

(MeCO)<sub>2</sub><sup>2-</sup>, shows no fold (Adamson, Daly & Forster, 1974).

### Problems with small crystals

The difficulties encountered in recording X-ray diffraction data for structure determination from very small crystals arise from the poor crystal quality as well as from the low intensity of the diffraction pattern. The poor crystal quality is indicated by the large mosaic spread, evaluated as *ca* 3° by *MADNES*. The diffraction spots become very extended in reciprocal space, and the derivation of good integrated intensities, especially for the weaker spots is difficult, as already described by Andrews *et al.* (1988). Improvements in the experimental strategy and processing software may allow somewhat better intensity measurements and structure determination in the future, but it is clear that even with the present procedures sufficient data can be recorded from such a small poor crystal to establish the chemical structure and stereochemistry.

We are grateful to SERC for financial support including a studentship (SJM), to SERC Daresbury Laboratory for experimental and computational facilities, and to Drs M. Z. Papiz, J. R. Helliwell, J. W. Campbell, I. J. Clifton, S. M. Clark, I. D. Glover and T. J. Greenough for advice or assistance in using these.

### References

- ADAMSON, G. W., DALY, J. J. & FORSTER, D. (1974). *Organomet. Chem.* **71**, C17–C19.  
 ANDREWS, S. J., PAPIZ, M. Z., MCMEEKING, R., BLAKE, A. J., LOWE, B. M., FRANKLIN, K. R., HELLIWELL, J. R. & HARDING, M. M. (1988). *Acta Cryst.* **B44**, 73–77.  
 CROMER, D. T. & LIBERMAN, D. (1970). *J. Chem. Phys.* **53**, 1891–1898.  
 JACOB, C. & HEATON, B. T. (1988). Personal communication.  
 MESSERSCHMIDT, A. & PFLUGRATH, J. W. (1987) *J. Appl. Cryst.* **20**, 306–315.  
 SHELDRIK, G. M. (1976). *SHELX*. Program for crystal structure determination. Univ. of Cambridge, England.  
 SMITH, J. H. & WONACOTT, A. J. (1979). *CCP4* program suite. SERC Daresbury Laboratory, England.

*Acta Cryst.* (1990). **B46**, 195–208

## Hydration in Protein Crystals. A Neutron Diffraction Analysis of Carbonmonoxymyoglobin\*

BY XIAODONG CHENG

*Department of Physics, State University of New York at Stony Brook, Stony Brook, NY 11794, USA, and Center for Structural Biology, Department of Biology, Brookhaven National Laboratory, Upton, NY 11973, USA*

AND BENNO P. SCHOENBORN

*Center for Structural Biology, Department of Biology, Brookhaven National Laboratory, Upton, NY 11973, USA*

(Received 25 April 1989; accepted 9 November 1989)

### Abstract

In protein crystallography, it has been customary to ignore the contribution of bulk solvent by omitting the low-order diffraction data in refinement procedures. However, these data contain important information on both the structure of the solvent and the gross features of the unit-cell contents. The contribution of the solvent to the low-order structure-factor terms can be evaluated by dividing the solvent volume into shells extending outward from the surface of the protein. Two hydration layers in myo-

globin crystals were characterized, allowing a better evaluation of the surface structure of the protein, improved placement of bound water and ion molecules, and a better overall fit to the observed data. A reciprocal-space least-squares refinement program was modified to include restraints on the configuration of water binding in water-to-protein and water-to-water associations. The inclusion of the solvent contribution allows all structure factors to be used in the refinement procedure. Eighty-seven water and five ion molecules were localized in carbonmonoxymyoglobin. All water molecules that are bound to protein bind to polar or charged groups, and the final *R* factor is 11.5%.

\* Research carried out under contract DE-AC02-76CH00016 with the US Department of Energy.

### 1. Introduction

In this paper we attempt to correlate the general information on hydration of proteins with crystallographic analysis of the solvent. On the basis of electron microscopic studies, Kellenberger (1978) suggested that water-soluble proteins are surrounded by a tightly bound layer of water molecules. Small-angle neutron scattering showed that the observed radius of gyration of a protein is larger than that predicted for the protein alone (Ibel & Stuhrmann, 1975; Schoenborn, 1989), suggesting that water is sufficiently closely associated with the protein to 'tumble' with it in solution. However, it has been difficult to demonstrate the existence of water layers in protein crystals. Crystallographic analyses of various proteins only 'see' water molecules that are well localized in three dimensions, and the agreement between water sites found in different experiments is generally poor (Takano, 1977*a,b*; Phillips, 1980, 1984; Hanson & Schoenborn, 1981; Kossiakoff, 1985). These bound water molecules represent only a small fraction of the total solvent (Savage, 1986); they are often arranged in clusters, and do not seem to form a layer surrounding the protein. The extent to which a presumed water layer is unexchangeable was investigated by Lehmann *et al.* using ethanol-water (Lehmann, Mason & McIntyre, 1985) and dimethyl sulfoxide-water (Lehmann & Stansfield, 1989) solutions; non-polar ethanol and dimethyl sulfoxide bind to hydrophobic sites in preference to water molecules. Numerous non-crystallographic techniques have been used to estimate the hydration of proteins in attempting to elucidate or improve the structure of the solvent. Analyses of water structure are discussed in recent reviews (Griffin, 1986; Kossiakoff, 1985; Edsall & McKenzie, 1983; Finney, 1979; Kuntz & Kauzmann, 1974; Cooke & Kuntz, 1974).

In crystallographic analysis of proteins, bound water molecules are often treated as though they 'belong' to the protein, but the rest of the bulk solvent is ignored. To minimize the effect of this omission, the low-order diffraction terms are not used in refinement procedures. However, an analysis of the low-order diffraction data can provide information on the general distribution of the solvent and protein within the crystal. To use the low-order crystallographic data in structural refinement, an evaluation procedure for the solvent was developed that uses its average scattering density (Schoenborn, 1988). As demonstrated in §2, this procedure can be extended to evaluate the scattering density of the solvent as a function of distance expanding outward from the surface of the protein. Analysis shows that the water layers surrounding the protein have little mobility, as demonstrated (§3) by a low liquidity factor. This procedure further allows the calculation

of phases for the low-order diffraction data that enhances the localization of bound water, so that all data can be included in the refinement calculations in §4 and §5. In §6, a full description of the solvent structure of carbonmonoxymyoglobin (MBCO) is given.

### 2. Solvent shell and solvent structure factor

The measured structure factor ( $F_o$ ) is described by the two crystal components of solvent ( $F_s$ ) and protein ( $F_p$ ) (Schoenborn, 1988):

$$F_o = |F_s + F_p|. \quad (1)$$

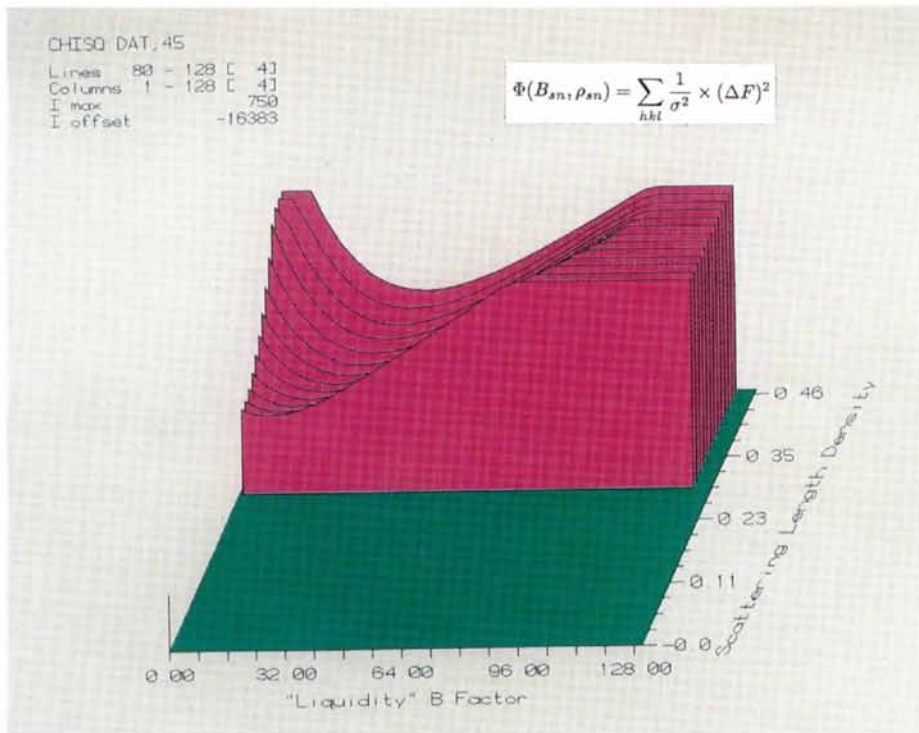
When the initial crystallographic analysis has reached the stage where a good description of the structure of the protein moiety has been obtained, the solvent component of the structure factors is given by the difference between the observed and calculated protein structure factors. Since the overall solvent has a relatively high disorder, here termed the liquidity factor, it contributes only to structure factors with low  $\sin\theta$  values and only low-order diffraction data need to be considered. In this analysis, the solvent component ( $F_s$ ) is expanded into  $n$  structure-factor terms ( $F_{sn}$ ), corresponding to  $n$  shells extending outward from the protein's surface:

$$F_s = \sum_n F_{sn}. \quad (2)$$

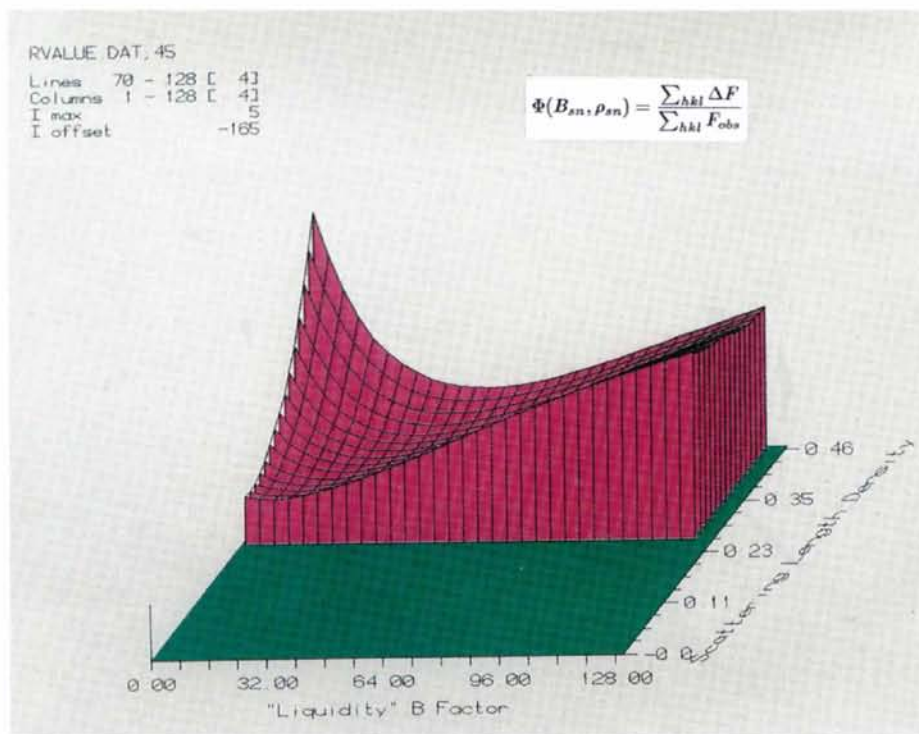
To achieve this the unit cell of the crystal is divided into a three-dimensional grid. A shell is made up of the grid points within a given range of distance from the surface of the protein: (a) grid points belong to the protein if the particular point falls within the van der Waals spheres of the atoms of protein; (b) grid points external to the protein belong to a given shell, depending on the distance from the surface of the protein. The grid points belonging to a given shell form the coordinate loci used for the calculation of structure factors for that shell. The solvent structure factors ( $F_{sn}$ ) for each shell and for each Miller index ( $hkl$ ) are calculated from the given coordinates and the average scattering density of the solvent ( $\rho_{sn}$ ) for each shell:

$$F_{sn} = \rho_{sn} \exp(-B_{sn} \sin^2 \theta / \lambda^2) \times \sum_{xyz} \exp[-2\pi i(hx + ky + lz)_n] \quad (3)$$

where  $xyz$  are the grid solvent coordinates for the  $n$ th shell. The scattering length of a solvent grid point is the volumetric scattering length ( $\rho_{sn}$ ), in  $\text{fm } \text{\AA}^{-3}$ , for the particular grid volume used, calculated from the scattering density of the solvent constituents (Table 1); all grid points within a shell have the same scattering densities.



(a)



(b)

Fig. 1. Plot of (a) the  $\chi^2$  and (b) the conventional  $R$  factor versus solvent-scattering density ( $\text{fm } \text{\AA}^{-3}$ ) and solvent liquidity ( $\text{\AA}^2$ ). The minimum  $R$  value is observed at the same point as the minimum  $\chi^2$ .

[To face p. 196

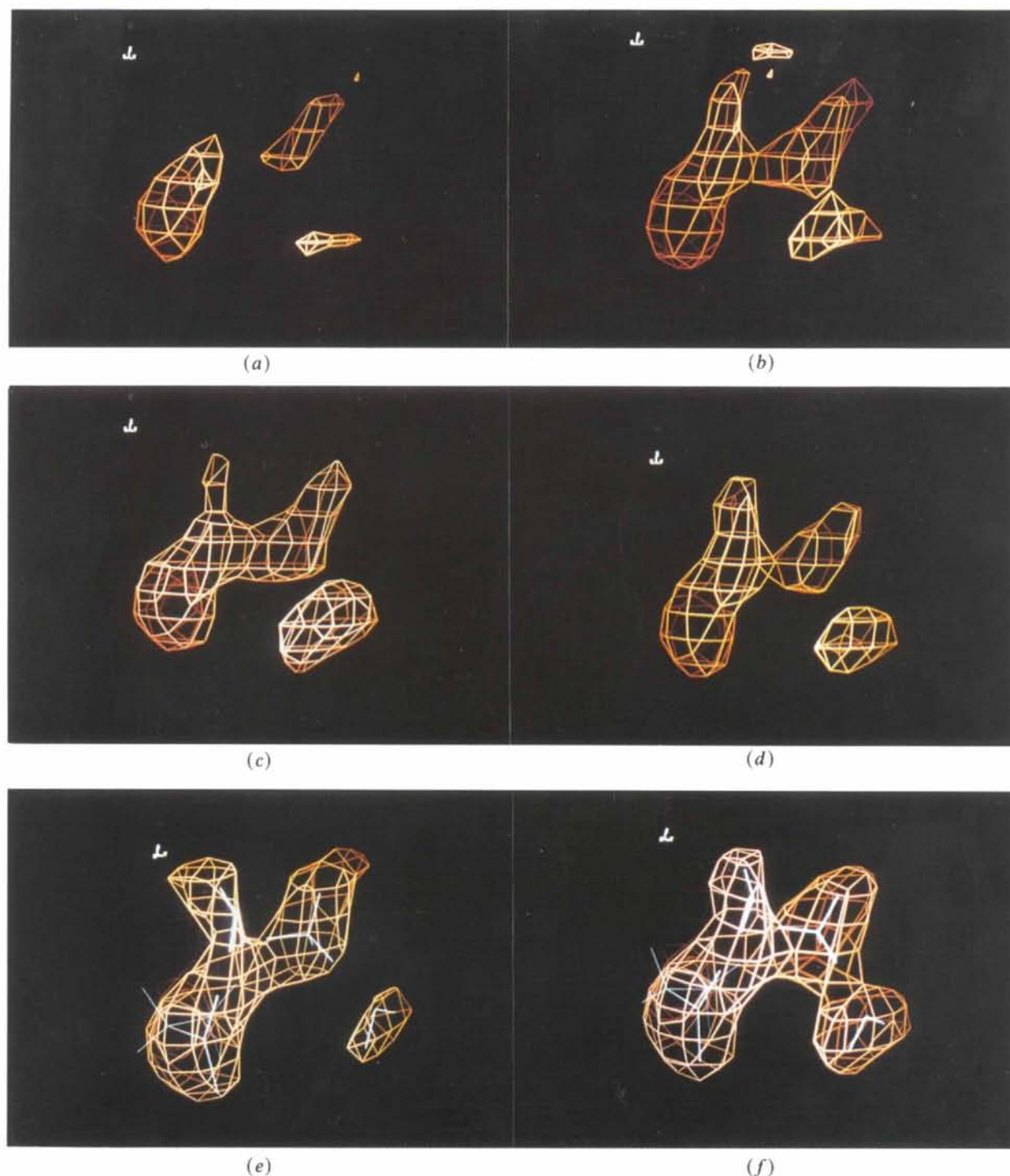


Fig. 4. Detail from the MBCO Fourier map of the surface of the protein, showing a feature on the lower left that corresponds to a COO<sup>-</sup> group of the propionic acid attached to pyrrole ring B of the heme. (a) The Fourier map was calculated with  $(2F_o - F_p)$  coefficients with protein phases ( $\Phi_p$ ) and corresponds to the classical case. (b) The total phases ( $\Phi_t$ ) [see equation (7.3)] including contribution from the whole uniform solvent were used with the same coefficients  $(2F_o - F_p)$  as in (a). (c) The total phases ( $\Phi_t$ ) using the 15 solvent-shell model were employed with the same coefficients  $(2F_o - F_p)$  as in (a). (d) The coefficients  $(2F_o - F_t)$  [see equation (7.2)] were used with the same phases ( $\Phi_t$ ) as in (c). The contours depicted above [from (a) to (d)] were interpreted to contain two water molecules ( $D_2O$ ) and an ammonium ( $ND_4^+$ ), and these molecules were added to the protein structure. (e) The coefficients  $(2F_o - F_p)$  with phases ( $\Phi_p$ ) were used after including the two water molecules and one ammonium ion in the calculation of the structure factors for  $F_p$ . (f) The coefficients  $(2F_o - F_t)$  with phases ( $\Phi_t$ ) for the revised solvent-shell model were used after adding the two water molecules and one ammonium ion to the protein structure.

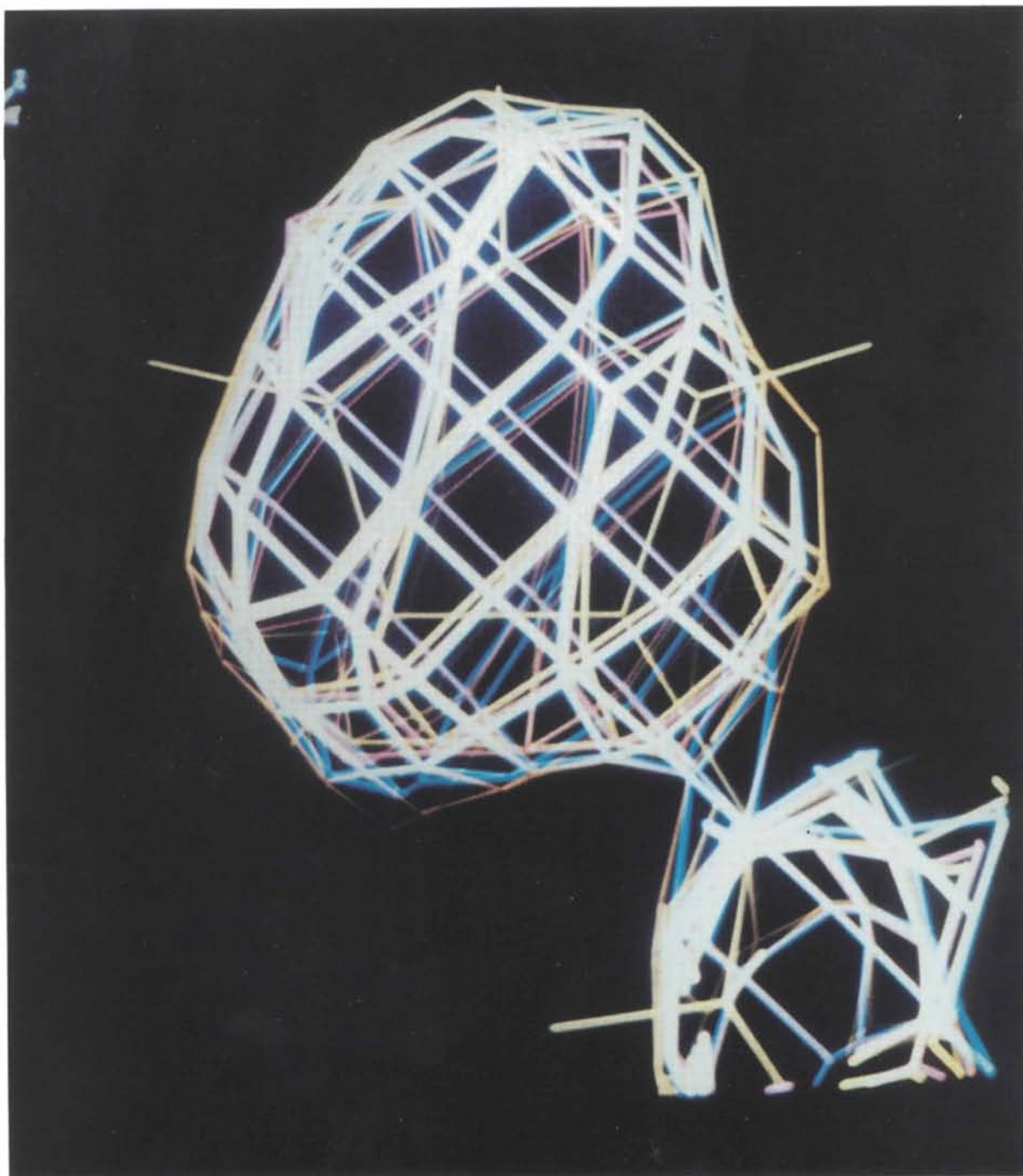


Fig. 5. Fourier sections of a histine (residue 36) in MCBO that are based on different solvent phases corresponding to Figs. 4(a) to 4(f), showing that the exterior solvent has little effect on the interior protein structure.

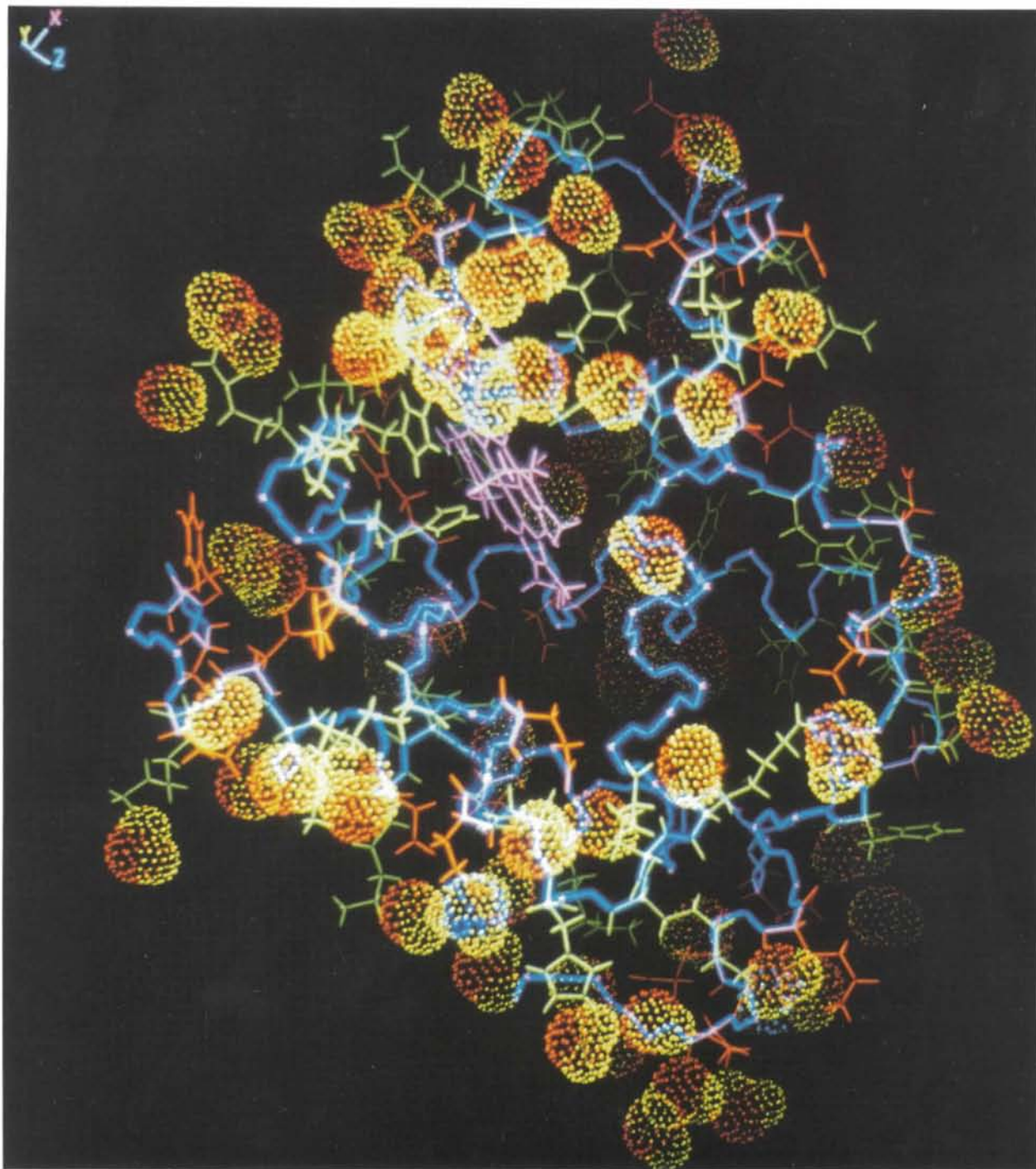


Fig. 7. A model of CO myoglobin showing the surface structure; 87 water and 5 ion sites are shown. Notice that the access path of CO (center of picture) is devoid of bound water. The heme group is in purple, the main chain in blue, acidic side chains in orange and basic side chains in green. The neutral and nonpolar residues are not shown; the purple dots in the main chain are the  $\alpha$  carbons. Water and ion molecules are indicated by space-filling dotted clusters.

Table 1. Average neutron-scattering densities of water, ammonium and sulfate ions

	Scattering density (fm Å <sup>-3</sup> )
H <sub>2</sub> O	-0.056
D <sub>2</sub> O	0.635
NH <sub>4</sub> <sup>+</sup>	-0.122
ND <sub>4</sub> <sup>+</sup>	0.814
SO <sub>4</sub> <sup>2-</sup>	0.257

Each shell has a liquidity factor ( $B_{sn}$ ) that describes its mobility (disorder). This liquidity factor is based on the root-mean-square displacement around a solvent grid point and is thus equivalent to the temperature factor. Each shell therefore has two variable parameters, the solvent-scattering density ( $\rho_{sn}$ ) and the liquidity factor ( $B_{sn}$ ). It now remains to find the best values for  $\rho_{sn}$  and  $B_{sn}$  for  $n$  shells that minimize  $\Delta F$ :

$$\Delta F(\rho_{sn}, B_{sn}) = |(F_o - F_p) - F_s(\rho_{sn}, B_{sn})|. \quad (4)$$

$F_p$  is calculated from the atomic scattering lengths and the coordinates of the protein atoms.  $F_s$  is calculated, as described above, for low-angle data until 3 Å resolution with initial liquidity factor  $B = 1$ . The best values for  $\rho_{sn}$  and  $B_{sn}$  are found by an iterative numerical procedure that changes these values in small increments (a 128 × 128 two-dimensional search) and tests for the lowest  $\Delta F$ . The behavior of this minimization is measured by a weighting function (Schoenborn, 1988):

$$\chi^2 = \sum (1/\sigma^2)(\Delta F)^2 \quad (5)$$

where  $\sigma$  is the counting statistics of the observed reflections. The conventional  $R$  value can be used with similar purpose:

$$R = \sum |\Delta F| / \sum |F_o|. \quad (6)$$

The functions are plotted as a three-dimensional plot (Fig. 1). For the first step, the best overall solvent scattering density  $\rho_s$  and the liquidity factor  $B_s$  were determined. These overall uniform values for the whole solvent area were used as initial values for all shells. Keeping all except shell  $l$  at their initial values of  $\rho_s$  and  $B_s$ , the values of  $\rho_{sl}$  and  $B_{sl}$  which gave the minimum  $\chi^2$  or  $R$  values were determined. With  $\rho_{sl}$  and  $B_{sl}$  kept for shell  $l$ , and with all shells but  $l$  and  $m$  kept at  $\rho_s$  and  $B_s$ , the best  $\rho_{sm}$  and  $B_{sm}$  were determined in the same manner. Continuing in this way, the best values of  $\rho_{sn}$  and  $B_{sn}$  for each succeeding shell were determined. This optimization is well conditioned, since the number of structure factors that contribute to this function is relatively large (~1000) compared with the few shells (~20).

Table 2. Average solvent-scattering densities of myoglobin crystals in different solvents

The crystal solvent D<sub>2</sub>O concentrations listed are the initial D<sub>2</sub>O/H<sub>2</sub>O percentages used for setting up crystallization. The value is an upper limit; the actual D<sub>2</sub>O concentration will be less due to exchange. The observed scattering density depends on the actual D<sub>2</sub>O concentration and the amount of ammonium sulfate in solution (Schoenborn, 1988).

	Scattering density (fm Å <sup>-3</sup> )	Crystal solvent
MB25	0.23 ± 0.02	~40% D <sub>2</sub> O
MB17	0.49 ± 0.02	~80% D <sub>2</sub> O
MBCO	0.56 ± 0.02	~90% D <sub>2</sub> O

### 3. Hydration layers in myoglobin crystals

This procedure is particularly powerful for neutron crystallographic analysis of proteins because the average neutron-scattering density of the solvent can be large and easily adjusted by mixing heavy and light water to obtain a given scattering density (Table 1). Neutron diffraction data for three myoglobin crystals (Table 2) with different solvent-scattering densities were collected and used to analyse the solvent liquidity and scattering density as described above. The crystals were grown from saturated ammonium sulfate solutions with MB25 in ~40% D<sub>2</sub>O, MB17 in ~80% D<sub>2</sub>O, and MBCO in ~90% D<sub>2</sub>O (MB25 and MB17 are aquamet-myoglobin crystals). The early crystallographic analysis provided the basic atomic positions of the protein (Hanson & Schoenborn, 1981; Schoenborn, 1971). The coordinates of the solvent grid were determined, using the known atomic positions of the protein. An asymmetric unit of the cell contains a total of 120 000 points at ~0.6 Å grid spacing with 44 978 points in the solvent region; in this case, a hydrogen radius of 1.6 Å was used to determine the surface of the protein, which shows that nearly 38% of the cell volume is solvent. Solvent structure factors ( $F_{sn}$ ) for fifteen shells with a grid size of ~0.6 Å were calculated to a resolution of 3.0 Å from the coordinates of all the solvent grid points, with initial scattering densities corresponding to the composition of the solvent that was used in crystallization (Table 2).

A minimization calculation of  $\chi^2$  and the  $R$  factor was carried out by changing the magnitude of the solvent-shell scattering densities and liquidity factors in small steps as discussed above. The best observed scattering densities and liquidity factors of the shells for each of the three data sets were thus determined in this manner. Fig. 2(a) shows the results with the best  $\rho_{sn}$  and  $B_{sn}$  as a function of distance from the surface of the protein. These data show that there is a layer of solvent molecules with low mobility, as given by the low liquidity factors (~25 Å<sup>2</sup>), that is located at ~2.3 Å from the center of the surface

atoms (hydrogen/deuterium). The scattering densities are highest at this distance, suggesting that the bulk mass density of the innermost layer is higher than those of solvent layers further out from the protein. The uniform scattering density for each shell, however, includes the scattering density of any ion or fraction of an ion located within that shell.

To test the effect of a finite grid size, calculations were made with a grid size down to  $\sim 0.3$  Å. The computations related to the structure factors require about 12 h of CPU time on a VAX 11/780 of the 378 205 solvent grid points from the total of 960 000 points. All three data sets yield essentially the same scattering data for the solvent region and are

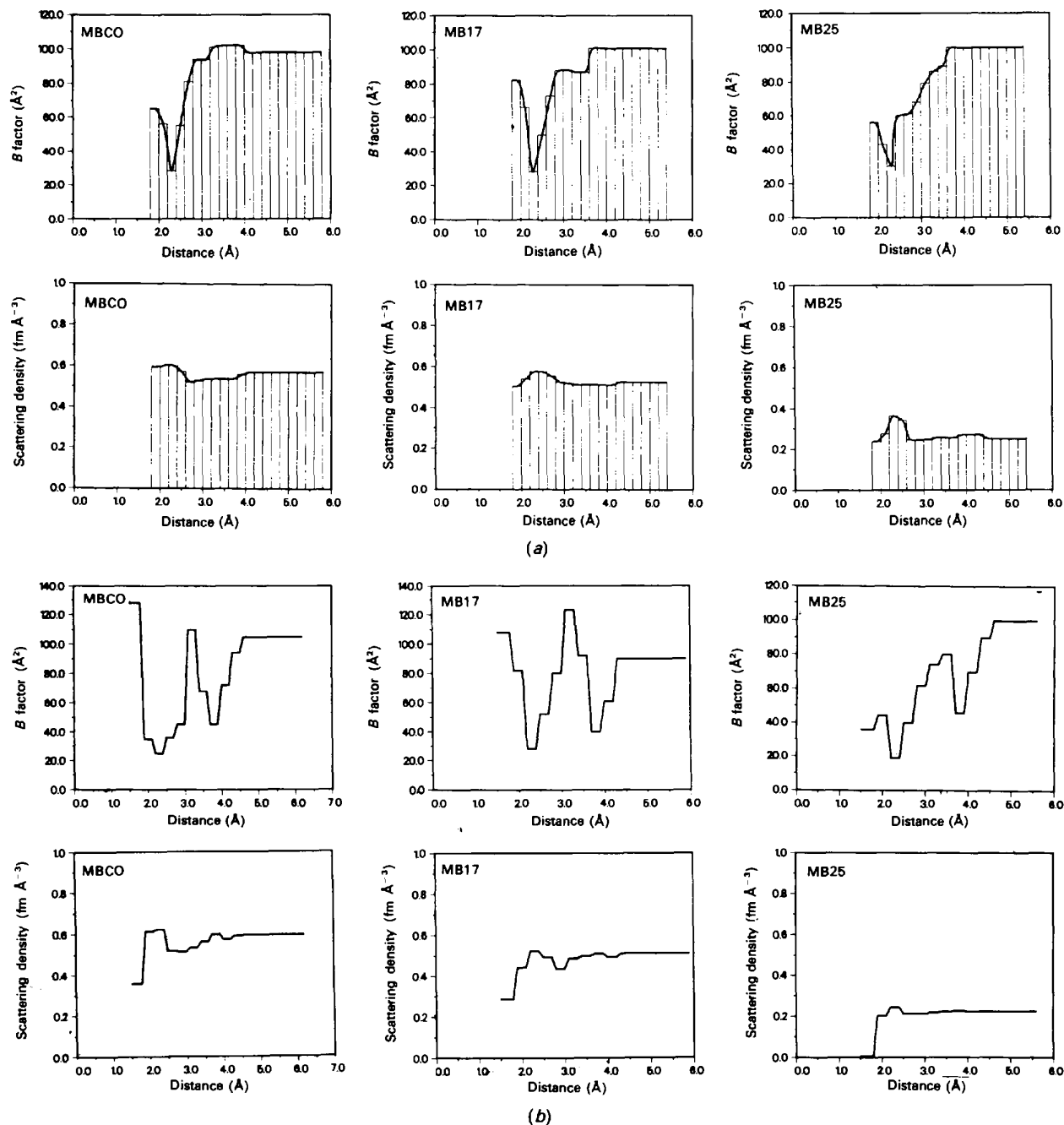


Fig. 2. Graphic presentation of the best scattering densities and  $B$  factors for three different myoglobin derivatives with different solvent compositions (see Table 2). (a) The initial radius for hydrogens is  $1.6$  Å and the shells extend outward from that radius with thickness  $0.2$  Å and a grid size of  $0.6$  Å. (b) The results based on a shell thickness of  $0.3$  Å and a grid size of  $0.3$  Å.



Table 3. *Shell-solvent analysis on MBCO*

The first row of data (total) are the best values of the overall solvent used to initiate the shell analysis. The content of solvent is about 40%. The shell thickness and the grid size are both 0.3 Å; the volume of one grid point is 0.035 Å<sup>3</sup>. The shells (11–15) of the last row of data, at 4.6 to 6.1 Å from the surface of the protein, are treated as one large shell. The radii used in calculations to determine the surface of the protein are C: 1.7, N: 1.5, O: 1.4, S: 1.85, Fe: 1.45, H: 1.6, D: 1.6 Å.

	Distance (Å)	Grid points (number)	Best $\rho_{sn}$ (fm Å <sup>-3</sup> )	Best $B_{sn}$ (Å <sup>2</sup> )
Total	1.6–6.1	378205	0.58	86
Shell 1	1.6–1.9	117962	0.36	128
2	1.9–2.2	72212	0.61	35
3	2.2–2.5	52679	0.62	25
4	2.5–2.8	40289	0.52	36
5	2.8–3.1	30754	0.51	45
6	3.1–3.4	23417	0.53	110
7	3.4–3.7	15915	0.56	68
8	3.7–4.0	11468	0.59	45
9	4.0–4.3	7017	0.57	72
10	4.3–4.6	3941	0.58	94
11–15	4.6–6.1	2551	0.59	105

independent of the grid size from 0.6 to 0.3 Å. Fig. 2(b) shows the results based on a grid size of ~0.3 Å and a shell thickness of 0.3 Å. Table 3 gives the detailed characteristics of the shell analysis on MBCO. The best overall values for the total solvent are listed in the first line of Table 3. The best average density in the overall solvent region was 0.58 fm Å<sup>-3</sup>, which agrees well with the D<sub>2</sub>O/H<sub>2</sub>O ratio of ~90% used for setting up crystallization. The innermost layer of the shell is very well defined and is separated from the second one. A second layer of water molecules is indicated by the second peak at 3.7 to 4.0 Å with liquidity factors of ~45 Å<sup>2</sup> (from Fig. 2b).

A possible picture of solvent structure which emerges from this analysis is shown schematically in Fig. 3. There is a well-defined layer of water molecules at the expected van der Waals distance from the protein surface (D/H atoms). Each water molecule in the first layer is hydrogen bonded to water molecules in the second layer. The ideal O—H...O bond length in ice (Pauling, 1967) is 2.76 Å and O—H...O angle 180°. The observed ~3.9 Å distance of the second layer from the protein surface agrees with the predicted value of 4.03 Å, Fig. 3. The scattering densities of the water layers are about the same, but the liquidity factors for the second water layer are higher suggesting that there is less order.

#### 4. Phasing, solvent structure determination and refinement

With the best parameters for the solvent shells known, modified phases for the observed structure

factors can be calculated to improve the Fourier map features on the surface of the protein that depict bound water and ion molecules. The phases for  $F_o$ 's are obtained from the calculated structure factors ( $F_i$ ) for protein and solvent shells using the expanded structure-factor form:

$$F_i = F_p + F_s(\rho_{sn}^{\text{best}}, B_{sn}^{\text{best}}) \quad (7.1)$$

and,

$$|F_i| = [(A_p + A_s)^2 + (B_p + B_s)^2]^{1/2} \quad (7.2)$$

where  $A_p$ ,  $A_s$  and  $B_p$ ,  $B_s$  are the real and imaginary terms of the protein and solvent components respectively.

The phase for  $F_o$  is given by:

$$\Phi_i = \text{tg}^{-1}[(B_p + B_s)/(A_p + A_s)]. \quad (7.3)$$

The magnitude and phase of the structure factor of the protein and solvent components are:

$$|F_p| = [(A_p^2 + (B_p)^2)]^{1/2} \quad (8.1)$$

$$\Phi_p = \text{tg}^{-1}(B_p/A_p) \quad (8.2)$$

and,

$$|F_s| = [(A_s)^2 + (B_s)^2]^{1/2} \quad (8.3)$$

$$\Phi_s = \text{tg}^{-1}(B_s/A_s). \quad (8.4)$$

The Fourier maps calculated use the observed structure factors ( $F_o$ ) and the calculated phases ( $\Phi_i$ ) based on all the known features of the particular protein. These Fourier maps are interpreted with a macromolecular graphics modeling program (*PS300 FRODO*). To show the effect of solvent phasing, the Fourier map calculated with coefficients  $2|F_o| - |F_p|$  based on phases ( $\Phi_p$ ) from the parent protein structure only, was compared with maps based on the

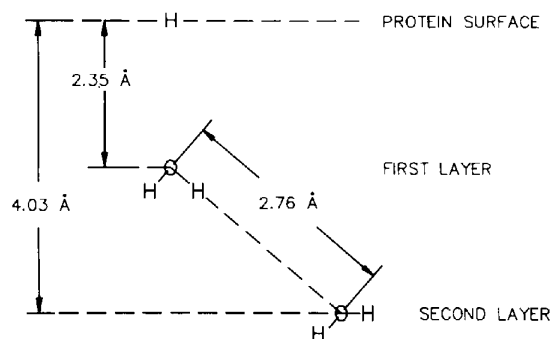


Fig. 3. Water molecules in the first and second hydration layers at the surface of the protein. 2.35 Å is the sum of van der Waals radii for the H atom (protein) and the O atom (water). 2.76 Å and the linear bond angle are the hydrogen-bond distance and angle in ice. The bisector of the H—O—H bond angle of 105° is taken perpendicular to the protein surface. (Not shown in this sketch are salt ions which can also be bound at the surface of the protein.)

total phases including the solvent structure contribution ( $\Phi_s$ ). The neutron density map ( $2|F_o| - |F_p|$ ,  $\Phi_p$ ) corresponds to the classical case (no water molecules are included in the calculated structure factors), and the map ( $2|F_o| - |F_p|$ ,  $\Phi_s$ ) corresponds to the solvent-phased density. The small Fourier section in Fig. 4 (see §6) shows the effects of including solvent phasing on the information obtainable, and demonstrates how the improved phasing model greatly facilitates the localization and orientation of water and ion molecules. The solvent-phased Fourier map shows features that are completely absent in the conventional map and, generally, yields a greatly improved topography of the surface of the protein.

These new calculated structure factors ( $F_s$ ) that represent the solvent-shell model can also be used to modify the observed structure factors ( $F_o$ ) so that they exclude the contribution of the solvent and represent only protein. These modified terms ( $F_{om}$ ) are then used to refine the protein structure in a modern reciprocal-space least-squares refinement (*PROLSQ*) that includes all structure factors:

$$F_{om} = |F_o - F_s(\rho_{sn}^{best}, B_{sn}^{best})|. \quad (9)$$

The inclusion of the low-order data in the refinement is advantageous, since these data are generally the strongest with the best counting statistics and they condition the least-squares refinement well.  $|F_o|$  is scaled to match  $|F_p|$  for the high-angle data ( $< 3 \text{ \AA}$ ), the same scale factor is used for low-angle data ( $> 3 \text{ \AA}$ ) to scale  $|F_{om}|$  to match  $|F_p|$ .

Through an inspection of the Fourier map, bound-water and ion molecules were localized and added to the protein coordinates. The D and O atoms of water were placed at stereochemically reasonable positions as indicated by the Fourier features. Another set of modified  $F_{om}$  were calculated with all of the newly located solvent molecules treated as a part of the protein. The refinement of water molecules poses some new problems. Water molecules do not always have full occupancy and, in addition, hydrogen atoms are exchangeable. In our case, the occupancy factor ( $Q$ ) for all three atoms in a water molecule is determined by the oxygen atom only, while the scattering length of an exchangeable hydrogen is given by the relation  $b = (6.67 + 3.74)f - 3.74$ , where 6.67 and  $-3.74$  are the scattering lengths, in fm, of deuterium and hydrogen, respectively, and  $f$  is the fractional hydrogen-exchange rate. The permissible range of  $f$  is from 0 (hydrogen) to 1 (deuterium). The refinement is initiated with zero scattering ( $b = 0.0$  when  $f = 0.355$ ). On completion of several cycles of refinement of the protein and water structure, the whole solvent-analysis procedure is repeated, followed by a Fourier synthesis and calculation of another set of modified  $F_{om}$  which are then used in the next round of refinement.

## 5. Water restraints in the least-squares refinement

The restrained least-squares refinement (Hendrickson & Konnert, 1981; Hendrickson, 1985) was modified to allow the refinement of the water structure by including restraints on the binding configurations in water-to-protein and water-to-water associations. The *PROLSQ* procedure of Hendrickson & Konnert (1981) was expanded to include a set of water-repulsive terms. These restrictions on water structure are based on results of Savage & Finney (1986) that describe stereochemical restraints in ices and crystalline hydrates. These interactions are characterized by interatomic potential-energy functions  $U(d)$  which can be approximated (Hendrickson, 1985) by:

$$U(d) - U(d_{min})|_{d < d_{min}} \approx (1/\sigma_v^4)(d - d_{min})^4 \quad (10)$$

where  $d_{min}$  locates the minimum energy, only repulsive contacts ( $d < d_{min}$ ) are considered, and  $\sigma_v$  is the standard deviation to be permitted.  $\sigma_v$  varies in principle with the type of contact, but in practice the values are quite uniform. The value of  $d_{min}$  depends on which atomic elements are in contact and what type of contact they form:

(a) O...O repulsion of hydrogen-bonded structures (O—H...O)

$$d_{min}|_{\alpha \geq 120} \approx -(\alpha/120) + 4.0$$

where  $\alpha = \text{O—H...O}$  angle ( $^\circ$ ), and  $\sigma_v = 0.65$  for water-to-protein and 0.80 for water-to-water contacts.

(b) H...O repulsion of hydrogen bonds (O—H...O)

$$d_{min}|_{\alpha \geq 120} \approx -(\alpha/70.0) + 4.0$$

where  $\alpha = \text{O—H...O}$  angle ( $^\circ$ ) and  $\sigma_v = 0.65$ .

(c) H...H' non-bonded repulsion (O—H...H'—O')

$$d_{min}|_{180 \leq \chi \leq 250} \approx -(\chi/220) + 3.1$$

where  $\chi = \alpha + \beta$ ,  $\alpha = \text{O—H...H'}$  angle ( $^\circ$ ),  $\beta = \text{H...H—O'}$  angle ( $^\circ$ ) and  $\sigma_v = 0.65$ .

(d) H<sub>2</sub>...O<sub>1</sub> oxygen-remote-neighbor hydrogen non-bonded interaction (H<sub>2</sub>—O<sub>2</sub>...H<sub>1</sub>—O<sub>1</sub>)

$$d_{min}|_{80 \leq \varphi \leq 113} \approx -(3\varphi/200) + 4.2$$

where  $\varphi = \text{H}_2\text{—O}_2\text{...H}_1$  angle ( $^\circ$ ) and  $\sigma_v = 0.65$ .

Beyond the range of angles ( $\alpha \geq 120$ ,  $180 \leq \chi \leq 250$ ,  $80 \leq \varphi \leq 113^\circ$ ),  $d_{min}$  will take the value of the van der Waals contact distance. Restraints have been applied to possible bonds formed by water characterized by distance limits.

## 6. Application to carbonmonoxymyoglobin crystals

The structure of MBCO was originally determined from neutron crystallographic data (Norvell, Nunes & Schoenborn, 1975) using the Diamond real-space refinement (Hanson & Schoenborn, 1981). The CO myoglobin crystal was grown from a 70% saturated

Table 4. Summary of the refinement

Every run involved several cycles. Each cycle refines either atomic coordinates  $xyz$ , or together with the refinement of one of the three factors,  $B$  (individual isotropic thermal parameter),  $f$  (hydrogen-exchange rate of H/D atom) and  $Q$  (occupancy factor of solvent molecule). Each cycle reads parameter shifts from the previous cycle and uses a conjugate-gradients procedure to solve the new parameter shifts. Finally, it tests the new shifts on the  $R$  value, searching for an optimal shift-damping factor against a sample of data. In our case, we simply alternate the application of  $xyz$ ,  $B$ ,  $f$  and  $Q$  shifts in successive cycles. The total number of cycles of each run is shown in parentheses in column 1; the number of cycles on each parameter is shown in columns 3 to 6. Only occupancy factors of solvent molecules are refined.

Run No.	$R$ (%)	$xyz$	$B$	$f$	$Q$	Atoms	Comments
1 (5 cycles)	29.3→26.7	2	1	2	—	2543	MBCO structure from real-space refinement (Hanson & Schoenborn, 1981), original observed structure factor $F_o$ used.
2 (8)	25.4→19.3	4	1	3	—	2543	Overall solvent contribution included, a modified observed structure factor $F_{om}$ used.
3 (16)	19.0→16.3	5	4	4	3	2780	Solvent-shell analysis performed, 79 water molecules added, a new $F_{om}$ used.
4 (16)	17.6→14.8	5	3	3	5	2941	Solvent-shell analysis performed, a total of 127 water and 3 ammonium ions included, a new $F_{om}$ used.
5 (28)	15.2→12.7	7	7	7	7	2821	Solvent-shell analysis performed, a total of 86 water, 3 ammonium and one sulfate ions included, a new $F_{om}$ used.
6 (8)	12.9→11.5	4	4	—	—	2829	With a relaxed weighting of the stereochemical conditions.

ammonium sulfate solution at pH 5.7, and was soaked in  $D_2O$  mother-liquor solution for several months. We re-examined this structure using the solvent-modified structure-factor terms in the expanded least-squares refinement (*PROLSQ*), as described above. The starting coordinates of the protein resulted in the following *PROLSQ* parameters: 2543 atoms, 8397 distances, 213 planes, 182 chiral centers, 763 torsion angles, 47 018 possible contacts, and 3156 non-bonded contacts in the protein itself. The total number of variables is 10 173 and there are 6520 structure-factor observations ( $>2\sigma$  for  $\infty$  to 2 Å and  $>7\sigma$  for 2 to 1.8 Å). Five cycles of refinements (two cycles on  $xyz$ , one cycle on individual isotropic thermal  $B$  factors, and two cycles on exchange ratios of H/D atoms) with the original observed  $F_o$  reduced the  $R$  from 29.3 to 26.7% (Table 4). This produced an r.m.s. shift in atomic positions for all atoms of 0.36 Å. Comparison of the coordinates with the starting set shows that most of the large shifts were in surface side-chain atoms.

The overall solvent contribution to the structure factors was then calculated using a  $\sim 0.6$  Å grid spacing and data to 3.0 Å resolution (2245 reflections). The protein structure was further refined after the overall solvent contribution was included in the

Table 5. Reflections to a resolution of 12 Å of MBCO

The definitions of each column are given in formulas (7) to (9).

$H$	$K$	$L$	Protein			Solvent		Total		$F_{om}$
			$F_o$	$F_p$	$\Phi_p$	$F_s$	$\Phi_s$	$F_t$	$\Phi_t$	
0	2	0	332.	972.	340.	758.	177.	331.	298.	973.
0	0	1	661.	1397.	14.	902.	207.	560.	352.	1492.
0	1	1	690.	857.	102.	941.	251.	484.	186.	901.
0	0	2	562.	881.	0.	1395.	180.	514.	180.	833.
1	2	-1	295.	352.	56.	237.	300.	327.	15.	328.
1	2	0	113.	267.	311.	176.	136.	94.	301.	286.
1	2	1	291.	504.	77.	233.	280.	304.	60.	492.
1	1	-2	392.	323.	100.	176.	9.	367.	71.	346.
1	1	-1	695.	639.	48.	1352.	225.	715.	222.	659.
1	1	0	1906.	76.	7.	1865.	141.	1814.	139.	70.
1	1	2	270.	396.	41.	379.	250.	193.	330.	428.
1	0	-2	478.	408.	0.	101.	180.	307.	360.	579.
1	0	-1	1270.	542.	0.	1781.	180.	1239.	180.	511.
1	0	2	346.	27.	0.	264.	180.	237.	180.	82.
2	2	0	228.	63.	91.	312.	266.	249.	265.	84.
2	2	1	588.	868.	204.	323.	17.	550.	208.	906.
2	1	-2	273.	470.	75.	230.	274.	264.	58.	479.
2	1	-1	196.	954.	187.	1317.	17.	412.	40.	1141.
2	1	0	369.	164.	269.	321.	335.	416.	314.	135.
2	1	1	448.	894.	317.	617.	166.	468.	277.	879.
2	1	2	426.	740.	34.	492.	204.	268.	52.	890.
2	0	-1	993.	248.	0.	1196.	180.	948.	180.	203.
2	0	0	1305.	330.	180.	1575.	0.	1246.	360.	270.
2	0	1	794.	97.	0.	747.	180.	650.	180.	47.
2	0	2	582.	421.	180.	168.	180.	589.	180.	414.
3	2	-1	142.	438.	349.	287.	148.	199.	21.	390.
3	2	0	484.	713.	242.	260.	74.	462.	235.	735.
3	1	-2	324.	42.	100.	346.	89.	387.	91.	22.
3	1	-1	150.	707.	79.	795.	256.	96.	234.	658.
3	1	0	431.	691.	313.	510.	146.	227.	282.	871.
3	1	1	393.	508.	128.	197.	320.	318.	120.	582.
3	0	-2	212.	846.	180.	724.	360.	122.	180.	937.
3	0	-1	209.	654.	0.	900.	180.	246.	180.	691.
3	0	0	147.	146.	180.	294.	360.	148.	360.	147.
3	0	1	283.	644.	180.	259.	360.	385.	180.	541.
4	1	-2	176.	343.	282.	261.	111.	95.	256.	417.
4	1	-1	473.	706.	244.	273.	39.	474.	258.	706.
4	1	0	347.	148.	325.	351.	19.	454.	3.	92.
4	0	-1	229.	1029.	0.	1329.	180.	301.	180.	1101.
4	0	1	545.	709.	0.	86.	180.	623.	360.	631.
5	0	-1	269.	274.	180.	50.	360.	224.	180.	320.
5	0	0	417.	572.	180.	155.	360.	418.	180.	572.

observed structure factors; *i.e.*, modified observed structure factors were calculated using the best scattering densities and liquidity factors of the total solvent, and then used in the refinement. The  $R$  value improved rapidly to 19.3% after eight cycles of refinement (run No. 2). A solvent-shell analysis was then performed. Solvent structure factors ( $|F_s|$ ,  $\Phi_s$ ) were calculated, based on the results listed in Table 3 and equations (7)–(9). A comparison of  $F_o$ ,  $F_p$ ,  $F_s$ ,  $F_t$  and  $F_{om}$  for low-resolution data ( $>12$  Å) is given in Table 5, showing the effect of the solvent's contribution. For the low-order data, the contributions from the protein and the solvent are of the same order of magnitude and show large phase shifts.  $F_t$ , which includes contributions from both the protein and solvent parts, is now directly comparable with the observed total structure factor ( $F_o$ ).

Based upon the protein coordinates and best solvent-shell scattering densities and liquidity factors from the shell analysis, a new neutron density map ( $2|F_o| - |F_p|$ ,  $\Phi_t$ ) was calculated and analyzed on an interactive graphics system, *FRODO PS300*. At this stage, the 79 largest features with densities greater

than  $0.56 \text{ fm } \text{\AA}^{-3}$  on the surface of the protein were fitted with water molecules, localizing the atoms to the neutron density while maintaining ideal molecular geometry and suitable hydrogen-bond donors and acceptors. These individual water molecules were then added to the protein coordinates and refined, applying the previously described water restraints. The coordinates were refined against the modified observed neutron structure factors ( $F_{om}$ ). The occupancies of the oxygen atoms of water and hydrogen exchange rates of the H/D atoms were refined alternately with coordinates, and after 16 cycles an  $R$  factor of 16.3% was obtained. The process of the refinement is shown in Table 4.

After run No. 3, another solvent-shell analysis was performed and a new neutron density map was calculated. The water sites were carefully analyzed and two water sites with high densities ( $>0.70 \text{ fm } \text{\AA}^{-3}$ ) were changed to ammonium ion sites. Inspection of the remaining Fourier features localized another ammonium ion and 50 more water sites, to a total of 127 water molecules. All three ammonium sites are adjacent to negative charges. Fig. 4 shows the Fourier representation of a  $\text{COO}^-$  group of one of the two propionic acids of the heme showing the effects of different phasing used to calculate Fourier maps. The contours depicted in Fig. 4(b) were interpreted as containing two  $\text{D}_2\text{O}$  and one  $\text{ND}_4^+$ . The inclusion of solvent in this analysis led to significant changes in the surface configuration of the protein, particularly in the placement of bound water and ions, without changing the interior features of the protein. An example of this is shown in Fig. 5 which depicts superimposed Fourier sections of a histidine residue C1 (36) that are based on phases calculated with different solvent phasing models corresponding to Figs. 4(a) to 4(f).

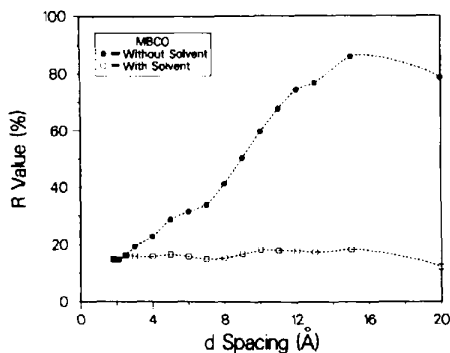


Fig. 6. A graph of the conventional  $R$  value as a function of  $d$  spacing ( $n\lambda = 2d\sin\theta$ ). The top curve represents the  $R$  value without the solvent correction (the classical case), while the bottom curve shows the  $R$  value with the correction from the solvent-shell model.

Table 6. Resolution breakdown

An average of the value of  $\sigma_F = \frac{1}{2}(F_o - F_c)$  was used as weighting sigmas instead of individual counting statistics for the observed structure factors. Columns 3 to 5 are the refinement at the  $R = 12.7\%$  stage, and columns 6 to 8 are the refinement at the present stage after eight cycles with relaxed weighting.

Resolution (Å)	No. of reflections	$\langle F_o - F_c \rangle$	Shell $R$ (%)	Sphere $R$ (%)	$\langle F_o - F_c \rangle$	Shell $R$ (%)	Sphere $R$ (%)
6.00-∞	323	43.5	9.7	9.7	40.4	9.0	9.0
5.00-6.00	221	37.6	12.0	10.5	35.4	11.3	9.4
4.00-5.00	485	35.6	10.1	10.3	33.1	9.4	9.1
3.00-4.00	1216	32.0	10.7	10.5	29.9	10.0	9.3
2.50-3.00	1414	29.0	13.6	11.4	27.5	12.9	10.2
2.10-2.50	1845	27.3	15.3	12.3	25.9	14.5	11.1
1.80-2.10	1016	25.9	16.5	12.7	23.7	15.1	11.5
1.80-∞	6520	30.1	-	12.7	28.2	-	11.5

Five more water molecules were added after run No. 4, while 45 water sites were rejected due to low occupancy ( $<0.3$ ) and one Fourier feature close to the side chain of Lys H21 (145) was interpreted as a sulfate ion. The proof that this site is indeed an  $\text{SO}_4^{2-}$  ion must wait until the analysis of the equivalent non-deuterated ( $\text{H}_2\text{O}$ ) crystal is complete. In the  $\text{H}_2\text{O}$  case, an  $\text{SO}_4^{2-}$  molecule would have the same density while a water molecule with the same occupancy would have a much lower density ( $-0.056 \text{ fm } \text{\AA}^{-3}$ , Table 1). Finally, one more water molecule and another ammonium ion joined the family of solvent molecules after run No. 5.

Fig. 6 shows the effect of including the low-order data on the conventional  $R$  value of a structural refinement. The inclusion of the solvent model greatly improved the overall  $R$  value to 12.7% (Table 6) for all reflections. Atoms, including hydrogen and oxygen atoms of solvent molecules, in this model have an r.m.s. bond distance deviation from 'ideality' (Hendrickson, 1985) of  $0.006 \text{ \AA}$  for non-hydrogen bond lengths (1-2 neighbor, notation as defined by Hendrickson),  $0.022 \text{ \AA}$  for hydrogen-related bond distances,  $0.039 \text{ \AA}$  for angle-related distances (1-3 neighbor), and  $0.027 \text{ \AA}$  for intraplanar distances (1-4 neighbor). For repulsive water contacts (notation as defined in §5) this left an r.m.s. deviation from 'ideal' distances of  $0.50 \text{ \AA}$  for water-to-protein  $\text{O}\cdots\text{O}$  distances,  $0.80 \text{ \AA}$  for water-to-water  $\text{O}\cdots\text{O}$  distances,  $0.50 \text{ \AA}$  for  $\text{H}\cdots\text{O}$  distances and  $0.44 \text{ \AA}$  for  $\text{H}\cdots\text{H}'$  distances. Eight more cycles with a relaxed weighting of the stereochemical conditions brought  $R$  to 11.5% (Tables 4 and 6). The weighting factor related to bonding distances, for example, is determined by  $\sigma_D$  (Hendrickson, 1985), the standard deviation for the distances of a particular type, which was relaxed from  $0.01$  to  $0.03 \text{ \AA}$  for non-hydrogen 1-2 bonds, from  $0.03$  to  $0.05 \text{ \AA}$  for hydrogen-related 1-2 bonds, from  $0.04$  to  $0.06 \text{ \AA}$  for 1-3 and 1-4 neighbors, etc. Atoms shifted another  $0.22 \text{ \AA}$  (r.m.s.), and no significant shifts were seen in the final two cycles.

Table 7. *Water-binding configuration: water-to-protein and water-to-water*

The hydrogen bonds (O or N...D—O) between protein-to-water and water-to-water are listed. The number after the first slash is the non-deuterium distance (O or N...O) (Å); the number after the second slash is the angle at the deuterium atom (°). The number in parentheses is the residue number of water or ion; the symbol # represents the neighboring symmetry-related molecule.

Residue number and type	Main chain N C	Side chain
1 Val <i>NA1</i>	N—D...O(254)/2-80/165 O(254)...D—O(255)/2-63/151 O(255)—D...O(221)/3-21/141... (133 Lys)	
2 Leu <i>NA2</i>	D O	
3 Ser <i>A1</i>	D O	
4 Glu <i>A2</i>	N—D...O(253)/2-60/150 O(253)...O(#174 SO <sub>4</sub> <sup>2-</sup> )/1-89... (#145 Lys)	O <sub>11</sub> ...D—O(286)/2-47/160
5 Gly <i>A3</i>	D O...O(225)/2-24 O(225)—D...O(226)/2-97/162	
6 Glu <i>A4</i>	D O	O <sub>11</sub> '...D—O(289)/3-03/124 O <sub>12</sub> '...D—O(289)/3-59/144
7 Trp <i>A5</i>	D O	
8 Gln <i>A6</i>	D O	
9 Leu <i>A7</i>	D O	
10 Val <i>A8</i>	D O	
11 Leu <i>A9</i>	D O	
12 His <i>A10</i>	D O	
13 Val <i>A11</i>	D O	
14 Trp <i>A12</i>	D O	
15 Ala <i>A13</i>	D O	
16 Lys <i>A14</i>	D O...D—O(192)/2-72/170... (#59 Glu)	
17 Val <i>A15</i>	D O	
18 Glu <i>A16</i>	D O	O <sub>12</sub> '...D—O(159)/2-54/164 O(159)...D—O(264)/2-57/156... (77 Lys)
19 Ala <i>AB1</i>	D O	
20 Asp <i>B1</i>	D O	
21 Val <i>B2</i>	D O	
22 Ala <i>B3</i>	D O	
23 Gly <i>B4</i>	D O	
24 His <i>B5</i>	D O	
25 Gly <i>B6</i>	D O	
26 Gln <i>B7</i>	D O	
27 Asp <i>B8</i>	D O	O <sub>12</sub> '...D—O(287)/2-43/151... (118 Arg) O(287)—D...O(191)/3-16/144 O <sub>12</sub> '...D—O(#256)/2-97/159
28 Ile <i>B9</i>	D O	
29 Leu <i>B10</i>	D O	
30 Ile <i>B11</i>	D O	
31 Arg <i>B12</i>	D O	N <sub>12</sub> <sup>+</sup> ...O(198)/3-29 N <sub>12</sub> <sup>+</sup> ...O(232)/2-42 O(232)—D...O(229)/2-76/154
32 Leu <i>B13</i>	D O	
33 Phe <i>B14</i>	D O	
34 Lys <i>B15</i>	D O	
35 Ser <i>B16</i>	D O	
36 His <i>C1</i>	D O	
37 Pro <i>C2</i>	— O	
38 Glu <i>C3</i>	D O...D—O(175)/3-12/143	
39 Thr <i>C4</i>	D O	
40 Leu <i>C5</i>	D O	
41 Glu <i>C6</i>	D O...D—O(251)/2-98/177	O <sub>11</sub> '...D—O(157)/2-84/157 O(157)...D—O(#264)/2-56/176... (#77 Lys)
42 Lys <i>C7</i>	D O...D—O(250)/2-65/158	N <sup>+</sup> —D...O(214)/3-50/140 O(214)—D...O(268)/2-93/177... (98 Lys)
43 Phe <i>CD1</i>	D O	
44 Asp <i>CD2</i>	N—D...O(252)/3-22/174 D O...D—O(184)/2-79/154... (48 His) O(184)—D...O(208)/2-89/170	O <sub>12</sub> '...D—O(215)/3-25/140 O(215)...D—O(216)/3-85/153
45 Arg <i>CD3</i>	N—D...O(205)/3-43/166...O <sub>101</sub> (154 Heme) D O...O(188)/2-61/165... (60 Asp)	N <sub>11</sub> <sup>+</sup> —D...O(199)/3-38/158... (64 His) N <sub>12</sub> <sup>+</sup> —D...O(222)...O(199)/3-61/143... (67 Thr) N <sub>12</sub> <sup>+</sup> ...O(242)/3-67 O(242)...D—O(241)/2-67/163
46 Phe <i>CD4</i>	D O	
47 Lys <i>CD5</i>	D O	
48 His <i>CD6</i>	D O	N <sub>11</sub> —D...O(184)/2-95/132... (44 Asp) N <sub>12</sub> —D...O(234)/3-06/150
49 Leu <i>CD7</i>	D O	
50 Lys <i>CD8</i>	D O	
51 Thr <i>D1</i>	D O	
52 Glu <i>D2</i>	D O	
53 Ala <i>D3</i>	N—D...O(233)/3-01/151	O <sub>12</sub> '...D—O(193)/2-95/153
54 Glu <i>D4</i>	N—D...O(168)/2-86/168	O <sub>11</sub> '...D—O(274)/2-93/133 O(274)...D—O(263)/3-02/147 O(263)...D—O(284)/2-80/159 O(284)—D...O(285)/3-14/152 O(284)...D...O(285)/3-14/152 O(263)...D—O(235)/3-07/152 O(235)...D—O(257)/2-69/137 O(257)...D—O(256)/2-84/154... (#27 Asp)

Table 7 (cont.)

Residue number and type	Main chain		Side chain
	N	C	
55 Met D5	D	O	
56 Lys D6	D	O...D—O(207)/3-29/157 O(207)...D—O(270)/3-82/164	N <sup>+</sup> —D...O(260)/3-27/157
57 Ala D7	D	O	
58 Ser E1	D	O	
59 Glu E2	D	O...D—N(209 ND <sub>4</sub> <sup>+</sup> )/2-90/169 N—D...O(240)/3-02/159	O <sub>1</sub> <sup>-</sup> ...N <sup>+</sup> (209 ND <sub>4</sub> <sup>+</sup> )/2-89 O <sub>2</sub> <sup>-</sup> ...D—O( # 231)/2-45/178...(119 His) O <sub>4</sub> <sup>-</sup> ...D—O( # 192)/2-69/157...(16 Lys) O <sub>61</sub> <sup>-</sup> ...D—O(188)/2-71/158...(45 Arg)
60 Asp E3	N	D...O(239)/2-74/154	
61 Leu E4	D	O	
62 Lys E5	D	O	
63 Lys E6	D	O...D—O(210)/2-68/160	
64 His E7	D	O	
65 Gly E8	D	O	
66 Val E9	D	O	
67 Thr E10	D	O	
68 Val E11	D	O	
69 Leu E12	D	O	
70 Thr E13	D	O	
71 Ala E14	D	O	
72 Leu E15	D	O	
73 Gly E16	D	O	
74 Ala E17	D	O	
75 Ile E18	D	O	
76 Leu E19	D	O	
77 Lys E20	D	O	
78 Lys EF1	D	O	N <sup>+</sup> —D...O(264)/3-17/127
79 Lys EF2	D	O	N <sup>+</sup> —D...O(161)/3-10/162
80 Gly EF3	N	D...O(181)/3-04/161...(82 His) D O...D—N <sup>+</sup> (218 ND <sub>4</sub> <sup>+</sup> )/2-80/145...(141 Asp)	N <sup>+</sup> —D...O(190)/3-09/158
81 His EF4	D	O	
82 His EF5	D	O	
83 Glu EF6	D	O	N <sub>2</sub> —D <sub>2</sub> ...O(288)/3-35/118
84 Ala EF7	D	O	N <sub>81</sub> —D <sub>81</sub> ...O(181)/2-70/178...(80 Gly)
85 Glu EF8	N	D...O(176)/3-24/161	O <sub>21</sub> <sup>-</sup> ...D—O(176)/2-95/150
86 Leu F1	D	O	
87 Lys F2	D	O	
88 Pro F3	—	O	N <sup>+</sup> —D...O(238)/2-69/170
89 Leu F4	D	O	
90 Ala F5	D	O	
91 Gln F6	D	O	
92 Ser F7	D	O	
93 His F8	D	O	
94 Ala F9	D	O	
95 Thr FG1	D	O	
96 Lys FG2	D	O	
97 His FG3	D	O	
98 Lys FG4	D	O	N <sub>81</sub> —D <sub>81</sub> ...O(203)/3-01/157 N <sup>+</sup> —D...O(213)/3-29/147 N <sup>+</sup> ...O(268)/2-78 O(268)...D—O( # 171)/3-53/173 O(268)—D...O( # 226)/3-20/165
99 Ile FG5	D	O	
100 Pro G1	—	O	
101 Ile G2	D	O	
102 Lys G3	N	D...O(162)/3-08/147 O(162)—D...O(290)/3-30/174	
103 Tyr G4	D	O	
104 Leu G5	D	O	
105 Glu G6	D	O	O <sub>2</sub> <sup>-</sup> ...D—O(280)/2-58/130 O(280)...D...O(282)/2-76/170...(147 Lys) O(280)—D...O(177)/3-36/175
106 Phe G7	D	O...D—O(262)/2-95/136	
107 Ile G8	D	O	
108 Ser G9	D	O	
109 Glu G10	D	O	
110 Ala G11	D	O	
111 Ile G12	D	O	
112 Ile G13	D	O	
113 His G14	D	O	
114 Val G15	D	O	
115 Leu G16	D	O	
116 His G17	D	O	
117 Ser G18	D	O	
118 Arg G19	D	O	
119 His GH1	D	O	N <sub>61</sub> —D <sub>61</sub> ...O(261)/2-78/158
120 Pro GH2	—	O	
121 Gly GH3	D	O	
122 Asp GH4	D	O	
123 Phe GH5	D	O	
124 Gly GH6	N	D...N <sup>+</sup> (178 ND <sub>4</sub> <sup>+</sup> )/3-05/144...(126 Asp)	N <sub>2</sub> —D <sub>2</sub> ...O(287)/2-83/159...(27 Asp) N <sub>82</sub> —D <sub>82</sub> ...O(287)/3-26/2-50/133 N <sub>81</sub> —D <sub>81</sub> ...O(231)/2-38/149...( # 59 Glu) O(231)—D...O <sub>82</sub> (122 Asp)/3-11/131
125 Ala H1	D	O	O <sub>82</sub> <sup>-</sup> ...D—O(231)/3-11/131...(119 His)

Table 7 (cont.)

Residue number and type	Main chain		Side chain
	N	C	
126 Asp H2	D	O	O <sub>11</sub> ...D—O(172)/3-17/140 O <sub>11</sub> ...D—O(173)/2-31/152 O <sub>22</sub> ...D—O(172)/3-07/146 O <sub>22</sub> ...D—N <sup>+</sup> (178 ND <sub>4</sub> <sup>+</sup> )/2-34/155...(124 Gly) O(172)...D—O(186)/2-78/150 O(173)...D—O(171)/2-73/159
127 Ala H3	D	O	
128 Gln H4	D	O	O <sub>1</sub> ...D—O(189)/3-05/143 O(189)—D—O(223)/3-06/153...O <sub>21</sub> (#154 Heme) O(189)...D—O(261)/2-70/158...(113 His) O(223)...D—O(196)/2-91/153 O(196)...D—O(261)/2-59/145
129 Gly H5	D	O	
130 Ala H6	D	O	
131 Met H7	D	O	
132 Asn H8	D	O	
133 Lys H9	D	O	N <sub>22</sub> ...D <sub>222</sub> ...O(182)/3-27/124...(136 Glu) N <sup>+</sup> ...O(183)/3-31/156...(126 Asp)
134 Ala H10	D	O	
135 Leu H11	D	O	
136 Glu H12	D	O	O <sub>22</sub> ...D—O(182)/2-61/155...(132 Asn) O <sub>21</sub> ...D—O(221)/4-20/145
137 Leu H13	D	O	
138 Phe H14	D	O	
139 Arg H15	D	O	
140 Lys H16	D	O	
141 Asp H17	D	O	O <sub>21</sub> ...D—O(227)/3-27/156 O <sub>22</sub> ...D—N <sup>+</sup> (218 ND <sub>4</sub> <sup>+</sup> )/2-91/146...(80 Gly) N <sup>+</sup> (218 ND <sub>4</sub> <sup>+</sup> )—D—O(276)/2-63/155 O(276)—D—O(247)/3-01/139
142 Ile H18	D	O	
143 Ala H19	D	O	
144 Ala H20	D	O	
145 Lys H21	D	O	N <sup>+</sup> —D—O(237)/3-39/138 N <sup>+</sup> ...O(174 SO <sub>4</sub> <sup>2-</sup> )/3-18
146 Tyr H22	D	O	
147 Lys H23	D	O	N <sup>+</sup> —D—O(282)/2-60/160 O <sub>21</sub> ...D—O(164)/2-74/175...(#152 Gln) O(164)...D—O(#162)/3-47/174...(#102 Lys) O <sub>21</sub> ...D—O(167)/2-97/124 O <sub>22</sub> ...D—O(167)/3-40/151
148 Glu H24	D	O	
149 Leu HC1	D	O	
150 Gly HC2	D	O	
151 Tyr HC3	D	O	
152 Gln HC4	D	O	N <sub>22</sub> ...D—O(220)/3-41/166 O <sub>21</sub> ...D—O(#164)/3-24/128...(#148 Glu) O(163)...D—O(236)/3-38/141 O <sub>212</sub> ...D—O(271)/2-40/173 O <sub>212</sub> ...D—N <sup>+</sup> (283 ND <sub>4</sub> <sup>+</sup> )/2-97/152 O <sub>21</sub> ...D—O(195)/3-14/175 O <sub>21</sub> ...D—O(205)/3-21/155...(45 Arg) N <sup>+</sup> (283)—D—O(197)/2-60/145 O(197)...D—O(195)/3-50/144 O(195)—D—O(252)/3-11/158...(44 Asp) O(195)...D—O(203)/2-61/149...(97 His) O(252)—D—O(250)/3-21/170...(42 Lys) O(205)—D—O(252)/3-05/140 O <sub>21</sub> ...D—O(#223)/3-38/156
153 Gly HC5	D	O	
154 Heme	O <sub>1</sub>	O <sub>2</sub>	O <sub>1</sub> ...D—O(163)/3-65/147

The 87 water sites and 5 ion sites that were determined (Fig. 7) are described in Table 7.\* In all, the 87 water and 5 ion molecules correspond to about 30% of the solvent in the crystal. Of the 87 water molecules, 37 are bound to side-chain atoms, 16 are bound to main-chain atoms, and 12 are in

\* Atomic coordinates and structure factors have been deposited with the Protein Data Bank, Brookhaven National Laboratory (Reference: 2MB5, R2MB5SF), and are available in machine-readable form from the Protein Data Bank at Brookhaven or one of the affiliated centres at Melbourne or Osaka. The data have also been deposited with the British Library Document Supply Centre as Supplementary Publication No. SUP 37032 (as microfiche). Free copies may be obtained through The Technical Editor, International Union of Crystallography, 5 Abbey Square, Chester CH1 2HU, England.

bridges between protein atoms (ten intra- and two intermolecular). Twenty-two water molecules are bound only to other water molecules. All water molecules bound to protein are hydrogen bonded to polar or charged groups. Twelve main-chain amide deuteriums form hydrogen bonds with water or ion molecules. Fifteen peptide carbonyl groups are hydrogen bonded to water or ion molecules. Only hydrophilic side chains form hydrogen bonds to water or ion molecules. Analysis of the polar nature of the amino-acid side chains by residue type (Table 8) suggests that the affinity for water occurs in the order acidic > basic > polar > nonpolar. Six out of the seven (86%) Asp residues are hydrogen bonded, and 11 out of 14 (79%) Glu residues are hydrogen

Table 8. *Water affinity of amino-acid side chains and main chains*

Amino acid	Side chain			Relative hydrophilicity		Peptide backbone	
	Number <sup>a</sup>	Number <sup>b</sup> observed	This study (%) <sup>c</sup>	Sweet & Eisenberg (OMH scale) <sup>d</sup>	Wolfenden <i>et al.</i> (kcal mol <sup>-1</sup> ) <sup>e</sup>	ND % (no.) <sup>f</sup>	CO % (no.) <sup>f</sup>
Asp	7	6	86	1.31	13.31	29 (2)	29 (2)
Glu	14	11	79	1.22	12.58	29 (4)	21 (3)
Arg	4	3	75	0.59	22.31	25 (1)	25 (1)
His	12	7	58	0.64	12.62	—	—
Lys	19	10	53	0.67	11.91	5 (1)	26 (5)
Asn	1	1	100	0.92	12.07	—	—
Gln	5	2	40	0.91	11.77	—	—
Thr	5	1	20	0.28	7.27	—	—
Tyr	3	—	—	-1.67	8.50	—	—
Ser	6	—	—	0.55	7.45	—	—
Gly	11	—	—	0.67	0	18 (2)	18 (2)
Phe	6	—	—	-1.92	3.15	17 (1)	—
Ala	18	—	—	0.40	0.45	6 (1)	—

Notes: (a) Number of residues existing in CO myoglobin. (b) Number of residues that form hydrogen bonds to water. (c) Number in column three divided by number in column two gives the percentage in this column. (d) The hydrophilicity (OMH) scale was derived by Sweet & Eisenberg (1983) to optimize the match between similar sequences in three-dimensional protein structures. The scale has been normalized with a mean of 0 and a standard deviation of 1.0. (e) Measured by Wolfenden *et al.* (1979, 1981) from the distribution coefficient between H<sub>2</sub>O and vapor for a small molecule representing only the side chain of the amino acid at pH 7; the value is expressed relative to that of Gly. (1 kcal = 4.187 kJ.) (f) Values given are the percentage of peptides which form hydrogen bonds to water, and the values in parentheses are the number of peptides.

bonded to water, *etc.* There are no waters bound to the side chains of Tyr and Ser residues, but they are involved in intraprotein hydrogen bonding. The hydrophilicity of polar side chains is expressed relative to that of nonpolar, say Gly, chains and given in Table 8. (The hydrophobicity of an amino-acid side chain is a measure of its lack of affinity for water while the hydrophilicity is the converse.) Comparisons with the refined optimal matching hydrophilicity scale (Sweet & Eisenberg, 1983) and the measurement of the hydrophilicity of the amino acid by using analogous small molecules (Wolfenden, Cullis & Southgate, 1979; Wolfenden, Andersson, Cullis & Southgate, 1981) shows a certain degree of agreement (correlation), and the values are also given in Table 8 (columns 5 and 6). Not surprisingly, the hydrophilicities of the peptides (Table 8) are related to the hydrophilicities of their side chains. Table 8 shows that there is a quantitative relationship for polar residues: the average hydrophilicity of a peptide is one third of that of its side chain. Asp CD2 (44) and Arg CD3 (45) have both amide deuterium and oxygen atoms of their peptides bound to water molecules. The negatively charged side chain and the carbonyl oxygen of Glu E2 (59) are bound to an ammonium ion, and its amide is bound to a water molecule. Gly EF3 (80) and Gly GH6 (124) are in contact with ions, through which they are bridged to the side chains of Asp H17 (141) and Asp H2 (126), respectively. The amide deuterium of Gly EF3 (80) is bound to a water molecule. The neighboring residues of Gly A3 (5), Ala D3 (53) and Phe G7 (106) are Gly residues; the former three peptides form bonds to water. It should be noted again that MBCO crystals were grown from a saturated ammonium sulfate solution at pH 5.7.

Table 7 describes the water hydrogen-bonding geometry: the distance between donor and acceptor, and the angle between donor, hydrogen and acceptor. The only criterion used to define hydrogen bonds was that the donor-to-acceptor distance had to be less than 3.5 Å. Distributions of the distances and angles for these hydrogen bonds that satisfy this condition are given in Fig. 8 and Table 9. Hydrogen-bonded structures are divided into three groups according to the target atom of the bond: bonds between water and main-chain atoms, water and side-chain atoms, and water and water. The O...O distances between water and main chains vary between ~2.60 and 3.40 Å, and hydrogen-bond angles vary between ~130 and 180°. The average values of distances and angles in this water and main-chain group are 2.96 Å and 158°. The O...O distances between water and water vary between

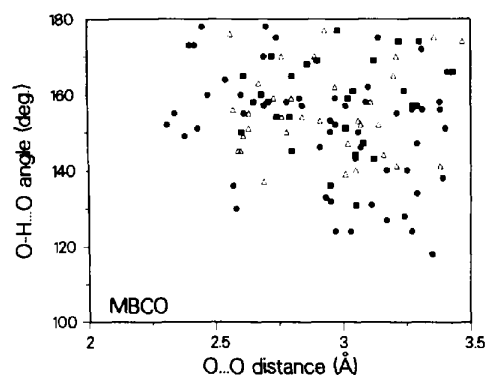


Fig. 8. Plot of O—H...O angle versus O...O distance: ■, hydrogen bonds between water and main-chain atoms; ●, hydrogen bonds between water and side-chain atoms; △, hydrogen bonds between water and water.



Table 9. Observed hydrogen-bond summary

The summary is based on a cutoff of a maximum donor-to-acceptor distance of 3.50 Å.

Group	O...O distance (Å)			O—D...O angle (°)		
	Min.	Max.	Ave.	Min.	Max.	Ave.
Water-to-main	2.60	3.43	2.96 ± 0.22	131	177	158 ± 12
Water-to-side	2.31	3.41	2.94 ± 0.39	118	178	151 ± 14
Water-to-water	2.56	3.47	2.92 ± 0.34	137	177	156 ± 11

Table 10. Distance geometry of the ions

The number in parentheses is the residue number; water and ion molecules are treated like a residue in the refinement.

Atom pair	Distance (Å)
N(283 ND <sub>1</sub> )...O <sub>β2</sub> (154 heme)	2.97
N(283) ...O(197 water)	2.60
N(283) ...O(271 water)	2.35
N(218 ND <sub>1</sub> )...O <sub>β2</sub> (141 Asp)	2.91
N(218) ...O(80 main)	2.80
N(218) ...O(276 water)	2.63
N(178 ND <sub>1</sub> )...O <sub>β2</sub> (126 Asp)	2.34
N(178) ...D(124 main)	2.30
N(178) ...D(172 water)	2.91
N(209 ND <sub>1</sub> )...O <sub>β1</sub> (59 Glu)	2.89
N(209) ...O(59 main)	2.90
S(174 SO <sub>4</sub> <sup>2-</sup> ) ...N <sup>+</sup> (145 Lys)	3.55
O(174) ...D(145 ND <sub>1</sub> )	2.18
O(174) ...O(253 water)	1.89

~2.56 and 3.47 Å, and hydrogen-bond angles vary between ~140 and 180°. The average values of distances and angles in this group are 2.92 Å and 156°. We found a wider range of distances and angles between water and side chains; O...O distances vary between ~2.30 and 3.40 Å, and angles vary between ~120 and 180°. The average values of distances and angles in this group are 2.94 Å and 151°. The three groups have almost identical averages. The average hydrogen-bond distance is longer than that found in ice (Pauling, 1967) (2.76 Å), and the average angle deviates from those for linear hydrogen bonds. The environment of water is varied and some water molecules form networks containing up to seven molecules. D<sub>2</sub>O (199), for example, is surrounded by side chains of three residues: its oxygen atom is directed towards the side chain of Arg CD3 (45) and the N<sub>δ1</sub>—D<sub>δ1</sub> group on imidazole His E7 (64); one of its deuterium atoms is directed towards the side chain of Thr E10 (67). This water molecule site is close to the location for a sulfate ion postulated in an early X-ray study by Takano (1977a,b) on metmyoglobin; the neutron density in this site is much higher than that for a sulfate ion can be (0.257 fm Å<sup>-3</sup>, Table 1). The water network between Glu D4 (54) and neighboring symmetry-related residue Asp B8 (27) contains seven molecules. Another seven-membered water network connects one of the two propionic acids of the heme to side chains of Lys C7 (42), Asp CD2 (44), Arg CD3 (45) and His FG3

(97). A ring of four water molecule sits above the plane of His A10 (12) with approach distances from ~3.2 to ~3.7 Å.

The stereochemistry of four ammonium ions and one sulfate ion is listed in Table 10. Three ammonium ions are in bridges between protein atoms. Ammonium ion (178) is in contact with the side chain of Asp H2 (126) and the main-chain amide group of residue Gly GH6 (124). Ammonium ion (209) contacts directly to the side chain of Glu E2

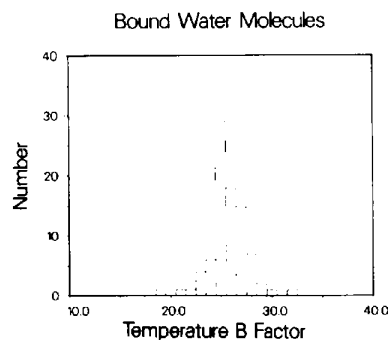


Fig. 9. Histogram of temperature factors of observed discrete water molecules.

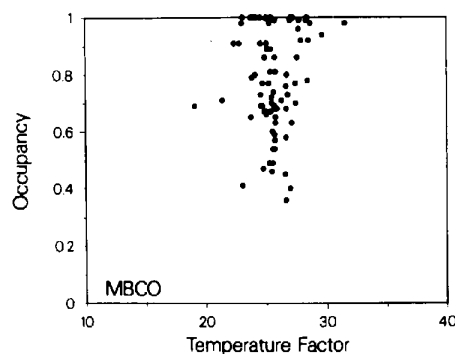


Fig. 10. Plot of the observed solvent-occupancy values against their respective temperature factors.

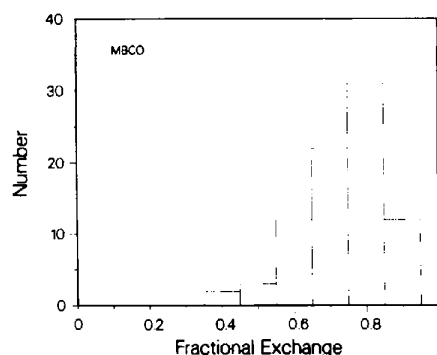


Fig. 11. Distribution of the fractional hydrogen-exchange rate of H/D atoms of observed bound water and ion molecules.

(59), and one ammonium D binds to the main-chain oxygen of same residue. Ammonium ion (218) is in contact with the side chain of Asp H17 (141), with one ammonium D bound to the main-chain oxygen of residue Gly EF3 (80) and another bound to a water molecule. One of the two propionic acids of the heme is bound to one ammonium ion (283), with two ammonium D atoms bound to two water molecules. One sulfate ion is in a channel between symmetry-related protein molecules. Intermolecular contacts are made by the side chain of Lys H21 (145) in contact with the sulfate ion, and by contacts *via* a water molecule and main-chain amide group of residue Glu A2 (4).

A histogram of temperature factors of discrete water molecules is shown in Fig. 9. There is a clear peak at  $B = 25 \text{ \AA}^2$ , corresponding to the root-mean-square vibration of  $0.56 \text{ \AA}$ . Occupancy parameters were refined and are plotted against the temperature-factor values (Fig. 10). All solvent sites with occupancy factors  $Q$  of less than 0.3 were eliminated. The distribution of the observed hydrogen exchange of bound water molecules is shown in Fig. 11. There is a clear peak at  $f = 0.80$ , measured as the fractional deuterium present. This value agrees well with the initial overall solvent-scattering density  $D_2O/H_2O$  that was used in crystallization of the MBCO crystals.

The detailed structure of MBCO and the comparisons of MBCO with the X-ray structures of metmyoglobin (Takano, 1977a), deoxymyoglobin (Takano, 1977b), oxymyoglobin (Phillips, 1980) and carbonmonoxymyoglobin (Kuriyan, Wilz, Karplus & Petsko, 1986) will be described elsewhere. The original MBCO data was collected in collaboration with A. C. Nunes and J. C. Norvell. The original myoglobin coordinates were provided by J. C. Kendrew and H. C. Watson.

We are grateful to Dr A. S. Brill for helpful discussions and critically reading the manuscript. We thank Drs A. D. Woodhead and S. W. White for critically reading the manuscript. This research was done under the auspices of the Office of Health and

Environmental Research and calculations were performed under the supercomputing program of the US Department of Energy.

#### References

- COOKE, R. & KUNTZ, I. D. (1974). *Annu. Rev. Biophys. Bioeng.* **3**, 95–126.
- EDSALL, J. T. & MCKENZIE, A. A. (1983). *Adv. Biophys.* **16**, 53–183.
- FINNEY, J. L. (1979). *Water: A Comprehensive Treatise*, Vol. 6, edited by F. FRANKS, p. 47. New York: Plenum.
- GRIFFIN, J. F. (1986). Editor. *Am. Crystallogr. Assoc. Trans.* **22**.
- HANSON, J. C. & SCHOENBORN, B. P. (1981). *J. Mol. Biol.* **153**, 117–146.
- HENDRICKSON, W. A. (1985). *Methods Enzymol.* **115**, 252–270.
- HENDRICKSON, W. A. & KONNERT, J. H. (1981). *Biomolecular Structure, Function, Conformation and Evolution*, Vol. 1, edited by R. SRINIVASAN, pp. 43–57. Oxford: Pergamon.
- IBEL, K. & STUHRMANN, H. B. (1975). *J. Mol. Biol.* **93**, 255–265.
- KELLENBERGER, E. (1978). *Trends Biochem. Sci.* **3**, N135–137.
- KOSSIAKOFF, A. A. (1985). *Annu. Rev. Biochem.* **54**, 1195–1227.
- KUNTZ, I. D. & KAUFMANN, W. (1974). *Adv. Protein Chem.* **28**, 239–345.
- KURIYAN, I., WILZ, S., KARPLUS, M. & PETSKO, G. A. (1986). *J. Mol. Biol.* **192**, 133–154.
- LEHMANN, M. S., MASON, S. A. & MCINTYRE, G. J. (1985). *Biochemistry*, **24**, 5862–5869.
- LEHMANN, M. S. & STANSFIELD, R. F. D. (1989). *Biochemistry*, **28**, 7028–7033.
- NORVELL, J. C., NUNES, A. C. & SCHOENBORN, B. P. (1975). *Science*, **190**, 568–570.
- PAULING, L. (1967). *The Chemical Bond*. Ithaca: Cornell Univ. Press.
- PHILLIPS, S. E. V. (1980). *J. Mol. Biol.* **142**, 531–554.
- PHILLIPS, S. E. V. (1984). *Neutrons in Biology, Basic Life Sciences*, Vol. 27, edited by B. P. SCHOENBORN, pp. 305–322. New York: Plenum.
- SAVAGE, H. (1986). *Water Science Reviews*, Vol. 2, edited by F. FRANKS, pp. 1–82. Cambridge Univ. Press.
- SAVAGE, H. F. J. & FINNEY, J. L. (1986). *Nature (London)*, **322**, 717–720.
- SCHOENBORN, B. P. (1971). *Cold Spring Harbor Symp. Quant. Biol.* **36**, 569–575.
- SCHOENBORN, B. P. (1988). *J. Mol. Biol.* **201**, 741–749.
- SCHOENBORN, B. P. (1989). Unpublished results.
- SWEET, R. M. & EISENBERG, D. (1983). *J. Mol. Biol.* **171**, 479–488.
- TAKANO, T. (1977a). *J. Mol. Biol.* **110**, 537–568.
- TAKANO, T. (1977b). *J. Mol. Biol.* **110**, 569–584.
- WOLFENDEN, R., ANDERSSON, L., CULLIS, P. M. & SOUTHGATE, C. C. B. (1981). *Biochemistry*, **20**, 849–855.
- WOLFENDEN, R. V., CULLIS, P. M. & SOUTHGATE, C. C. B. (1979). *Science*, **206**, 575–577.