

Because of these differences, the conditional probability distributions which will be of most use may vary from one case to the other. Similarly, the phasing procedures which will be most effective may vary from one case to the other.

In this paper we have focused on the similarities among these cases, for the following reasons. In a direct-methods probabilistic approach, the derivation of a joint probability distribution is often a lengthy initial task. As shown here, this analysis need only be done once, if it is formulated in a general way. It is also an easy task to translate a distribution derived for a specific case into more general terms. Consequently, much of the already available theoretical foundation for either SIR, SAS or partial/complete structure data may be reformulated so that it can be used in any example of isomorphous data sets. For example, the joint probability distribution of a triplet of isomorphous data sets (Fortier, Weeks & Hauptman, 1984*b*) can be translated easily into any case of interest, such as the case of a native protein and a single heavy-atom derivative for which Friedel-pair data are available. Thus, as in the algebraic approach presented by Karle (1984, 1985), general formulae can be used on a large variety of combinations of the various cases. Finally, while there is still little experience in the use of direct methods in macromolecular structure determination, much valuable experience, both practical and theoretical, has been gained in the use of direct methods for *ab initio* phasing of isomorphous data sets in small molecules. In particular, much can be learned from the vast

amount of expertise that has been gained in the applications of direct methods to the problem of partial structure expansion in the *DIRDIF* system (Beurskens *et al.*, 1981).

Financial assistance from the Natural Sciences and Engineering Research Council of Canada is gratefully acknowledged.

References

- BEURSKENS, P. T., BOSMAN, W. P., DOESBURG, H. M., GOULD, R. O., VAN DEN HARK, TH. E. M., PRICK, P. A. J., NOORDIK, J. H., BEURSKENS, G. & PARTHASARATHY, V. (1981). Tech. Rep. 1981/2. Crystallography Laboratory, Toernooiveld, 6525ED Nijmegen, The Netherlands.
- BEURSKENS, P. T., PRICK, P. A. J., DOESBURG, H. M. & GOULD, R. O. (1979). *Acta Cryst.* **A35**, 765-772.
- FAN HAI-FU, HAN FU-SON, QIAN JIN-ZI & YAO JIA-XING (1984). *Acta Cryst.* **A40**, 489-495.
- FORTIER, S., WEEKS, C. M. & HAUPTMAN, H. (1984*a*). *Acta Cryst.* **A40**, 544-548.
- FORTIER, S., WEEKS, C. M. & HAUPTMAN, H. (1984*b*). *Acta Cryst.* **A40**, 646-651.
- GIACOVAZZO, C. (1983*a*). *Acta Cryst.* **A39**, 585-592.
- GIACOVAZZO, C. (1983*b*). *Acta Cryst.* **A39**, 685-692.
- GIACOVAZZO, C., CASCARANO, G. & ZHENG CHAO-DE (1988). *Acta Cryst.* **A44**, 45-51.
- GLUSKER, J. P. & TRUEBLOOD, K. N. (1985). *Crystal Structure Analysis: a Primer*. Oxford Univ. Press.
- HAUPTMAN, H. (1982*a*). *Acta Cryst.* **A38**, 289-294.
- HAUPTMAN, H. (1982*b*). *Acta Cryst.* **A38**, 632-641.
- KARLE, J. (1984). *Acta Cryst.* **A40**, 374-379.
- KARLE, J. (1985). *Acta Cryst.* **A41**, 182-189.
- KARLE, J. & HAUPTMAN, H. (1958). *Acta Cryst.* **11**, 264-269.
- SIM, G. A. (1959). *Acta Cryst.* **12**, 813-815.
- SRINIVASAN, R. & PARTHASARATHY, S. (1976). *Some Statistical Applications in X-ray Crystallography*. New York: Pergamon.

Acta Cryst. (1989). **A45**, 254-258

The Probability Distributions of X-ray Intensities in Fiber Diffraction: Largest Likely Values for Fiber Diffraction *R* Factors

BY GERALD STUBBS

Department of Molecular Biology, Vanderbilt University, Nashville, TN 37235, USA

(Received 11 March 1988; accepted 14 September 1988)

Abstract

R factors in fiber diffraction are generally lower than in conventional crystallography, because of the cylindrical averaging of fiber diffraction data. The probability distributions for fiber diffraction intensities, analogous to Wilson's distributions for crystal diffraction intensities, are derived, and from these the largest likely values of *R* are estimated. These values depend on the size and symmetry of the diffracting particle and on the resolution of the analysis, and range from 0.586 for systems for very high symmetry (as in crystal

diffraction) to much lower values for systems of low symmetry.

Introduction

The *R* factor, $R = \sum ||F_{\text{obs}}| - |F_{\text{calc}}|| / \sum |F_{\text{obs}}|$, has been used for many years as an index of the quality of crystallographic structure determinations. It is also widely quoted in descriptions of structures determined by refinement of models against fiber diffraction data, although in fiber diffraction $|F|$ must be replaced by $I^{1/2}$. (In fiber diffraction, $|F|$ is not gen-

erally equal to $I^{1/2}$ because of the cylindrical averaging of the data.) Wilson (1949) derived the probability distribution of X-ray intensities in conventional crystallography. For crystals without a center of symmetry, he showed that the probability that I will lie between I and $I + dI$ is

$$P(I) dI = \sum^{-1} \exp(-I/\Sigma) dI$$

or, equivalently

$$P(F) dF = 2F \sum^{-1} \exp(-F^2/\Sigma) dF$$

where $\Sigma = \sum f_j^2$, and f_j is the scattering factor of the j th atom in the unit cell. He went on to show (Wilson, 1950) that the largest likely value for the R factor from an acentric crystal is $2 - 2^{1/2} = 0.586$.

In fiber diffraction, the R factor for a structure in which the atomic coordinates are uncorrelated with the true atomic coordinates is expected to be lower than in crystallography. This is because in a fiber specimen the diffracting particles are randomly oriented about the fiber axis, so the diffraction pattern is cylindrically averaged. The diffracted intensity at reciprocal-space radius R on layer line l is

$$I(R, l) = \sum_n \mathbf{G}_{n,l}(R) \mathbf{G}_{n,l}^*(R) \quad (1)$$

(Waser, 1955; Franklin & Klug, 1955), where n is the order of the Bessel functions J_n that contribute to the complex Fourier-Bessel structure factor \mathbf{G} (Klug, Crick & Wyckoff, 1958). For a helical structure, n is restricted by the selection rule $l = tn + um$, where m is an integer and there are u subunits in t turns of the helix. The number of significant \mathbf{G} terms in (1), N , depends on the symmetry and dimensions of the diffracting particle, and on the value of (R, l) . For example, for tobacco mosaic virus (TMV) at 2.9 Å resolution, N can be as large as 8. Equation (1) is the analog of the crystallographic equation

$$I(h, k, l) = \mathbf{F}_{hkl} \mathbf{F}_{hkl}^*$$

It is useful to define a $2N$ -dimensional vector \mathcal{G} , whose components are the real and imaginary components of the \mathbf{G} terms contributing to a particular intensity $I(R, l)$. From (1), the norm of \mathcal{G} , \mathcal{G} , is equal to $I^{1/2}$. Just as the largest likely R factor for an acentric crystal, for which the structure factor \mathbf{F} is two-dimensional, is smaller than for a centric crystal, where \mathbf{F} is one-dimensional (Wilson, 1950), the R factor for a fiber, with \mathcal{G} multidimensional, is even lower. Qualitatively, it is easier to predict the sum of several data, such as \mathcal{G} , than to fit individual data points, such as F .

Recent developments in fiber diffraction analysis have created a need for a quantitative understanding of fiber diffraction R factors. TMV has been refined to an R factor of 0.096 at 2.9 Å resolution (Namba, Pattanayek & Stubbs, 1989), and the filamentous bacteriophage Pf1 to 0.24 at 4 Å (Nambudripad &

Makowski, personal communication). In the future, structure determinations of a wide variety of filamentous viruses, cytoskeletal filaments and other macromolecular assemblies will depend on the further development of all the forms of analysis that are now used routinely in protein crystallography, including refinement and the understanding of indicators of the progress of refinement such as R factors. In contrast to conventional crystallography, fiber diffraction R factors for systems of different size and symmetry are not directly comparable to each other. One solution to this problem is to determine the R factor expected for a totally wrong determination of a particular structure, and to compare this with an experimental R factor in order to assess the reliability of a model.

In this paper, I will attempt to quantify these general statements, and to predict the values of the largest likely R factors in fiber diffraction. It will be shown that these values depend on the number of terms in (1), and therefore on the diameter and symmetry of the diffracting particles, and on the resolution of the data. The probability distributions of X-ray intensities in fiber diffraction will be derived, following the approach used by Wilson (1949), and used to calculate the largest likely values of R .

Theory

Definitions and preliminary results

Following Wilson (1950), we define

$$H(F) = \int_0^F F P(F) dF \quad (2)$$

and note that $H(\infty) = \langle |F| \rangle$. In order to avoid confusion with standard fiber diffraction notation, we use the symbol H in place of Wilson's G . Wilson showed that, for a random distribution of atoms,

$$R = 2 - 4\langle H(F) \rangle / \langle |F| \rangle. \quad (3)$$

The derivation is valid for all distributions of F , and is therefore applicable to fiber diffraction, using \mathcal{G} in place of F . In order to use (3) for fiber diffraction data, $P(\mathcal{G})$ and thence $H(\mathcal{G})$ must be derived.

To derive $P(\mathcal{G})$ in the next section, we will need to know the integral V_M over all points \mathbf{r} in M -dimensional space such that

$$r^2 = \sum_{i=1}^M x_i^2.$$

This corresponds to the circumference of a circle in two dimensions, and the surface area of a sphere in three. If we assume that the integral V_{M-1} over all points \mathbf{d} such that

$$d^2 = \sum_{i=1}^M x_i^2$$

is known, then V_M is the integral over x_M of all points such that $d^2 + x_M^2 = r^2$, weighted by V_{M-1} . Setting $d = r \sin \theta$, we see that this is the weighted integral along half the circumference of a circle of radius r , that is, $V_M = \int_0^\pi V_{M-1} r d\theta$. It may be shown by induction that $V_n = v_n d^{n-1}$, where v_n is a constant.

$$\begin{aligned} \therefore V_M &= v_{M-1} r \int_0^\pi d^{M-2} d\theta \\ &= v_{M-1} r^{M-1} \int_0^\pi \sin^{M-2} \theta d\theta. \end{aligned} \quad (4)$$

The integral in (4) can be evaluated using Wallis's formula [Abramowitz & Stegun (1972); equation (6.1.49)]. Table 1 contains expressions for V_M for $M = 2N$ up to $N = 8$.

The probability distributions of fiber diffraction intensities

The probability distribution of \mathcal{G} may be derived from the contribution of each atom j to the components A_i of \mathcal{G} (where each A_i is the real or the imaginary part of a \mathbf{G}). With no loss of generality, we may assume that A_i is real. Then

$$A_i = \sum_j f_j J_n(2\pi R r_j) \cos[-n\varphi_j + 2\pi l(z/c)j]$$

where r_j , φ_j and z_j are the cylindrical polar coordinates of atom j , f_j is its scattering factor, c is the axial repeat of the diffracting particle and J_n is the Bessel function of order n .

Setting the cosine argument to θ_j , and assuming that the atoms are sufficiently randomly distributed for cross terms to cancel, we find

$$A_i^2 = \sum_j f_j^2 J_n^2(2\pi R r_j) \cos^2(\theta_j).$$

For random values of φ_j and z_j , and therefore of θ_j ,

$$\begin{aligned} \langle A_i^2 \rangle &= \frac{1}{2} \sum_j f_j^2 J_n^2(2\pi R r_j) \\ &= \frac{1}{2} \sum_j \text{where } \Sigma = \sum_j f_j^2 J_n^2(2\pi R r_j). \end{aligned}$$

This derivation of A_i differs from that of Wilson (1949) in one vital respect: whereas Wilson averaged the cosine terms over reflections, here they are averaged over atomic coordinates. This is necessary because in fiber diffraction θ_j is not randomly distributed in reciprocal space unless the number of \mathbf{G} terms in the whole diffraction pattern is very large; for a given \mathbf{G} , θ_j is constant. This difference may impose limits on the validity of the derivation for structures with very few atoms, but it will not affect the analysis for fiber diffraction from macromolecules.

In principle, Σ depends on the values of r_j , but in practice it may be assumed that the contributing atoms are evenly distributed within the known radial limits of the diffracting particle. A more precise esti-

mate of Σ may be obtained if the radial density distribution is known and the atomic structure factors can be assumed to be equal; the radial density distribution can usually be calculated from the equatorial data, either with the aid of a heavy-atom derivative (Caspar, 1956; Franklin, 1956), or by the minimum wavelength principle (Bragg & Perutz, 1952; Finch, 1965). Many useful calculations, however, including that of the largest likely values of R , can be made without knowing the value of Σ .

Assuming, by the central limit theorem, that A_i is normally distributed, we find

$$P(A_i) dA_i = (\pi \Sigma)^{-1/2} \exp(-A_i^2/\Sigma) dA_i. \quad (5)$$

This is equation (12) of Wilson (1949). Combination of the $2N$ orthogonal parts of \mathcal{G} gives the joint probability distribution

$$P(\mathcal{G}) d\mathcal{G} = (\pi \Sigma)^{-N} \exp(-\mathcal{G}^2/\Sigma) d\mathcal{G}. \quad (6)$$

The step from (5) to (6) assumes that Σ is the same for all Bessel orders n , that is, that the average value of $J_n(2\pi R r_j)$ does not depend on n . This is generally true, but there are small deviations near the first maximum of J_n . The derivations will therefore be somewhat less accurate for diffraction patterns dominated by first maxima. Although, for simplicity, many diffraction patterns have in the past been analysed as if they were so dominated, first maxima do not in fact make up a significant part of high-resolution fiber diffraction patterns.

In order to describe observed intensity distributions and calculate $H(F)$ from (2), $P(\mathcal{G}) d\mathcal{G}$, the probability that $|\mathcal{G}|$ lies between \mathcal{G} and $\mathcal{G} + d\mathcal{G}$, must be derived from $P(\mathcal{G})$. Now

$$P(\mathcal{G}) d\mathcal{G} = \int P(\mathcal{G}) d\mathcal{G}$$

where the integral is over all points between \mathcal{G} and $\mathcal{G} + d\mathcal{G}$. Hence

$$P(\mathcal{G}) d\mathcal{G} = V_M P(\mathcal{G}) d\mathcal{G}, \quad (7)$$

where V_M (4) is the area of the $(M-1)$ -dimensional surface formed by all points at a distance \mathcal{G} from the origin (Table 1). From (6) and (7),

$$P(\mathcal{G}) d\mathcal{G} = v_M (\pi \Sigma)^{-N} \mathcal{G}^{M-1} \exp(-\mathcal{G}^2/\Sigma) d\mathcal{G}. \quad (8)$$

This expression describes the distribution of amplitudes in a fiber diffraction pattern; the corresponding expression for intensities is

$$P(I) dI = \frac{1}{2} v_M (\pi \Sigma)^{-N} I^{N-1} \exp(-I/\Sigma) dI.$$

Largest likely R factors

In order to determine the largest likely value of R from (3), H and $\langle H \rangle$ must be determined. From (2) and (8),

$$H(\mathcal{G}) = v_M (\pi \Sigma)^{-N} \int_0^{\mathcal{G}} \mathcal{G}^M \exp(-\mathcal{G}^2/\Sigma) d\mathcal{G}$$

and

$$\begin{aligned} \langle H(\mathcal{G}) \rangle &= \int_0^\infty H(\mathcal{G}) P(\mathcal{G}) d\mathcal{G} \\ &= v_M^2 (\pi \Sigma)^{-M} \int_0^\infty \mathcal{G}^{M-1} \exp(-\mathcal{G}^2/\Sigma) \\ &\quad \times \int_0^{\mathcal{G}} \mathcal{G}^M \exp(-\mathcal{G}^2/\Sigma) d\mathcal{G} d\mathcal{G}. \end{aligned}$$

For small values of N , these integrals may be evaluated by parts (Wilson, 1950), but in general it is straightforward and convenient to evaluate them numerically. The values of R determined from (3) do not depend on Σ . Table 1 includes these values for M between 2 and 16, covering the range encountered in most fiber diffraction analyses currently in progress.

Applications

The R factors in Table 1 allow us to estimate the value of R to be expected for a completely wrong structure in a fiber diffraction analysis. In practice, the number of overlapping terms in (1), and therefore the value of M , varies across the diffraction pattern, so a weighted average of the R values in the table must be used. A reasonable procedure is to assume that a Bessel-function term of order n becomes significant when the argument $2\pi rR$ is equal to $n-2$, with r equal to the maximum radius of the diffracting particle. For TMV, with a maximum radius of 90 Å and 49 subunits in three turns of the viral helix, this weighted average is 0.40 between 10 and 5 Å resolution, and 0.34 between 10 and 3 Å. For Pf1, with a maximum radius of 30 Å and 27 subunits in five turns (Makowski, 1984), the corresponding figures are 0.48 and 0.41.

Helical assemblies with smaller repeating units generally have smaller maximum radii, but they often have lower symmetry than larger assemblies. These two properties have opposite effects on the number of terms contributing to the diffracted intensity in (1). As a hypothetical example, one might take a non-crystalline fiber in which the maximum radius was 10 Å and the asymmetric unit repeated ten times in one turn of the helix. In such a case, the largest likely R factor between 10 and 3 Å resolution would be 0.41.

Although primarily intended for continuous diffraction from non-crystalline fiber specimens, the approach described here may equally well be applied to crystalline fiber data in which there is a significant number of overlapping reflections. For example, the diffracted data from chondroitin 4-sulfate (Winter, Arnott, Isaac & Atkins, 1978) include five intensities derived from three overlapping reflections ($M = 5$ or

Table 1. *The integral, V_M , over all points \mathbf{x} in M -dimensional space such that $r^2 = \sum_{i=1}^M x_i^2$ and the largest likely R factors for N overlapping G terms, with a total of M components*

For M even, $M = 2N$. Odd values of M occur when one of the G terms has only a real component, for example, on the equator; in this case, $M = 2N - 1$.

M	V_M	R
2 (circle)	$2\pi r$	0.586
3 (sphere)	$4\pi r^2$	0.475
4	$2\pi^2 r^3$	0.409
5	$\frac{8}{3}\pi^2 r^4$	0.364
6	$\pi^3 r^5$	0.332
7	$\frac{16}{15}\pi^3 r^6$	0.306
8	$\frac{1}{3}\pi^4 r^7$	0.286
9	$\frac{32}{105}\pi^4 r^8$	0.269
10	$\frac{1}{12}\pi^5 r^9$	0.255
11	$\frac{64}{945}\pi^5 r^{10}$	0.243
12	$\frac{1}{60}\pi^6 r^{11}$	0.232
13	$\frac{128}{10395}\pi^6 r^{12}$	0.223
14	$\frac{1}{360}\pi^7 r^{13}$	0.215
15	$\frac{256}{135135}\pi^7 r^{14}$	0.207
16	$\frac{1}{2520}\pi^8 r^{15}$	0.201

6, depending on whether one of the reflections is from a centric zone), 37 from two reflections ($M = 3$ or 4), and 11 from single reflections ($M = 2$). In this case, the weighted averaged largest likely R factor would be 0.47.

The results described here should allow more objectivity in comparing results of fiber diffraction analyses and assessing the value of atomic parameter refinements against fiber diffraction data. It is clear that the R factors to be expected in fiber diffraction refinement are considerably lower than those in conventional crystallography; nonetheless, even for systems of relatively low symmetry (and therefore large numbers of overlapping terms), the largest likely R factors are significantly higher than typical levels of error in the data. Refinement of molecular models against fiber diffraction data is therefore practical, but must be evaluated by comparing the R factors obtained from the model with the appropriate values derived here.

I thank Rick Millane and Lee Makowski for helpful comments on the manuscript of this paper. This work

was supported by NSF grant BBS8717949 and NIH grant GM33265.

References

- ABRAMOWITZ, M. & STEGUN, I. A. (1972). *Handbook of Mathematical Functions*. New York: Dover.
- BRAGG, W. L. & PERUTZ, M. F. (1952). *Proc. R. Soc. London Ser. A*, **213**, 425-435.
- CASPAR, D. L. D. (1956). *Nature (London)*, **177**, 928.
- FINCH, J. T. (1965). *J. Mol. Biol.* **12**, 612-619.
- FRANKLIN, R. E. (1956). *Nature (London)*, **177**, 938-930.
- FRANKLIN, R. E. & KLUG, A. (1955). *Acta Cryst.* **8**, 777-780.
- KLUG, A., CRICK, F. H. C. & WYCKOFF, H. W. (1958). *Acta Cryst.* **11**, 199-213.
- MAKOWSKI, L. (1984). In *Biological Macromolecules and Assemblies*. Vol. 1. *Virus Structures*, edited by F. A. JURNAK & A. MCPHERSON, pp. 203-253. New York: Wiley.
- NAMBA, K., PATTANAYEK, R. & STUBBS, G. (1989). *J. Mol. Biol.* In the press.
- WASER, J. (1955). *Acta Cryst.* **8**, 142-150.
- WILSON, A. J. C. (1949). *Acta Cryst.* **2**, 318-321.
- WILSON, A. J. C. (1950). *Acta Cryst.* **3**, 397-399.
- WINTER, W. T., ARNOTT, S., ISAAC, D. H. & ATKINS, E. D. T. (1978). *J. Mol. Biol.* **125**, 1-19.

Acta Cryst. (1989). **A45**, 258-260

R Factors in X-ray Fiber Diffraction. I. Largest Likely *R* Factors for *N* Overlapping Terms

BY R. P. MILLANE

*The Whistler Center for Carbohydrate Research, Smith Hall, Purdue University,
West Lafayette, Indiana 47907, USA*

(Received 2 June 1988; accepted 19 September 1988)

Abstract

Simple expressions are obtained for the largest likely *R* factor in X-ray fiber diffraction recently derived by Stubbs [*Acta Cryst.* (1989), **A45**, 254-258]. These generalize the largest likely *R* factors obtained by Wilson [*Acta Cryst.* (1950), **3**, 397-399] for centric and acentric crystals. Expressions are obtained in terms of special functions and as finite series that simplify the calculation of *R* factors. These may be useful for further analysis and understanding of the effects of particle diameter and symmetry and diffraction data resolution on the reliability of structure determinations.

1. Introduction

The *R* factor is used routinely in crystallography to measure the reliability of structure determinations. Interpretation of the *R* factor obtained for a particular structure determination is aided by comparing it with the value for a completely wrong structure, *i.e.* a structure that is uncorrelated with the correct structure. This is referred to as the 'largest likely *R* factor', and Wilson (1950) showed that its value is $2\sqrt{2} - 2 = 0.828$ for a centric crystal and $2 - \sqrt{2} = 0.586$ for an acentric crystal.

Recent advances in data collection and structure refinement in X-ray (and neutron) fiber diffraction analysis (Millane, 1988) have led to determinations of the structures of complex fibrous molecules and assemblies (Millane, Walker, Arnott, Chandrasekaran & Ratliff, 1984; Namba & Stubbs, 1985; Park, Arnott, Chandrasekaran, Millane & Campagnari, 1987;

Stark, Glucksman & Makowski, 1988), and the *R* factor is used as a measure of the reliability of these structures also. The molecules in a fiber specimen are randomly oriented about the fiber axis so that the diffraction pattern is cylindrically averaged. The measured intensity is therefore equal to the sum of a number of different intensity terms diffracted by a single molecule. The number of terms in the sum depends on the maximum diameter and symmetry of the molecule, and the position in reciprocal space at which the intensity is measured. Since the measured intensities are sums of individual structure intensities, the *R* factor is in general smaller than in conventional crystallography. Stubbs (1989) has recently determined the largest likely *R* factor in fiber diffraction analysis as a function of the number of overlapping intensity terms. This allows the maximum value of the *R* factor for a particular structure determination to be estimated by averaging the values over the recorded diffraction pattern where the number of overlapping terms varies. This can be applied to both continuous diffraction from non-crystalline specimens and Bragg diffraction from polycrystalline specimens. The values obtained allow the *R* factor to be used for an objective assessment of the quality of structures determined by fiber diffraction.

Here, simple analytical and algebraic forms of Stubbs's (1989) expression for the largest likely fiber diffraction *R* factor are derived. These may be useful for further theoretical analysis of the dependence of the *R* factor on the number of overlapping intensity terms, the particle size and symmetry, and resolution of the diffraction data.