

Rapid Communication

New expression of bimodal phase distributions in direct-method phasing of protein single-wavelength anomalous diffraction data*

Zhang Tao(张涛)^{a)b)}, Gu Yuan-Xin(古元新)^{b)}, Zheng Chao-De(郑朝德)^{b)}, and Fan Hai-Fu(范海福)^{b)†}

^{a)}Research Institute of Magnetic Materials, School of Physical Sciences and Technology, Lanzhou University, Lanzhou 730000, China

^{b)}Beijing National Laboratory for Condensed Matter Physics and Key Laboratory of Soft Matter Physics, Institute of Physics, Chinese Academy of Sciences, Beijing 100190, China

(Received 21 April 2010; revised manuscript received 29 April 2010)

One of the essential points of the direct-method single-wavelength anomalous diffraction (SAD) phasing for proteins is to express the bimodal SAD phase distribution by the sum of two Gaussian functions peaked respectively at $\varphi_{\mathbf{h}}'' + |\Delta\varphi_{\mathbf{h}}|$ and $\varphi_{\mathbf{h}}'' - |\Delta\varphi_{\mathbf{h}}|$. The probability for $\Delta\varphi_{\mathbf{h}}$ being positive (P_+) can be derived based on the Cochran distribution in direct methods. Hence the SAD phase ambiguity can be resolved by multiplying the Gaussian function peaked at $\varphi_{\mathbf{h}}'' + |\Delta\varphi_{\mathbf{h}}|$ with P_+ and multiplying the Gaussian function peaked at $\varphi_{\mathbf{h}}'' - |\Delta\varphi_{\mathbf{h}}|$ with $P_- (= 1 - P_+)$. The direct-method SAD phasing has been proved powerful in breaking SAD phase ambiguities, in particular when anomalous-scattering signals are weak. However, the approximation of bimodal phase distributions by the sum of two Gaussian functions introduces considerable errors. In this paper we show that a much better approximation can be achieved by replacing the two Gaussian functions with two von Mises distributions. Test results showed that this leads to significant improvement on the efficiency of direct-method SAD-phasing.

Keywords: direct methods, single-wavelength anomalous diffraction, OASIS program, proteins**PACC:** 6110M, 8715

1. Introduction

The single-wavelength anomalous diffraction (SAD) method is now becoming the first choice of *de novo* protein-structure determination. However there is the phase ambiguity intrinsic to the SAD phasing, i.e., the phase of reflections derived from SAD experiments is not unique, but rather a phase doublet, which can be expressed as

$$\varphi_{\mathbf{h}} = \varphi_{\mathbf{h}}'' \pm |\Delta\varphi_{\mathbf{h}}|, \quad (1)$$

where $\varphi_{\mathbf{h}}$ is the phase of reflection with reciprocal vector \mathbf{h} ; $\varphi_{\mathbf{h}}''$ is the phase of

$$F_{\mathbf{h}}'' = |F_{\mathbf{h}}''| \exp(i\varphi_{\mathbf{h}}'') = i \sum_{j=1}^N \Delta f_j'' \exp(i2\pi\mathbf{h} \cdot \mathbf{r}_j),$$

which is the structure factor contributed from the imaginary-part scattering of the heavy-atom substructure;

$|\Delta\varphi_{\mathbf{h}}|$ can be calculated from

$$|\Delta\varphi_{\mathbf{h}}| = \left| \cos^{-1} \left(\frac{|F_{\mathbf{h}}^+|^2 - |F_{\mathbf{h}}^-|^2}{4\langle F_{\mathbf{h}} \rangle |F_{\mathbf{h}}''|} \right) \right|, \quad (2)$$

where $\langle F_{\mathbf{h}} \rangle \simeq (|F_{\mathbf{h}}^+| + |F_{\mathbf{h}}^-|)/2$.

A number of procedures have been proposed to break the phase ambiguity. Ramachandran & Raman^[1] in 1956 proposed that between the two equally possible phases of the doublet, one can make a choice of that phase which is closer to the phase of the real-part scattering of the heavy-atom substructure. In 1981 Hendrickson and Teeter^[2] used a similar but improved method to solve the protein crambin from the sulfur-SAD data collected with Cu-K α x-rays. In this work, the phase of a particular reflection was determined as follows. If the Sim^[3] distribution (Eq. (3)) clearly favoured one of the maxima in the bimodal SAD-phase distribution (Eq. (4)), then the unimodal distribution of Eq. (3) was used directly. Otherwise,

*Project supported by the Innovation Foundation of the Chinese Academy of Sciences and by the National Basic Research Program of China (Grant No. 2002CB713801).

†Corresponding author. E-mail: fanhf@cryst.iphy.ac.cn

© 2010 Chinese Physical Society and IOP Publishing Ltd

<http://www.iop.org/journals/cpb> <http://cpb.iphy.ac.cn>

equations (3) and (4) were multiplicatively combined to give the distribution of the phase $\varphi_{\mathbf{h}}$. SAD phas-

ing algorithms of many modern programs are based on the same principle of the second treatment

$$P_{\text{Sim}}(\varphi_{\mathbf{h}}) = N \exp[\chi \cos(\varphi_{\mathbf{h}} - \varphi'_{\mathbf{h}})] \simeq N \exp \left\{ \chi \cos \left[\left(\varphi_{\mathbf{h}} - \left(\varphi''_{\mathbf{h}} - \frac{\pi}{2} \right) \right) \right] \right\}, \quad (3)$$

$$P_{\text{SAD}}(\varphi_{\mathbf{h}}) = N' \exp\{-[\Delta F - 2|F''_{\mathbf{h}}| \sin(\varphi_{\mathbf{h}} - \varphi'_{\mathbf{h}})]^2/2E^2\} \\ \simeq N' \exp\{-[\Delta F - 2|F''_{\mathbf{h}}| \cos(\varphi_{\mathbf{h}} - \varphi''_{\mathbf{h}})]^2/2E^2\}. \quad (4)$$

In Eqs. (3) and (4), N and N' are the normalizing coefficients of the corresponding probability distribution; χ is related to structure-factor magnitudes of the whole unit cell, the heavy-atom substructure and the unknown part of the unit cell; $\varphi'_{\mathbf{h}}$ is the phase of $\mathbf{F}'_{\mathbf{h}} = \sum_{j=1}^N (f_j + \Delta f'_j) \exp(i2\pi\mathbf{h} \cdot \mathbf{r}_j)$, which is the structure factor contributed from the real-part scattering of the heavy-atom substructure (notice that $\varphi'_{\mathbf{h}} = \varphi''_{\mathbf{h}} - 90^\circ$ when there is only one kind of heavy atoms); $\Delta F = |F_{\mathbf{h}}^+| - |F_{\mathbf{h}}^-|$ is the Bijvoet difference; $|F''_{\mathbf{h}}|$ is the magnitude of $\mathbf{F}''_{\mathbf{h}} = i \sum_{j=1}^N \Delta f''_j \exp(i2\pi\mathbf{h} \cdot \mathbf{r}_j)$; $E = (\sigma_{\Delta F}^2 + E_0^2)$, where σ is the standard deviation and E_0 is the residual lack-of-closure error.^[4] It turns out from Eq. (4) that the SAD phase distribution is bimodal and, since $\Delta F \simeq 2|F''_{\mathbf{h}}| \cos \Delta\varphi_{\mathbf{h}}$,^[5] the maxima will be at $\varphi_{\mathbf{h}} = \varphi''_{\mathbf{h}} + |\Delta\varphi_{\mathbf{h}}|$ and $\varphi_{\mathbf{h}} = \varphi''_{\mathbf{h}} - |\Delta\varphi_{\mathbf{h}}|$. According to Eq. (3) the Sim distribution is unimodal and peaks at $\varphi_{\mathbf{h}} = \varphi''_{\mathbf{h}} - 90^\circ$. Consequently, to break the SAD phase ambiguity with either the heavy-atom substructure or the corresponding Sim distribution will always favour the peak at $\varphi_{\mathbf{h}} = \varphi''_{\mathbf{h}} - |\Delta\varphi_{\mathbf{h}}|$ in

Eq. (4). This is equivalent to forcing phases of the whole structure to approach that of the heavy-atom substructure. Thus large errors in phase estimation will be introduced, because the diffraction contribution of the heavy-atom substructure could never dominate that of the protein structure. Wang^[6] introduced in 1981 the iterative single-wavelength anomalous scattering (ISAS) method, in which both the alternative phases are used to calculate the electron density map and then the iterative solvent flattening and phase merging process is applied to break the phase ambiguity. The solvent flattening method is nowadays the most popular technique of improving the quality of electron density maps. However double-phase electron density maps could cause problems when the anomalous-scattering substructure is centrosymmetric. Fan and Gu^[7] proposed in 1985 the direct-method SAD phasing method. The main points are described below.

(i) Bimodal phase distributions from the SAD experiment (Eq. (4)) are approximated as the sum of two Gaussian functions

$$P(\varphi_{\mathbf{h}}) = \frac{1}{2\sigma_{\mathbf{h}}(2\pi)^{1/2}} \exp \left[-[\varphi_{\mathbf{h}} - (\varphi''_{\mathbf{h}} + |\Delta\varphi_{\mathbf{h}}|)]^2/2\sigma_{\mathbf{h}}^2 \right] \\ + \frac{1}{2\sigma_{\mathbf{h}}(2\pi)^{1/2}} \exp \left[-[\varphi_{\mathbf{h}} - (\varphi''_{\mathbf{h}} - |\Delta\varphi_{\mathbf{h}}|)]^2/2\sigma_{\mathbf{h}}^2 \right]. \quad (5)$$

(ii) The probability for $\Delta\varphi_{\mathbf{h}}$ being positive is given by

$$P_+ = \frac{1}{2} + \frac{1}{2} \tanh \left\{ \sin |\Delta\varphi_{\mathbf{h}}| \left[\sum_{\mathbf{h}'} m_{\mathbf{h}'} m_{\mathbf{h}-\mathbf{h}'} \kappa_{\mathbf{h},\mathbf{h}'} \sin (\Phi'_3 + \Delta\varphi_{\mathbf{h}'\text{,best}} + \Delta\varphi_{\mathbf{h}-\mathbf{h}'\text{,best}}) + \chi \sin \delta_{\mathbf{h}} \right] \right\}, \quad (6)$$

which is based on the product of the Sim distribution (Eq. (3)) and the Cochran distribution^[8] (Eq.(7))

$$P_{\text{Cochran}}(\varphi_{\mathbf{h}}) = N'' \exp \left[\sum_{\mathbf{h}'} \kappa \cos(\varphi_{\mathbf{h}} - \varphi_{\mathbf{h}'} - \varphi_{\mathbf{h}-\mathbf{h}'}) \right]. \quad (7)$$

The reader is referred to the original paper by Fan and Gu for more details.^[7]

(iii) The phase ambiguity is resolved by multiplying the first Gaussian function by P_+ and multiplying the second Gaussian function by $P_- (= 1 - P_+)$. It follows that

$$P(\varphi_{\mathbf{h}}) = \frac{P_+}{\sigma_{\mathbf{h}}(2\pi)^{1/2}} \exp \left[-[\varphi_{\mathbf{h}} - (\varphi''_{\mathbf{h}} + |\Delta\varphi_{\mathbf{h}}|)]^2/2\sigma_{\mathbf{h}}^2 \right] + \frac{1 - P_+}{\sigma_{\mathbf{h}}(2\pi)^{1/2}} \exp \left[-[\varphi_{\mathbf{h}} - (\varphi''_{\mathbf{h}} - |\Delta\varphi_{\mathbf{h}}|)]^2/2\sigma_{\mathbf{h}}^2 \right]. \quad (8)$$

Unlike the Sim distribution, the Cochran distribution may peak anywhere from 0 to 2π . Hence P_+ in Eq. (8) may enhance the first Gaussian function as well as the second. Besides, since the Cochran distribution is independent of the heavy-atom substructure, there will be no effect of whether or not the heavy-atom substructure is centrosymmetric. Finally, with the use of Cochran's distribution, the phase $\varphi_{\mathbf{h}}$ is not only directly related to the anomalous-scattering signal of the particular reflection \mathbf{h} , but also directly related to that of all available reflections in the reciprocal space. Estimations from direct-method SAD phasing are thus much more reliable. In practice, a number of typical applications^[9,10] showed that the direct-method SAD phasing is on average better than other methods, especially when the anomalous-scattering signals are weak. On the other hand, the algorithm of direct-method SAD phasing so far used in the program OASIS^[11] has its own disadvantages. The bimodal SAD phase distribution (Eq. (4)) expressed as the sum of two Gaussian functions (Eq. (5)) is not accurate. Because the experimental phase distribution is a circular function, while the Gaussian distribution

is a linear one. Errors thus introduced are not negligible, especially when the two peaks in the bimodal distribution are not sharp enough and too close to each other. This disadvantage can be effectively eliminated by replacing the two Gaussian functions with two von Mises distributions.

2. Replacing the Gaussian distribution with von Mises distribution

The von Mises distribution is also known as the circular normal distribution, which is expressed in our case as

$$P(\varphi_{\mathbf{h}}|\mu, \kappa) = \frac{\exp[\kappa \cos(\varphi_{\mathbf{h}} - \mu)]}{2\pi I_0(\kappa)}, \quad (9)$$

where μ and κ are analogous to the mean and the inverse variance of the Gaussian distribution; $I_0(\kappa)$ is a modified Bessel function of the order 0. Replacing the two Gaussian functions in Eq. (5) with the von Mises distribution expressed as Eq. (9), we have

$$P_{\text{von Mises}}(\varphi_{\mathbf{h}}) = \frac{1}{2\pi I_0(\kappa)} \exp[\kappa \cos(\varphi_{\mathbf{h}} - \varphi_{\mathbf{h}}'' - |\Delta\varphi_{\mathbf{h}}|)] + \frac{1}{2\pi I_0(\kappa)} \exp[\kappa \cos(\varphi_{\mathbf{h}} - \varphi_{\mathbf{h}}'' + |\Delta\varphi_{\mathbf{h}}|)]. \quad (10)$$

The initial value of $|\Delta\varphi_{\mathbf{h}}|$ is calculated according to Eq. (2), while the initial value of κ can be derived from the standard deviation of the "lack of closure error".^[12] A least squares refinement process is then used to refine values of $|\Delta\varphi_{\mathbf{h}}|$ and κ against the target Eq. (4). The subroutine LMDIF from the MINPACK Project^[13] is used for this purpose. 72 points with increments of 5° in the range from 0° to 360° were selected. The sum (Eq. (11)) of squares of the difference between Eqs. (10) and (4) for the 72 points is minimized by the least squares process

$$\sum_{j=0}^{72} \left[P_{\text{von Mises}}\left(\frac{5\pi \times j}{180}\right) - P_{\text{SAD}}\left(\frac{5\pi \times j}{180}\right) \right]^2. \quad (11)$$

Resultant values $|\Delta\varphi_{\mathbf{h}}|_{\text{best}}$ and κ_{best} are substituted into Eq. (10) to give

$$P_{\text{von Mises,best}}(\varphi_{\mathbf{h}}) = \frac{1}{2\pi I_0(\kappa_{\text{best}})} \exp[\kappa_{\text{best}} \cos(\varphi_{\mathbf{h}} - \varphi_{\mathbf{h}}'' - |\Delta\varphi_{\mathbf{h}}|_{\text{best}})] + \frac{1}{2\pi I_0(\kappa_{\text{best}})} \exp[\kappa_{\text{best}} \cos(\varphi_{\mathbf{h}} - \varphi_{\mathbf{h}}'' + |\Delta\varphi_{\mathbf{h}}|_{\text{best}})]. \quad (12)$$

By applying Eq. (6) to Eq. (12) we have

$$P_{\text{von Mises,best}}(\varphi_{\mathbf{h}}) = \frac{P_+}{\pi I_0(\kappa_{\text{best}})} \exp[\kappa_{\text{best}} \cos(\varphi_{\mathbf{h}} - \varphi_{\mathbf{h}}'' - |\Delta\varphi_{\mathbf{h}}|_{\text{best}})] + \frac{1 - P_+}{\pi I_0(\kappa_{\text{best}})} \exp[\kappa_{\text{best}} \cos(\varphi_{\mathbf{h}} - \varphi_{\mathbf{h}}'' + |\Delta\varphi_{\mathbf{h}}|_{\text{best}})]. \quad (13)$$

Now the best phase $\varphi_{\mathbf{h},\text{best}}$ and figure of merit $m_{\mathbf{h}}$ for a particular reflection \mathbf{h} can be calculated via the iteration of Eqs. (6), and (13)–(17), with the initial P_+ set to 1/2

$$m_{\mathbf{h}} \sin \varphi_{\mathbf{h},\text{best}} = \int_0^{2\pi} \sin \varphi_{\mathbf{h}} P_{\text{von Mises,best}}(\varphi_{\mathbf{h}}) d\varphi_{\mathbf{h}} = A, \quad (14)$$

$$m_{\mathbf{h}} \cos \varphi_{\mathbf{h},\text{best}} = \int_0^{2\pi} \cos \varphi_{\mathbf{h}} P_{\text{von Mises,best}}(\varphi_{\mathbf{h}}) d\varphi_{\mathbf{h}} = B, \quad (15)$$

$$\varphi_{\mathbf{h},\text{best}} = \arctg\left(\frac{A}{B}\right), \quad (16)$$

$$m_{\mathbf{h}} = (A^2 + B^2)^{1/2}. \quad (17)$$

Numerical integration is used to calculate Eqs. (14) and (15).

3. Test data

SAD data from two known proteins were used as examples in testing the new algorithm. The sample

proteins are summarized in Table 1. Both examples are difficult SAD phasing cases. The sample azurin^[14] has a low Bijvoet ratio ($\langle|\Delta F|\rangle/\langle F\rangle$) 1.45%. Furthermore the overall data completeness is only 60%. The sample xylanase^[15] belongs to one of the most difficult cases. The sample's Bijvoet ratio 0.56% is one of the lowest Bijvoet ratios, which have so far been successfully treated by SAD phasing. Besides, the low solvent content 37% is not beneficial to the subsequent phase improvement by density modification techniques.

Table 1. Test samples.

protein	number of residues /per AU	high resolution limit/Å	space group	x-ray wavelength /Å	number of anomalous scatterers/per AU	expected Bijvoet ratio/%	reference
azurin	129	1.9	P4 ₁ 22	0.97	1 (copper)	1.45	[14]
xylanase	303	1.8	P2 ₁	1.49	5 (sulfur)	0.56	[15]

4. Test and results

4.1. Comparison of accuracy between different approximations

The accuracy of the von Mises approximation (Eq. (12)) was compared with that of the Gaussian approximation (Eq. (5)). The sample azurin was used in this test. For the majority of reflections, both approximations are good, but the von Mises approximation is better. There remained an indispensable portion of reflections, for which the Gaussian approximation led to large errors while the von Mises approximation still gave pretty good results. Phase probability distributions of four reflections belonging to this portion are shown in Fig. 1. As is seen, von Mises approximations are all reasonably well fitted into the experimental curves and are much better than Gaussian approximations. The difference in accuracy between the two approximations would have definite effects on the SAD phasing results, which are shown in the next section.

4.2. Comparison between results of iterative SAD phasing based on different approximations

Both samples azurin and xylanase were used in this test. They were subject to the following steps in each iteration cycle during the test: i) SAD phasing by OASIS;^[11] ii) density modification by DM;^[16] iii) model building by RESOLVE^[17,18] for the first two cycles; iv) model building by ARP/wARP^[19] from the third cycle onward. Results are listed in Table 2. As is seen, the von Mises approximation significantly speeds up the convergence of iteration. For the sample azurin, in 3 cycles of iteration based on the von Mises approximation the size of structure model grew up from 50 to 108 of the total 129 residues. 99 of the 108 residues have been assigned into the sequence. Meanwhile the phase error decreased from 71.0° to 36.4°. In contrast, in 3 cycles of

iteration based on the Gaussian approximation the model size grew up from 54 to only 59 residues, while the phase error decreased just from 73.8° to 63.4° . For the sample xylanase, in 4 cycles of iteration based on the von Mises approximation, the model size grew up from 0 to 298 of the total 303 residues. All the 298 residues have been assigned into the sequence. Meanwhile the phase error decreased from 80.0° to 26.9° . In contrast, in 4 cycles of iteration based on the Gaussian approximation the model size changed from 51 to 35 residues, none of which has been assigned into the sequence. The phase error decreased just from 77.0° to 52.0° .

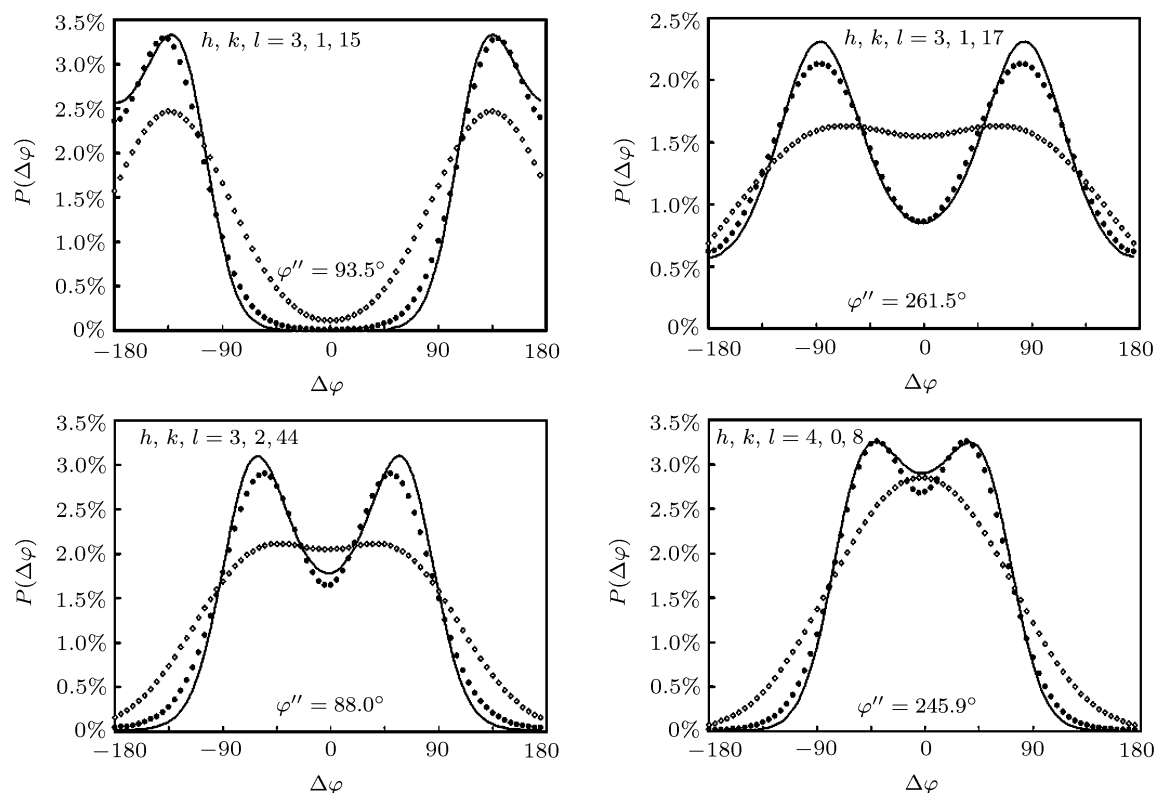


Fig. 1. Experimental bimodal phase distributions approximated by the sum of two unimodal distributions. Solid line is the experimental bimodal phase distribution calculated from Eq. (4). The symbol \bullet indicates the approximation by the sum of two von Mises distributions calculated from Eq. (12). The symbol \diamond is the approximation by the sum of two Gaussian distributions calculated from Eq. (5).

Table 2. Test results of iterative SAD phasing.

cycle	azurin						xylanase					
	von Mises approximation			Gaussian approximation			von Mises approximation			Gaussian approximation		
	number of residues built in the model	number of residues assigned into the sequence	phase error ($^\circ$) of the model	number of residues built in the model	number of residues assigned into the sequence	phase error ($^\circ$) of the model	number of residues built in the model	number of residues assigned into the sequence	phase error ($^\circ$) of the model	number of residues built in the model	number of residues assigned into the sequence	phase error ($^\circ$) of the model
1	50	0	71.0	54	0	73.8	42	0	80.0	51	0	77.0
2	74	0	64.6	53	0	72.3	98	0	74.6	95	0	76.4
3	108	99	36.4	59	26	63.4	23	0	53.4	5	0	57.4
4							298	298	26.9	35	0	52.0

5. Conclusion

The above comparison shows that great improvement on the efficiency of SAD phasing can be achieved by replacing the Gaussian distribution with the von Mises distribution to approximate the bimodal SAD phase distribution. The new algorithm described in this paper will be incorporated into the next version of the program OASIS.^[11]

References

- [1] Ramachandran J N and Raman S 1956 *Curr. Sci.* **25** 346
- [2] Hendrickson W A and Teeter M M 1981 *Nature* **290** 107
- [3] Sim G A 1959 *Acta Cryst.* **12** 813
- [4] Ten Eyck L F and Arnone A 1976 *J. Mol. Biol.* **100** 3
- [5] Blundell T L and Johnson L N 1976 *Protein Crystallography* (London: Academic Press Inc.) p. 177
- [6] Wang B C 1985 *Methods in Enzymology* **115** 90
- [7] Fan H F and Gu Y X 1985 *Acta Cryst. A* **41** 280
- [8] Cochran W 1955 *Acta Cryst.* **8** 473
- [9] Watanabe N, Kitago Y, Tanaka I, Wang J W, Gu Y X, Zheng C D and Fan H F 2005 *Acta Cryst. D* **61** 1533
- [10] Yao D Q, Li H, Chen Q, Gu Y X, Zheng C D, Lin Z J, Fan H F, Watanabe N and Sha B D 2008 *Chin. Phys. B* **17** 1
- [11] Zhang T, Wu L J, He Y, Wang J W, Zheng C D, Hao Q, Gu Y X and Fan H F 2009 *OASIS4.0—A direct-methods program for SAD/SIR phasing and reciprocal-space fragmentation extension* Institute of Physics, Chinese Academy of Sciences, Beijing 100190, China (program available at <http://cryst.iphy.ac.cn>)
- [12] Blow D M and Crick F H C 1959 *Acta Cryst.* **12** 794
- [13] Moré J J, Sorensen D C, Hillstrom K E and Garbow B S 1984 *The MINPACK Project in Sources and Development of Mathematical Software* ed. Cowell W J (Upper Saddle River: Prentice-Hall) pp. 88–111.
- [14] Dodd F, Hasnain S S, Abraham Z H, Eady R R and Smith B E 1995 *Acta Cryst. D* **51** 1052
- [15] Ramagopal U A, Dauter M and Dauter Z 2003 *Acta Cryst. D* **59** 1020
- [16] Cowtan K D and Main P 1993 *Acta Cryst. D* **49** 148
- [17] Terwilliger T C 2003 *Acta Cryst. D* **59** 38
- [18] Terwilliger T C 2003 *Acta Cryst. D* **59** 45
- [19] Perrakis A, Morris R and Lamzin V S 1999 *Nature Struct. Biol.* **6** 458