuclease*, which crystallizes in $P4_1$ symmetry with unit-cell dimensions $a_1 = 48.18$, $a_2 = 48.18$ and $a_3 = 63.51$ Å. The density $\rho^{\mathrm{mod}}(\mathbf{r})$ of this model is convoluted using the substitution

$$F_k^{\mathrm{mod}} \rightarrow \begin{cases} \exp[-(\pi W/\lambda_k)^2]F_k^{\mathrm{mod}} & \text{if } \lambda_k/2 \geq R \\ 0 & \text{otherwise,} \end{cases} \quad (3.1)$$

where $F_k^{\mathrm{mod}}$ denotes the structure factors of $\rho^{\mathrm{mod}}(\mathbf{r})$ and where $R = W = 2$ Å. The amplitudes $A_k$ defined by (2.4) are then taken to be equal to $|F_k^{\mathrm{mod}}|$ and the density estimate $\rho^{\mathrm{est}}(\mathbf{r})$ is taken to be equal to the density $\rho^{\mathrm{mod}}(\mathbf{r})$. The subdomain $\mathscr{D}$ of the unit cell where $\rho^{\mathrm{est}}(\mathbf{r})$ is assumed to accurately describe the electron density is defined on the grid $\mathbf{r}_n$, given by (2.6) as

$$\chi_n = \begin{cases} 1 & \text{if } \rho^{\mathrm{est}}(\mathbf{r}_{n'}) \leq \rho_{\mathrm{threshold}}, \forall n' \text{ such that } |\mathbf{r}_{n'} - \mathbf{r}_n| \leq d \\ 0 & \text{otherwise,} \end{cases}$$
$$(3.2)$$

where $\rho_{\mathrm{threshold}}$ and $d$ are parameters.

The cost function $E$ is then defined using (2.5) with $N_1 = N_2 = N_3 = 32$. Note that this cost function only incorporates the knowledge of the structure-factor amplitudes $A_k$ and the knowledge of the electron density restricted to domain $D$. The minimization of this cost function with respect to phase angles $\varphi_k$ is expected to recover the electron density in the whole unit cell.

Our test is performed for domains $\mathscr{D}$ of various sizes defined by (3.2), with $d = 2$ Å and (a) $\rho_{\mathrm{threshold}} = 0.25$, (b) $\rho_{\mathrm{threshold}} = 0.3$, (c) $\rho_{\mathrm{threshold}} = 0.35$, (d)

Table 1. *Average phase difference (°) between the structure factors $F_k^{\mathrm{mod}}$ of the density $\rho^{\mathrm{mod}}(\mathbf{r})$ and the structure factors of $\rho^{\mathrm{mod}}(\mathbf{r})$ restricted to domain $\mathscr{D}$, calculated in various resolution ranges (Å) for the cases (a), (b), (c), (d) and (e) considered in the text*

The upper and lower numbers for each case correspond to centrosymmetric and noncentrosymmetric phases. The average is calculated using a weight given by the amplitude $A_k$ defined by (2.4).

| Resolution range $\lambda_k/2$ | ∞–8 | 8–4 | 4–2.66 | 2.66–2 |
|---|---|---|---|---|
| Number of phases | 7 | 21 | 35 | 51 |
| | 13 | 120 | 338 | 667 |
| (a) | 138.0 | 124.6 | 99.2 | 92.2 |
| | 148.7 | 133.5 | 108.3 | 93.8 |
| (b) | 138.0 | 134.9 | 96.2 | 89.9 |
| | 132.3 | 129.7 | 109.7 | 91.0 |
| (c) | 23.25 | 163.7 | 119.4 | 96.8 |
| | 107.7 | 125.5 | 111.0 | 97.2 |
| (d) | 23.3 | 147.3 | 120.3 | 121.3 |
| | 52.0 | 112.5 | 117.6 | 96.8 |
| (e) | 23.3 | 105.3 | 111.4 | 110.1 |
| | 50.5 | 102.5 | 117.9 | 98.6 |

$\rho_{\mathrm{threshold}} = 0.38$ and (e) $\rho_{\mathrm{threshold}} = 0.41$ Å$^{-3}$. The domains thus defined respectively occupy a fraction (a) $v = 0.447$, (b) $v = 0.502$, (c) $v = 0.565$, (d) $v = 0.613$ and (e) $v = 0.665$ of the volume of the unit cell, where $v$ is defined by

$$v = (N_1 N_2 N_3)^{-1} \sum_{n_\alpha = 0}^{N_\alpha - 1} \chi_n. \quad (3.3)$$

In all five cases, the domain $\mathscr{D}$ of 'known density' contains mostly interprotein void and the partial knowledge of the density restricted to this domain fails to provide reliable phase estimates by conventional Fourier methods. This can be seen in Table 1, which lists the average phase difference between $F_k^{\mathrm{mod}}$ and the structure factors of the density $\rho^{\mathrm{mod}}(\mathbf{r})$ restricted to domain $\mathscr{D}$, which is obtained from the discrete Fourier transform of $\chi_n \times \rho^{\mathrm{mod}}(\mathbf{r}_n)$. This phase difference is larger than 90°, except for a few phases at low resolution.

In order to evaluate the reconstruction of the electron density in the whole unit cell, we calculate the quantities

$$\sigma_{\mathrm{in}}(T) = \left\{ (vN_1N_2N_3)^{-1} \right.$$

$$\left. \times \sum_{n_\alpha = 0}^{N_\alpha - 1} \chi_n \langle |\rho_{\underline{\varphi}_k}(\mathbf{r}_n) - \rho^{\mathrm{mod}}(\mathbf{r}_n)|^2 \rangle_T \right\}^{1/2} \quad (3.4)$$

and

$$\sigma_{\mathrm{out}}(T) = \left\{ [(1-v)N_1N_2N_3]^{-1} \right.$$

$$\left. \times \sum_{n_\alpha = 0}^{N_\alpha - 1} (1 - \chi_n) \langle |\rho_{\underline{\varphi}_k}\mathbf{r}_n) - \rho^{\mathrm{mod}}(\mathbf{r}_n)|^2 \rangle_T \right\}^{1/2}, \quad (3.5)$$

where $\langle\ \rangle_T$ denotes the average calculated during the Metropolis iteration at a given temperature $T$. The quantities $\sigma_{in}(T)$ and $\sigma_{out}(T)$ represent average discrepancies between the reconstructed density $\rho_{\varphi_k}(\mathbf{r})$ and the correct density $\rho^{mod}(\mathbf{r})$, respectively evaluated inside and outside domain $\mathscr{D}$. The quantity $\sigma_{in}(T)$ corresponds to the square root of the average value of the cost function $E$, calculated during the Metropolis interation and divided by $v$. The value of $\sigma_{out}(T)$ at a given temperature $T$ can be regarded as the average error obtained in the reconstruction of the density outside domain $\mathscr{D}$ when the uncertainty in the density estimate defined in $\mathscr{D}$ is of order $\sigma_{in}(T)$. The simulated annealing algorithm is applied to cases $(a)$–$(e)$ using parameters $T_0 = 3000\nu$, $\varepsilon = 0.01$, $N_{cycle} = 3000$ and $N_{step} = 900$. The quantities $\sigma_{in}(T)$ and $\sigma_{out}(T)$ are calculated every 40 cooling steps. Our results for $\sigma_{in}(T)$ and $\sigma_{out}(T)$ are shown in Fig. 1. Each dot corresponds to a fixed value of temperature.

Fig. 1 can be discussed in terms of the average errors $\Delta_{in}$ and $\Delta_{out}$ in atom positions inside and outside domain $\mathscr{D}$, which correspond to the average errors $\sigma_{in}(T)$ and $\sigma_{out}(T)$ in electron density. For a collection of C atoms with mean electron density $n_e = 0.43\ \text{Å}^{-3}$, the average density error $\sigma$ is given as a function of the average error $\Delta$ in atom positions by

$$\sigma = \Delta(6^{1/2}n_e^{1/2})/(2^{3/4}\pi^{3/4}W^{5/2}) \qquad (3.6)$$

for errors $\Delta$ small compared to the Gaussian width $W$ defined by (3.1). In cases $(c)$, $(d)$ and $(e)$, we see in Fig. 1 that the error $\sigma_{out}(T)$ in the density reconstructed outside domain $\mathscr{D}$ is always smaller than twice the uncertainty $\sigma_{in}(T)$ in the density estimate defined in $\mathscr{D}$, suggesting that the position of atoms outside domain $\mathscr{D}$ could be obtained with an error

Table 2. *Average difference* (°) *between the phases of structure factors* $F_k^{mod}$ *and the annealed phases obtained at temperature* $T$ *such that the quantity* $\sigma_{in}(T)$ *defined by* (3.4) *is equal to* 0.01 Å$^{-3}$

The upper and lower numbers for each case correspond to centro-symmetric and noncentrosymmetric phases. The average is calculated using a weight given by the amplitude $A_k$ defined by (2.4)

| $\lambda_k/2$ | $\infty$–8 | 8–4 | 4–2.66 | 2.66–2 |
|---|---|---|---|---|
| $(a)$ | 0.00 | 0.00 | 15.47 | 51.71 |
|       | 1.37 | 6.51 | 30.38 | 56.89 |
| $(b)$ | 0.00 | 0.00 | 7.41 | 28.25 |
|       | 1.01 | 5.36 | 23.15 | 44.50 |
| $(c)$ | 0.00 | 0.00 | 6.49 | 19.61 |
|       | 0.90 | 2.83 | 13.72 | 30.73 |
| $(d)$ | 0.00 | 0.00 | 1.40 | 18.79 |
|       | 0.82 | 2.64 | 12.20 | 26.74 |
| $(e)$ | 0.00 | 0.00 | 0.45 | 24.08 |
|       | 0.78 | 2.33 | 9.51 | 23.59 |

$\Delta_{out}$ smaller than twice the uncertainty $\Delta_{in}$ in the position of atoms in domain $\mathscr{D}$.

In cases $(a)$ and $(b)$, the annealing process reaches a stationary behavior characterized by a finite value of $\sigma_{in}(T)$ at low temperature, suggesting that the system cannot reach thermal equilibrium due to fast cooling and is trapped into a local minimum.

Let us now consider the phase recovery. Table 2 lists the average differences between true and recovered phases when $\sigma_{in} = 0.01\ \text{Å}^{-3}$, which, following (3.6), corresponds to an average error $\Delta_{in} \simeq 0.15$ Å in the position of atoms in domain $\mathscr{D}$. In cases $(d)$ and $(e)$, the phases are recovered to *ca* 25°. For all cases, we observe that the corresponding errors in structure factors $F_k$ are roughly independent of the wavlength $\lambda_k$ for noncentrosymmetric reflections.

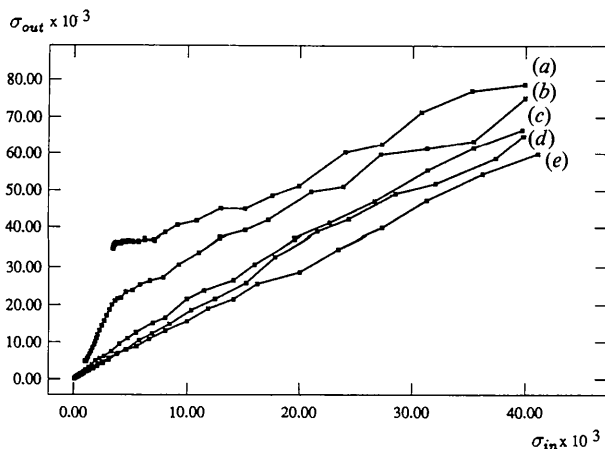We now consider the sensitivity of the phase recovery to noise in the amplitudes $A_k$ of structure



Fig. 1. Quantities $\sigma_{in}(T)$ and $\sigma_{out}(T)$ defined by (3.4) and (3.5) and evaluated for cases $(a)$–$(e)$ during the annealing procedure described in the text. The dots correspond to values of temperature with fixed ratio $T_{n+1}/T_n = 0.669$.
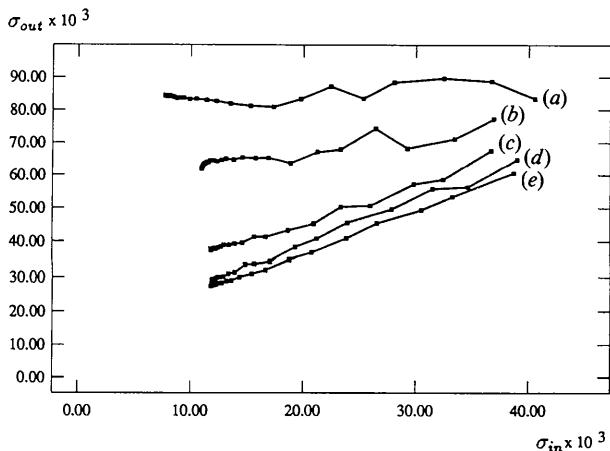


Fig. 2. Quantities $\sigma_{in}(T)$ and $\sigma_{out}(T)$ evaluated for cases $(a)$–$(e)$ during the annealing procedure in the presence of 10% noise in the data for the amplitudes $A_k$ defined by (2.4). The dots correspond to values of temperature with fixed ratio $T_{n+1}/T_n = 0.669$.

factors. Fig. 2 shows results for annealing processes performed for the same test case as before, but using amplitudes given by $A_k = u \times |F_k^{\text{mod}}|$, where $u$ is a random variable of mean 1 and variance 0.1. At low temperature, the annealing process reaches a stationary behavior characterized by a large value of $\sigma_{\text{in}}$. This is because the values assigned to amplitudes $A_k$ are incompatible with the density estimate $\rho^{\text{est}}(\mathbf{r})$. Again, in cases $(c)$, $(d)$ and $(e)$, the error $\sigma_{\text{out}}(T)$ in the density reconstructed outside domain $\mathscr{D}$ is always smaller or equal to twice the uncertainty $\sigma_{\text{in}}(T)$ in the density estimate defined in $\mathscr{D}$. Finally, we note that the effect of noise in amplitude $A_k$ on phase recovery is stronger in cases corresponding to a domain $\mathscr{D}$ of small size, as indicated by the important enhancement of $\sigma_{\text{out}}(T)$ for cases $(a)$ and $(b)$ in the presence of such noise.

We conclude from these simulations that the optimization procedure is practical for the treatment of protein diffraction data at relatively high resolution and can successfully handle cases in which the partial structural information is insufficient for the application of Fourier methods.

## 4. Spin glass analogy

We now address the questions of the nature and difficulty of the optimization problem presented here. Equation (2.8) for the cost function $E$ can be regarded as the interaction energy of a system of rotators (or spins) parametrized by angles $\varphi_k$. The first term describes a coupling between rotators, the second term describes the interaction of rotators with an external field and the third term is a constant. This system bears some similarity to the spin glass problem with long-range interactions, which has been extensively studied (Kirkpatrick & Sherington, 1978).

In the above spin glass system, every spin is coupled to all others. For each couple of spins, this interaction is randomly chosen to be either ferromagnetic (favoring spin alignment) or antiferromagnetic (favoring spin anti-alignment). This problem is characterized by the existence of numerous energy minima which are similar in energy and which are distant from each other and/or separated by high-energy barriers in configuration space (Fu & Anderson, 1986). At low temperature, the relaxation of the system to thermal equilibrium is, therefore, significantly slowed down, i.e. the characteristic time, called relaxation time, needed to visit all energy minima according to Boltzmann's probability distribution is significantly increased (Kirkpatrick & Sherington, 1978). In the context of optimization, this slowing down of relaxation requires the use of a very slow annealing schedule, leading to an important increase of computational work. The relaxation time,

and hence the computational work, grows exponentially with the number of spins in the system. Thus, the spin glass problem is an NP-complete optimization problem (Barahona, 1982).

It is interesting to compare the optimization problem arising in phasing using spatial constraints to a spin glass problem and to ask in particular if the former problem exhibits the slowing down in relaxation characteristic of the spin glass problem. Clearly, the difficulty of the present optimization problem depends on the amount of information available in terms of the size of a region of known density. For example, in the extreme case where the density estimate $\rho^{\text{est}}(\mathbf{r})$ is assumed to be accurate in the whole unit cell, the cost function $E$ only possesses one minimum and the phase configuration corresponding to this minimum is simply given by conventional Fourier methods. We also note that the optimization procedure described in the previous section for test cases $(c)$, $(d)$ and $(e)$ is able to retrieve 1252 independent phases with relatively little computational work, suggesting that the present optimization problem is simpler than a spin glass problem when the region $\mathscr{D}$ of known density occupies a large enough volume in the unit cell. On the other hand, the stationary behavior of $\sigma_{\text{in}}$ observed at low temperature in cases $(a)$ and $(b)$ suggests that the stochastic process could not reach thermal equilibrium due to an increase in relaxation time for cases corresponding to a domain $\mathscr{D}$ of small size.

In order to investigate the question of slowing down in relaxation, we study the relaxation time and the sensitivity of phase recovery to noise in the density estimate in the case of a small test problem for regions $\mathscr{D}$ of various sizes. In this problem, we use the density $\rho^{\text{mod}}(\mathbf{r})$ of the model of *Staphylococcus nuclease* convoluted using the substitution

$$F_k^{\text{mod}} \rightarrow \begin{cases} F_k^{\text{mod}} & \text{if } |\mathbf{k}_\alpha| < 4, \ \alpha = 1, 2, 3 \\ 0 & \text{otherwise.} \end{cases} \quad (4.1)$$

Again, the amplitudes $A_k$ are taken to be equal to $|F_k^{\text{mod}}|$ and the density estimate $\rho^{\text{est}}(\mathbf{r})$ is taken to be equal to the density $\rho^{\text{mod}}(\mathbf{r})$. The grid $\mathbf{r}_n$ is given by (2.6) with $N_1 = N_2 = N_3 = 8$. The subdomain $\mathscr{D}$ of the unit cell is defined by (3.2) with $d = 3$ Å and $(a')$ $\rho_{\text{threshold}} = 0.08$, $(b')$ $\rho_{\text{threshold}} = 0.085$, $(c')$ $\rho_{\text{threshold}} = 0.09$, $(d')$ $\rho_{\text{threshold}} = 0.14$, $(e')$ $\rho_{\text{threshold}} = 0.17$, $(f')$ $\rho_{\text{threshold}} = 0.2$ and $(g')$ $\rho_{\text{threshold}} = 0.22$ Å$^{-3}$. The domains thus defined respectively occupy a fraction $(a')$ $v = 0.32$, $(b')$ $v = 0.336$, $(c')$ $v = 0.367$, $(d')$ $v = 0.406$, $(e')$ $v = 0.445$, $(f')$ $v = 0.516$ and $(g')$ $v = 0.578$ of the volume of the unit cell. The cost function $E(\varphi_1, \ldots, \varphi_M)$ is then defined by (2.5). This cost function depends on $M = 48$ independent phases. In the spin glass case, the system with 48 spins is large enough to lead to generic spin glass behavior (Mackenzie & Young, 1982).

The sensitivity $s$ is defined as

$$s = ([v/(1 - v)]\{[\mathrm{Tr}(B^{-1}C)/M'] - 1\})^{1/2}, \quad (4.2)$$

where $M'$ is the number of noncentrosymmetric phases in the problem, Tr denotes the trace operator and where $B$ in an $M' \times M'$ matrix is given by

$$B_{ij} = [\partial^2 E/(\partial_{\bar{\varphi}_i}\partial_{\bar{\varphi}_j})] \, (\bar{\varphi}_1^{\mathrm{mod}}, \, ..., \, \bar{\varphi}_M^{\mathrm{mod}}), \quad (4.3)$$

where the partial derivatives are taken with respect to noncentrosymmetric phases $\bar{\varphi}_i$ and $\bar{\varphi}_j$ only and where $\bar{\varphi}_1^{\mathrm{mod}}, ..., \bar{\varphi}_M^{\mathrm{mod}}$ denote the phases of $F_k^{\mathrm{mod}}$. The matrix $C$ is defined by (4.3) by substituting for $E$ the cost function defined by taking domain $\mathscr{D}$ to occupy the whole unit cell.

The sensitivity $s$ corresponds to the ratio $\sigma_{\mathrm{out}}(T)/\sigma_{\mathrm{in}}(T)$ of quantities $\sigma_{\mathrm{in}}(T)$ and $\sigma_{\mathrm{out}}(T)$ defined by (3.4) and (3.5) evaluated in the limit $T \to 0$ at thermal equilibrium, i.e. for a large number of Metropolis iterations. This can be verified by expressing this ratio in terms of the Boltzmann probability distribution given by (2.13) and by expanding the cost function $E$ up to second order in the phase variables $\bar{\varphi}_i$ around its minimum. A large value of sensitivity $s$ indicates that density $\rho(\mathbf{r})$ can still fluctuate for $\mathbf{r}$ outside region $\mathscr{D}$, even when $\rho(\mathbf{r})$ is simultaneously constrained to match $\rho^{\mathrm{est}}(\mathbf{r})$ within region $\mathscr{D}$ and to be consistent with structure-factor amplitudes $A_k$, i.e. the electron density $\rho(\mathbf{r})$ is undefined due to lack of information.

Our results for $s$ are shown in Fig. 3 as a function of size $v$ of the region $\mathscr{D}$. The divergence of $s$ for small values of $v$ reflects the failure of phase recovery due to the lack of information resulting from the small size of region $\mathscr{D}$ of known density.

In order to estimate the relaxation time, we evaluate the average distance $K(t)$ between configurations obtained in the stochastic process at 'times' $t_0$ and

$t_0 + t$, where the 'time' is defined as the number of Metropolis iterations divided by the number $M$ of phase variables. $K(t)$ is evaluated at thermal equilibrium for a given temperature. The distance between two configurations of the system of phases is defined as the average squared difference between the electron densities corresponding to each configuration. For a large 'time' interval $t$, the average distance $K(t)$ tends to an asymptotic value denoted by $K_0$. The relaxation time $\tau$ is estimated by fitting this average distance as

$$K(t) \simeq K_0[1 - \exp(-t/\tau)]. \quad (4.4)$$

In order to detect a slowing down in relaxation characteristic of spin glass systems, we evaluate the relaxation time $\tau$ at low temperature $T$ corresponding to density fluctuations of size $\sigma_{\mathrm{in}} = 0.01 \, \text{Å}^{-3}$ in domain $\mathscr{D}$. A large number of Metropolis iterations are performed at temperature $T$ prior to the evaluation of the average distance $K(t)$, in order to allow the energy to relax first.

Our results for relaxation time $\tau$ are shown in Fig. 4 as a function of the size $v$ of the region $\mathscr{D}$. The divergence of $\tau$ for small values of $v$ reflects the fact that the minimization of cost function $E$ becomes increasingly difficult as the domain $\mathscr{D}$ of known density is reduced. This can be understood as follows. In the case of a large domain $\mathscr{D}$, the cost function $E$ defined by (2.5) contains a large amount of information which strongly restricts the number of phase configurations with a low value of $E$. Thus, no spurious minimum of $E$ can compete with the minimum corresponding to the correct phases $\bar{\varphi}_1^{\mathrm{mod}}, ..., \bar{\varphi}_M^{\mathrm{mod}}$. As the temperature is reduced, the fluctuations of phase variables decrease progressively. The corresponding optimization problem is easy because it can be separated into two independent problems, one
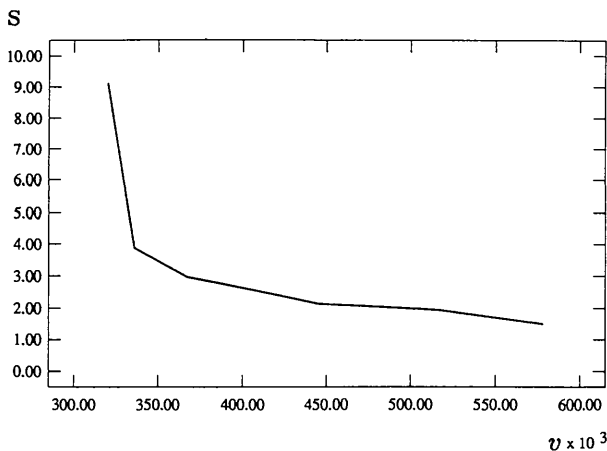


Fig. 3. Sensitivity $s$ defined by (4.2) for various sizes $v$ of domain $\mathscr{D}$ corresponding to cases $(a')–(f')$ described in the text.
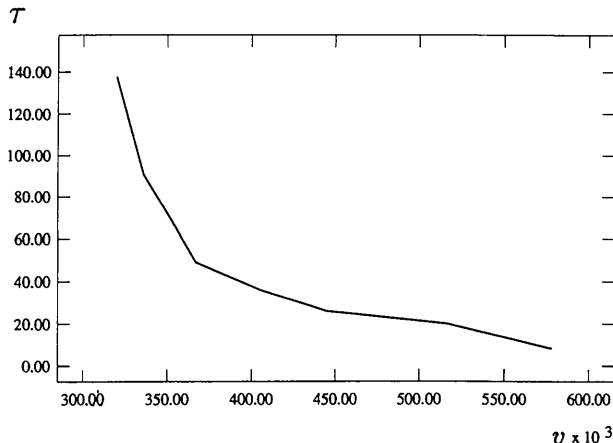


Fig. 4. Relaxation time $\tau$ defined by (4.4) for various sizes $v$ of domain $\mathscr{D}$ corresponding to cases $(a')–(f')$ described in the text. The uncertainty on the estimated values of $\tau$ is smaller than 3%.

consisting of the evaluation of the gross features of the electron density, which is solved at high temperature, and the other consisting of the evaluation of the fine details of the electron density, which is solved at low temperature (Kirkpatrick, Gelatt & Vecchi, 1983). In the case of a small domain $\mathscr{D}$, spurious minima of $E$ develop and compete with the minimum corresponding to the correct phases. As the temperature is reduced, the stochastic walk is confined to a particular minimum. This corresponds to the simultaneous freezing of a number of phase variables. In the vicinity of this transition, very slow cooling is required in order to allow the system to freeze in the configuration with lowest 'energy' $E$.

Note the similarity of the dependence of quantities $\tau$ and $s$ on size $v$ of domain $\mathscr{D}$. This suggests that the slowing down in relaxation characteristic of spin glass systems only occurs when phase recovery fails owing to the lack of information resulting from the small size of region $\mathscr{D}$ of known density. This also suggests that the optimization problem arising in phasing using spatial constraints is simpler than a spin glass problem when enough information is available for the unambiguous determination of the electron density.

### References

Barahona, F. (1982). *J. Phys. A*, **15**, 3241–3253.
Bricogne, G. (1976). *Acta Cryst.* **A32**, 832–847.
Cochran, W. (1951). *Acta Cryst.* **4**, 408–411.
Cura, V., Krishnaswamy, S. & Podjarny, A. D. (1992). *Acta Cryst.* **A48**, 756–764.
Fitzgerald, P. M. D. (1991). In *Crystallographic Computing 5*, edited by D. Moras, A. D. Podjarny & J. C. Thierry. IUCr/Oxford Univ. Press.
Fu, Y. & Anderson, P. W. (1986). *J. Phys. A: Math. Nucl. Gen.* **19**, 1605–1620.
Kirkpatrick, S., Gelatt, C. D. & Vecchi, M. P. (1983). *Science*, **220**, 671–680.
Kirkpatrick, S. & Sherington, D. (1978). *Phys. Rev. B*, **17**, 4384–4403.
Mackenzie, N. D. & Young, A. P. (1982). *Phys. Rev. Lett.* **49**, 301–304.
Read, R. J. (1986). *Acta Cryst.* **A42**, 140–149.
Rossmann, M. G. & Blow, D. M. (1962). *Acta Cryst.* **15**, 24–31.
Sim, G. A. (1959). *Acta Cryst.* **12**, 813–815.
Srinavasan, R. (1966). *Acta Cryst.* **20**, 143–145.
Szöke, A. (1993). *Acta Cryst.* **A49**, 853–866.
Vanderbilt, D. & Louie, S. G. (1984). *J. Comput. Phys.* **56**, 259–271.
Wang, B. (1985). *Methods Enzymol.* **115**, 90–112.

# On the Calculation of the Lattice Energy of Ionic Crystals using the Detailed Electron-Density Distribution. I. Treatment of Spherical Atomic Distributions and Application to NaF

By Zhengwei Su and Philip Coppens

*Department of Chemistry, State University of New York at Buffalo, Buffalo, New York 14214-3094, USA*

*Dedicated to Professor E. F. Bertaut on his 80th birthday*

## Abstract

Ewald's method of accelerated convergence [Ewald (1921). *Ann. Phys.* (*Leipzig*), **64**, 253–287] is generalized to calculate the electrostatic potential of a crystal in which the atoms have overlapping spherical densities. The algorithm is applied to the cubic NaF crystal. The potentials at the Na and F nuclei are calculated for the free-ion model and for the results from a $\kappa$ refinement of the experimental data of Howard & Jones [*Acta Cryst.* (1977), **A33**, 776–783]. The $\kappa$ refinement indicates an incomplete charge transfer but gives an electrostatic energy close to that of the point-charge model with full charge transfer and a lattice energy that is in good agreement with the experimental value.

## Introduction

Although calculations of the lattice energy of ionic crystals often produce good agreement with experimental values as determined in a Born–Haber cycle, they are generally based on a point-charge model, which does not properly describe the charge distribution in the crystals. In this article, we describe methods to evaluate the energy for a crystal consisting of spherical ions and apply the results to NaF,