graphite monochromator and $\omega/2\theta$ scan mode in the range $1 \le \theta \le 25°$. Table 2 gives the figures of merit of the best sets for all three methods; that of MAGEX89 gave an E map showing 22 atoms and is reproduced in Fig. 1. Fourier calculations revealed the remaining atoms and least-squares refinement gave a final residual of 0·055 for the observed reflexions.

## Discussion

It will be seen from Table 1 that for the trial structures the MAGEX89 method performed somewhat better than the other two and the computer resources used by MAGEX89 were 9% less than those of RANTAN. Not too much should be made of that since by modifying parameters all three methods are probably capable of solving all the structures. What we do say is that it is worthwhile having MAGEX89 available. While any individual method may not succeed for a particular structure, the probability of failure is far lower with many methods available.

The version of MAGEX89 we have used has been programmed for a PDP11/44 but should be able to run on most standard personal computers. It can handle all 230 space groups in the standard orientations, including alternative settings, as given in International Tables for Crystallography (1987).

### References

DEBAERDEMAEKER, T., TATE, C. & WOOLFSON, M. M. (1985). Acta Cryst. A41, 286-290.
DEBAERDEMAEKER, T., TATE, C. & WOOLFSON, M. M. (1988). Acta Cryst. A44, 353-357.
HULL, S. E. & IRWIN, M. J. (1978). Acta Cryst. A34, 863-870.
HULL, S. E., VITERBO, D., WOOLFSON, M. M. & ZHANG, S. H. (1981). Acta Cryst. A37, 566-572.
International Tables for Crystallography (1987). Vol. A. Dordrecht: Kluwer.
WHITE, P. S. & WOOLFSON, M. M. (1975). Acta Cryst. A31, 53-56.
YAO, J. X. (1981). Acta Cryst. A37, 642-644.
ZHANG, S. H., LUO, B. S., CHEN, S. K. & YAO, J. W. (1989). Acta Phys. Chim. Sin. 5, 536-540.
ZHANG, S. H. & WOOLFSON, M. M. (1982). Acta Cryst. A38, 683-685.

# A Method for Multiple Superposition of Structures

By A. SHAPIRO AND J. D. BOTHA

Department of Mathematics, Applied Mathematics and Astronomy, University of South Africa, PO Box 392, Pretoria 0001, South Africa

A. PASTORE

European Molecular Biology Laboratory, Meyerhofstrasse 1, Postfach 1022.09, 6900 Heidelberg, Germany

AND A. M. LESK

Department of Haematology, University of Cambridge Clinical School, MRC Centre, Hills Road, Cambridge CB2 2QH, England

## Abstract

The study of families of protein structures is important in analysing the results of NMR structure determinations and in investigating mechanisms of molecular evolution at the level of conformation. A method is discussed for finding the transformations that mutually superpose an arbitrary number of structures in the least-squares sense given specified atom-to-atom correspondence.

## 1. Introduction

Superposition has become an important tool for comparing protein structures and for deriving an 'average' structure from a family of conformations of a protein. The problem arises regularly in the determination of the three-dimensional structure of a protein in solution using nuclear magnetic resonance, which typically produces an ensemble of conformations (Wuethrich, 1986). Sets of interproton distances

obtained from measuring cross relaxation between nuclei less than 5 Å apart can be converted into a 3D structure by several methods, but all suffer from underdetermination as the number of observable distances is far smaller than the number of degrees of freedom. As a result, a family of structures, all consistent with experimental data, is produced (Kaptein, Zuiderweg, Scheek, Boelens & van Gunsteren, 1985; Clore, Gronenborn, Bruenger & Karplus, 1985; Havel, Kuntz & Crippen, 1983; Braun & Gō, 1985; Bassolino et al., 1988; Nilges, Clore & Gronenborn, 1988).

Given a family of conformations, a method to measure their similarity and dispersion is necessary. Techniques in common use rely on calculations of minimal root-mean-square deviations of atomic positions to achieve optimal superposition. For a pair of structures, techniques for finding the global minimum root-mean-square deviation in atomic position, with respect to rigid-body relative motions, are known (von Neumann, 1937; Kabsch, 1976, 1978; Golub & Van Loan, 1983). However, for more than two structures, most methods so far described involve combinations of pairwise superpositions. In practice, NMR spectroscopists usually either (a) calculate an average structure and then superpose all the structures on it or (b) superpose all the structures on one member of the family chosen as a reference.

Three previous reports address directly the problem of multiple structural superposition (Gerber & Müller, 1987; Sutcliffe, Haneef, Carney & Blundell, 1987; Kearsley, 1990). These methods provide solutions of multiple superposition problems in many practical cases, especially if the structures are quite similar. However, it is of interest to note that another general solution, free of simplifying assumptions and heuristic devices, exists. This is discussed in § 2.

## 2. Problem formulation and iterative procedures

In this paper we discuss an optimization problem of simultaneous rotation of $M$ rigid molecules to maximal similarity. That is, let $(x_{ik})$ be a set of $3 \times 1$ vectors such that for every $k = 1, \ldots, M$ the set $(x_{ik})$, $i = 1, \ldots, N$, represents atomic coordinates of a given molecule. Following Gerber & Müller (1987) we take the following criterion for a match (similarity) between rotated molecules:

$$E = \sum_{k<l}^{M} v_{kl} \sum_{i=1}^{N} w_i \|T_k x_{ik} - T_l x_{il}\|^2, \qquad (1)$$

where $T_1, \ldots, T_M$ are rotation (orthogonal) matrices and $v_{kl}$ and $w_i$ are given positive weights. (Here $\|x\|^2 = x'x$ denotes the squared length of a vector $x$ and $x'$ stands for the transpose of $x$.) An optimal similarity is achieved by minimization of the criterion $E$ as function of the orthogonal matrices $T_1, \ldots, T_M$.

It will be convenient to formulate the obtained optimization problem in a matrix form. Consider $3 \times N$ data matrices $X_k = (x_{1k}, \ldots, x_{Nk})$, $k = 1, \ldots, M$, the $N \times N$ diagonal matrix $W = \mathrm{diag}\,(w_i)$ and the $3 \times 3$ matrices $S_{kl} = v_{kl} X_k W X_l'$, $k, l = 1, \ldots, M$. Then

$$E = \sum_{k<l}^{M} v_{kl}\, \mathrm{tr}\,(T_k X_k - T_l X_l) W (T_k X_k - T_l X_l)'$$

$$= -2 \sum_{k<l}^{M} \mathrm{tr}\, T_k\, S_{kl} T_l'$$

$$+ \text{a sum independent of } T_j.$$

Consequently, the problem of minimization of $E$ is equivalent to maximization of

$$G = \sum_{k<l}^{M} \mathrm{tr}\, T_k S_{kl} T_l'. \qquad (2)$$

For $M = 2$ this optimization problem has a closed-form solution. Namely, one has to choose $T_1$ and $T_2$ in such a way that the matrix $T_1 S_{12} T_2'$ is diagonal and positive semidefinite. That is, $T_1$ and $T_2$ are formed by orthonormal eigenvectors of the matrices $S_{12} S_{12}'$ and $S_{12}' S_{12}$, respectively. This result has a long history. It was derived by von Neumann (1937) and has been rediscovered many times since. Notice that the optimal solution is not unique. We can always premultiply $T_1$ and $T_2$ by an orthogonal matrix without changing the trace of $T_1 S_{12} T_2'$. In particular, we can replace $T_1 S_{12} T_2'$ by $T_2' T_1 S_{12}$. The orthogonal matrix $T = T_2' T_1$ is then optimal if and only if $TS_{12}$ is symmetric and positive semidefinite (cf. Ten Berge, 1977).

For $M \geq 3$ a closed-form solution is not available and an iterative procedure is required. A simple and surprisingly efficient algorithm was proposed by Ten Berge (1977). The idea is to maximize $G$ with respect to one orthogonal matrix at a time keeping other orthogonal matrices fixed. That is, let matrices $T_2, \ldots, T_M$ be fixed and consider $G$ as a function of $T_1$ alone,

$$G = \mathrm{tr}\, T_1 S_{1\cdot} + \text{a sum independent of } T_1,$$

where

$$S_{k\cdot} = \sum_{l \neq k} S_{kl} T_l', \qquad k = 1, \ldots, M.$$

We can now maximize $G$ with respect to $T_1$ by the diagonalization procedure and similarly for $T_2$ etc.

Ten Berge's algorithm can be described as follows. Choose initial values of the matrices $T_k$, $k = 1, \ldots, M$. Usually these initial values are taken to be the identity matrices. Maximize $G$ with respect to $T_1$ and substitute the obtained optimal value. Do the same, in turn, for $T_2, \ldots, T_M$. Repeat the procedure until the increment in $G$ becomes less than a prescribed precision given by a positive number $\varepsilon$, say $\varepsilon = 10^{-4}$.

Ten Berge's algorithm converges to a stationary point where all matrices $S_k$, $k = 1, \ldots, M$, are symmetric and positive semidefinite. This is a necessary condition for optimality. Unfortunately, this necessary condition may be not sufficient (see Ten Berge, 1977, p. 270). Therefore a dual procedure was proposed by Shapiro & Botha (1988) in order to verify optimality of an obtained stationary point. Let

$$S(Z) = \begin{bmatrix} Z_1 & S_{12} & \ldots & S_{1M} \\ S_{21} & Z_2 & \ldots & S_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ S_{M1} & S_{M2} & \ldots & Z_M \end{bmatrix}$$

be a $3M \times 3M$ symmetric matrix considered as a function of the symmetric block-diagonal matrix $Z = \text{diag}(Z_1, \ldots, Z_M)$. Denote by $\lambda_1(Z) \geq \ldots \geq \lambda_{3M}(Z)$ the eigenvalues of $S(Z)$. Then for every $Z$ the number

$$f(Z) = \tfrac{1}{2}\left[ M \sum_{i=1}^{3} \lambda_i(Z) - \text{tr}\, S(Z) \right] \qquad (3)$$

gives an upper bound for the maximum of $G$ (Shapiro & Botha, 1988, theorem 1). By minimizing $f(Z)$ over all block-diagonal matrices, one obtains the corresponding least upper bound. It can be shown that this least upper bound is equal to the maximum of $G$ if $\lambda_3(Z_0) \neq \lambda_4(Z_0)$, where $Z_0$ is the minimizer of $f(Z)$.

It was found by Shapiro & Botha (1988) that a very good starting value for minimization of $f(Z)$ is constructed as follows. Let $P_{kl}D_{kl}Q'_{kl}$ be singular-value decompositions of the matrices $S_{kl}$, i.e. $D_{kl}$ are diagonal and positive semidefinite and $P_{kl}$ and $Q_{kl}$ are orthogonal. Then take

$$Z_k^* = -\sum_{\substack{l=1 \\ l \neq k}}^{M} P_{kl}D_{kl}P'_{kl}, \qquad (4)$$

$k = 1, \ldots, M$ and $Z^* = \text{diag}(Z_1^*, \ldots, Z_M^*)$.

It was found that for this value of $Z$ the corresponding upper bound $f(Z^*)$ was usually very close to the maximum value of $G$ produced by Ten Berge's algorithm. Therefore, for practical purposes, it was usually sufficient to evaluate the upper bound $f(Z^*)$ in order to confirm optimality of the calculated stationary point.

It may happen that some of the optimal orthogonal matrices calculated by Ten Berge's algorithm have negative determinant (equal to $-1$). This will be the case if the determinant of the corresponding matrices $S_{k.} = \sum_{l \neq k} S_{kl}$ is negative. Then the obtained optimal solution suggests reflections as well as rotations of the considered structures. Although such a situation is unlikely to happen in practice, it is possible from the theoretical point of view. In the remainder of this section we briefly discuss how to deal with this case.

Let us study first the case of $M = 2$. Again, in this case the problem has a closed-form solution. That is, consider the matrix $S_{12}$ and suppose that $\det S_{12} < 0$. Let $\bar{T}_1$ and $\bar{T}_2$ be orthogonal matrices such that the matrix $S_{12}^* = \bar{T}_1 S_{12} \bar{T}_2'$ is diagonal, $S_{12}^* = \text{diag}(d_1, d_2, d_3)$, and positive semidefinite. Since $\det S_{12} < 0$, we have here that the determinants of $\bar{T}_1$ and $\bar{T}_2$ have different signs, say $\det \bar{T}_1 = 1$ and $\det \bar{T}_2 = -1$. We can choose matrices $\bar{T}_1$ and $\bar{T}_2$ in such a way that the diagonal elements of $S_{12}^*$ are arranged in decreasing order, i.e. $d_1 \geq d_2 \geq d_3$.

We have to find rotation matrices $T_1$ and $T_2$ such that the trace of $T_1 S_{12} T_2'$ is maximized. It is claimed then that the optimal rotation matrices are given by the matrices $\bar{T}_1$ and $D\bar{T}_2$, where $D$ is the diagonal matrix with diagonal elements 1, 1 and $-1$, i.e. $D = \text{diag}(1, 1, -1)$. That is, the maximum of the trace of $T_1 S_{12} T_2'$, subject to $T_1$ and $T_2$ being orthogonal and $\det T_1 = \det T_2 = 1$, is given by $d_1 + d_2 - d_3$.

Indeed, consider a $3 \times 3$ orthogonal matrix $T$ with $\det T = -1$. Such a matrix can be represented in the form $T = QAQ'$, where $Q$ is an orthogonal matrix and $A$ is a block-diagonal matrix with the first block given by a $2 \times 2$ rotation matrix $B$ and the second block consisting of the element $-1$ (e.g. Curtis, 1979). We have then

$$\text{tr}\, S_{12}^* T = \text{tr}\, Q'S_{12}^* QA = \text{tr}\, Q_1'S_{12}^* Q_1 B - q_3'S_{12}^* q_3,$$

where $q_1$, $q_2$ and $q_3$ are column vectors of the matrix $Q$ and $Q_1$ is the $3 \times 2$ matrix formed by the first two columns of $Q$, i.e. $Q_1 = [q_1, q_2]$. Since $Q_1'S_{12}^* Q_1$ is symmetric and positive semidefinite, we know that the maximum of $\text{tr}\, Q_1'S_{12}^* Q_1 B$ is attained when $B$ is the identity matrix. Also by Ky Fan's inequality (Ky Fan, 1949) we have $\text{tr}\, Q_1'S_{12}^* Q_1$ is less than or equal to the sum of the two largest eigenvalues of $S_{12}^*$, that is $d_1 + d_2$. Similarly, $q_3'S_{12}^* q_3 \geq d_3$. All this implies

$$\text{tr}\, T_1 S_{12} T_2' = \text{tr}\, S_{12}^* T \leq d_1 + d_2 - d_3$$

where

$$T = \bar{T}_2 T_2' T_1 \bar{T}_1' \quad \text{and} \quad \det T = -1.$$

For $M \geq 3$ an obvious analogue of Ten Berge's iterative algorithm can then be derived. Unfortunately, it appears that in the case of negative determinants the Shapiro-Botha upper-bound procedure cannot be applied for verification of optimality of the calculated solution.

## 3. Numerical experimentation

Ten Berge's algorithm together with the Shapiro-Botha dual procedure was applied to five representative sets of coordinates of the protein molecule acyl phosphatase in an attempt to superpose them optimally. For each representative set there are 392 triples of coordinates that represent the positions of the 392 atoms that make up the whole structure. Only

the unweighted optimization problem was considered; that is, minimization of $E$ in (1) with respect to $\mathbf{T}_1, \ldots, \mathbf{T}_M$, with $v_{kl} = 1 (1 \le k < l \le M)$ and $w_i = 1 (1 \le i \le N)$.

First the representative sets were translated so that the centres of all lie at the origin. Hence for each set $\mathbf{X}_k = (\mathbf{x}_{1k}, \ldots, \mathbf{x}_{Nk})$, $1 \le k \le M$, $\sum_{i=1}^{N} \mathbf{x}_{ik} = \mathbf{0}$. Then the maximal value (denoted by $G_0$) of the function $G$ in (2) was obtained by means of the Ten Berge algorithm. Recall that the maximum of $G$ corresponds to the minimum of the criterion $E$ given in (1).

Note that this algorithm differs from algorithm 1 proposed by Sutcliffe *et al.* (1987, p. 378), which incidentally is identical to an algorithm proposed earlier by Kristof & Wingersky (1971). Both the Ten Berge and Kristof–Wingersky algorithms yield upon convergence stationary points which satisfy some necessary conditions for the optimality of $G$. However, as was indicated by Ten Berge (1977, pp. 270–272), the necessary condition satisfied by a stationary point obtained from the Ten Berge algorithm is stronger than that of Kristof & Wingersky.

The value of $G_0$, which was obtained in two iterations, is 605 606·0. This turned out to be very close to the optimal value of $G$. Evaluation of the upper bound $f(\mathbf{Z}^*)$ of $G$, with $f$ and $\mathbf{Z}^*$ defined as in (3) and (4) respectively, yields the value 605 607·0.

The program used for the above calculation also provides for an attempted minimization of the upper bound $f(\mathbf{Z})$ of $G$ by means of a Newton-like differentiation method. It uses $\mathbf{Z}^*$ as defined in (4) as starting value. However, this is not always possible, since the function $f$ is convex, but not everywhere differentiable. In fact, it is differentiable at a point $\mathbf{Z} = \mathrm{diag}(\mathbf{Z}_1, \ldots, \mathbf{Z}_M)$ if and only if $\lambda_3(\mathbf{Z}) \ne \lambda_4(\mathbf{Z})$ [see Shapiro & Botha (1988) for a more detailed discussion of the properties of $f$ and the relevant minimization technique].

In this particular case $f(\mathbf{Z})$ was minimized successfully. The difference $\lambda_3(\mathbf{Z}) - \lambda_4(\mathbf{Z})$ was sufficiently large throughout the minimization process with $|\lambda_4(\mathbf{Z})/\lambda_3(\mathbf{Z})| \ge 10^2$ (the eigenvalues were all negative). This means that $f(\mathbf{Z})$ was differentiable at each step and that the stationary point is a global minimizer, since $f(\mathbf{Z})$ is convex. Furthermore, the minimal value of $f$ turned out to be equal to $G_0$, which implies that the Ten Berge algorithm converged to the optimal solution. Hence, in this case the orthogonal matrices obtained from the Ten Berge algorithm optimally superpose the original five representative sets of coordinates. Also, it was found that these orthogonal matrices are all rotation matrices.

Finally, note that in this case the upper bound $f(\mathbf{Z}^*)$ is very close to the optimal value of both $f$ and $G$. This was also borne out by other experiments. This may therefore be a good way to check the performance of the Ten Berge algorithm, especially if one does not want to go through the process of minimizing $f$.

### References

BASSOLINO, D. A., HIRATA, F., KITCHEN, D. B., KOMINOS, D., PARDI, A. & LEVY, R. M. (1988). *Int. J. Supercomput. Appl.* **2**, 41–61.

BRAUN, W. & GŌ, N. (1985). *J. Mol. Biol.* **186**, 611–626.

CLORE, G. M., GRONENBORN, A. M., BRUENGER, A. T. & KARPLUS, M. (1985). *J. Mol. Biol.* **186**, 435–455.

CURTIS, C. W. (1979). *Linear Algebra*, 3rd ed., p. 270. Boston: Allyn and Bacon.

GERBER, P. R. & MÜLLER, K. (1987). *Acta Cryst.* A**43**, 426–428.

GOLUB, G. H. & VAN LOAN, C. F. (1983). *Matrix Computations*, pp. 425–426. Baltimore, MD: The Johns Hopkins Univ. Press.

HAVEL, T., KUNTZ, I. D. & CRIPPEN, G. M. (1983). *Bull. Math. Biol.* **45**, 655–720.

KABSCH, W. (1976). *Acta Cryst.* A**32**, 922–923.

KABSCH, W. (1978). *Acta Cryst.* A**34**, 827–828.

KAPTEIN, R., ZUIDERWEG, E. R. P., SCHEEK, R. M., BOELENS, R. & VAN GUNSTEREN, W. F. (1985). *J. Mol. Biol.* **182**, 179–182.

KEARSLEY, S. K. (1990). *J. Comput. Chem.* **11**, 1187–1192.

KRISTOF, W. & WINGERSKY, B. (1971). Proc. 79th Annual Convention of the American Psychological Association, pp. 89–90.

KY FAN (1949). *Proc. Natl Acad. Sci. USA*, **35**, 652–655.

NEUMANN, J. VON (1937). *Tomsk Univ. Rev.* **1**, 286–300.

NILGES, M., CLORE, G. M. & GRONENBORN, A. M. (1988). *FEBS Lett.* **229**, 317–324.

SHAPIRO, A. & BOTHA, J. D. (1988). *SIAM (Soc. Ind. Appl. Math.) J. Matrix Anal. Appl.* **9**, 378–383.

SUTCLIFFE, M. J., HANEEF, I., CARNEY, D. & BLUNDELL, T. L. (1987). *Protein Eng.* **1**, 377–384.

TEN BERGE, J. M. F. (1977). *Psychometrika*, **42**, 267–276.

WUETHRICH, K. (1986). *NMR of Proteins and Nucleic Acids.* New York: Wiley.