

Crystal Structure Analysis from a Viewpoint of Information Theory

BY SUKEAKI HOSOYA AND MASAYASU TOKONAMI

Institute for Solid State Physics, University of Tokyo, Azabu, Minato-ku, Tokyo, Japan

(Received 21 October 1966)

Information amount is defined for values of reflexion intensities and Patterson peaks. How the information recovery is made during the procedures of structure analysis is exemplified by a real structure 96R-SiC. The characteristic differences between the direct method, especially the statistical method, and the Patterson method are shown. In the usual statistical method the random distribution of atomic positions is assumed. However, limitations such as steric hindrance or molecular forms affect the intensity distribution. Influence due to this complexity is also included to some extent in the present example. From this standpoint, comments are made on several theoretical works on structure analysis.

Introduction

The information theory was first developed in communication engineering. Although the theory was later successfully applied to light optics, it has so far been used only in a few works in the field of X-ray structure analysis of crystals. The sampling theorem which is fundamental in the information theory was applied by Sayre (1952) to the phase problem. The familiar unit 'bit' devised in the information theory was used by Diamond (1963) for expressing the information amount included in each inequality relation among structure factors. The information amount included in an absolute value of a structure factor was discussed by one of the present authors in his review article on the phase problem (Hosoya, 1964). In the present paper, the notion of the information amount is defined in a more general form, and then applied to the information amount included in peaks of a Patterson function or of a vector set. Some numerical analyses are described for an existing model.

The general idea will be presented using the unitary structure factor U . As is well known, the crystal structure factor $F(\mathbf{h})$ is expressed by the atomic scattering factor f_j and the atomic coordinate \mathbf{r}_j of a j th atom as follows:

$$F(\mathbf{h}) = \sum_j f_j \exp(2\pi i \mathbf{h} \cdot \mathbf{r}_j).$$

If each f_j is proportional to the average atomic scattering factor \hat{f} as

$$f_j = a_j \hat{f},$$

all discussion on $U(\mathbf{h})$ is also valid for $F(\mathbf{h})$, where a_j is the atomic number of the j th atom. It is to be noted that the values of U corresponding to the reciprocal points located too far from the origin of reciprocal space cannot be observed.

Now, two kinds of space, \mathbf{x}_j space and \mathbf{I}^h space, are introduced as follows: \mathbf{x}_j space is defined by a set of J atomic coordinates (x_1, x_2, \dots, x_j) and \mathbf{I}^h space is defined by a set of H observed values $(I^{h1}, I^{h2}, \dots, I^{hH})$,

where $I^h = |U(\mathbf{h})|^2$ will be called the unitary intensity. The crystal structure analysis is eventually the procedure of finding out a point in \mathbf{x}_j space from the point in \mathbf{I}^h space given by the observations.

In the present work, consideration is given as to how the information obtained from the observation of I^h reduced the super-volume, including the solution point, in \mathbf{x}_j space, and how the amount of information can be geometrically defined in a multi-dimensional space. A real structure, 96R-SiC, is discussed for which the dimension J of \mathbf{x}_j space is 32, and the dimension H of \mathbf{I}^h space happens to be also 32.

Definition of information amount

At every stage when some information is given, a super-volume in \mathbf{x}_j space containing points which are possible solutions should shrink to a more limited volume. Then the information amount obtained at each stage may be expressed as $-\log W$ (Shannon, 1948), where W is the ratio of these two super-volumes before and after the shrinkage. When 2 is taken for the base of logarithms, the value of $-\log W$ is expressed in 'bits'.

The information amount can be defined for continuous as well as discontinuous variables. In practical calculations, however, each variable can be dealt with as a quantity quantized in discrete levels. For instance, atomic fractional coordinates are usually expressed by 2–3 decimal digits at an early stage of the analysis and then by 4–5 digits in a refining stage. In other words, the \mathbf{x}_j space can be considered as a set of sampling points forming a lattice. As to the diffraction intensity, the value I^h is also best taken as a quantity quantized in finite degrees of magnitude, especially because the measured values are always more or less affected by errors.

Suppose that the super-volume including the solution point is scanned in \mathbf{x}_j space. Corresponding to the change of I^h due to such a scanning, let $p(I_i^h)$ be the probability that the value I^h has the intensity of the i th degree. Once an \mathbf{h} reflexion has been known to have the value I_i^h , it gives the information amount $-\log$

$[p(I_i^h)]$, and therefore the expected value for the information amount from the h reflexion should be

$$H(h) = - \sum_i p(I_i^h) \log[p(I_i^h)]. \quad (1)$$

Expressions of this type will be used throughout the present work.

To be more general, the information amount given by two reflexions is

$$H(h1, h2) = - \sum_{ij} p(I_i^{h1}, I_j^{h2}) \log[p(I_i^{h1}, I_j^{h2})], \quad (2)$$

using a joint probability which can be defined analogously to $p(I_i^h)$. Those formula relevant to more than two reflexions are also given in a similar way. The redundancy in the two reflexions is thus defined as

$$R(h1, h2) = H(h1) + H(h2) - H(h1, h2). \quad (3)$$

The redundancy serves as a measure of the information recovery.

In the same way, the information amount given by a value at a point r of a vector set or a Patterson function can be defined as

$$H(r) = - \sum_i p(V_i^r) \log[p(V_i^r)], \quad (4)$$

where V_i^r is the value integrated over a small region specified by the position r in the Patterson space, and a suffix i is a degree of magnitude of this peak value.

Example with 96R-SiC

(1) The estimated number of possible structures

As is well known, the structure of SiC consists of a stacking of ABC layers, being essentially a one-dimensional structure. The structure factor of 96R-SiC can be expressed as

$$F(hkl) = f(hkl) \cdot U(l),$$

where $f(hkl)$ is the structure factor for a chemical unit of SiC and $U(l)$ is the unitary structure factor for the one-dimensional crystal:

$$U(l) = \sum_{j=1}^{32} \exp(2\pi i l z_j / 96), \quad (z_j = \text{integer}).$$

Because the approximate value of $|f(hkl)|$ can be easily calculated, $|U(l)|^2$ can be obtained from the observed $|F(hkl)|^2$. It is to be noted here that the maximum value of $|U(l)|$ in this example is normalized, for convenience, to be 32 instead of 1.

The x_j space for the 96R-SiC structure is, therefore, of 32 dimensions, in which only points with coordinates of multiples of $1/96$ need to be taken into account.

If there is no limitation at all for the coordinates z_j , $32^{96} = 2^{480} \approx 10^{144}$ kinds of structure may be possible. Actually, there are of course several limitations as follows.

(i) More than one layer cannot occupy an identical position. This limitation is expressed as

$$z_j - z_{j'} \neq 0 \pmod{96},$$

which makes possible structures decrease in number down to

$$96! / (32! 64!) \approx 2^{84.6} \approx 10^{25}.$$

(ii) Because of steric hindrance, a layer A cannot follow A , and this is the case with B and C , respectively. This limitation can be expressed as

$$z_j - z_{j'} \neq 1 \pmod{96},$$

which makes the possible number of structures decrease to $64! / (32!)^2 \approx 2^{60.7} \approx 10^{18}$.

(iii) The present structure has rhombohedral symmetry, and this means that the structure can be specified by an arrangement of ABC layers at 32 successive positions. This gives a limitation

$$z_j - z_{j'} \neq 32 \pmod{96},$$

which reduces the number of possible arrangements down to $2^{32} \approx 10^{9.6}$.

(iv) So far only permutation on a line has been considered, but permutation on a circle should further be considered because of the periodic nature of crystals. Moreover, it is usually impossible to distinguish two arrangements on a circle with the reverse order. It is, therefore, sufficient to consider a necklace permutation, which gives us possible structures of about $2^{32} / (96 \times 2) \approx 2^{24.4} \approx 10^{7.3} \approx 22,370,000$ arrangements or 24.4 bits.

Those structures which have survived the above limitations (i)–(iv) will hereafter be called *the big set*.

(v) We have, in addition, the experimental fact that all polytypes of SiC so far found have structures such that h does not follow h in terms of the Wyckoff–Jagodzinski h - k notation, except for the $2H$ wurtzite type (Jagodzinski, 1949; Krishna & Verma, 1965). This limitation gives the condition that $z_{j+1} - z_j \leq 4$ when z_j 's are arranged in increasing order. More details about these conditions expressed in z 's have been described elsewhere (Tokonami, 1966).

(vi) Among all kinds of necklace arrangements with 96 beads, there are some arrangements with a periodicity of 48 or other shorter lengths. However, these were excluded in advance in an experimental analysis because of a definite periodicity of the crystal being specified from the very beginning. Exactly speaking, this limitation should duly be added at every stage of the limitations (i) to (v). For instance, the possible number of structures is reduced from 2^{32} to $2^{32} - 2^{16}$ when the present limitation is added to condition (v). As seen from this example, this modification is always numerically negligible.

All arrangements of z_j 's satisfying conditions (i) to (vi) in the above will later be referred to as *the small set*. The number of arrangements in the small set, which have been counted one by one on an electronic computer PC-2 (the commercial name is FACOM 202), is 25780 (corresponding to 14.65 bits).

Summarizing the above, the information amount obtained by condition (i) is $480 - 84.6 \approx 395$ bits, that by (ii) is $84.6 - 60.7 \approx 24$ bits, that by (iii) is $60.7 - 32 \approx 29$

of the small set are listed in Table 2. The information amount $H(I)$ included in I^l for all members of each set calculated by use of formula (1) is listed in Table 3. The simultaneous distribution of I^{l_1} and I^{l_2} has been calculated for several pairs of l_1 and l_2 reflexions; then the information amount $H(I_1, I_2)$ [formula (2)] and the redundancy $R(I_1, I_2)$ [formula (3)] have also been calculated. Representative pairs which may be expected to have a heavy correlation are those with $l_1:l_2=1:2$ and $l_2=l_1+3$ (neighbouring reflexions). Some other pairs which may not be expected to have even a slight correlation were also chosen for comparison. According to the result shown in Table 4, correlation can hardly be found except for a pair of reflexions with $l_1=16$ and $l_2=32$ which are two larger available divisors of 96.

The value of $H(I)$ averaged over 32 l values is 3.6 and 3.2 bits respectively for the big set and for the small set. This value multiplied by the number of reflexions is far more than 24.4 and 14.65 bits to be recovered. It means that there should exist much redundancy in a whole set of I^l s. From the fact that redundancy hardly exists between any pair of I^l s except I^{16} and I^{32} , a large amount of redundancy should exist mostly among simultaneous distributions of three or more I^l s.

Table 3. Information amount in the case when the intensity is classified into 16 degrees

l	$H(I)$ in the big set	$H(I)$ in the small set
1	3.70	2.13
4	3.58	2.16
7	3.46	2.19
10	3.48	2.29
13	3.46	2.51
16	2.74	2.18
19	3.63	3.33
22	3.67	3.81
25	3.79	3.95
28	3.87	3.36
31	3.96	3.53
34	3.98	3.82
37	3.94	3.91
40	3.76	3.53
43	3.35	3.82
46	3.17	3.66
49	3.00	3.63
52	3.12	3.76
55	3.64	3.89
58	3.88	3.92
61	3.98	3.86
64	3.66	3.47
67	3.95	3.30
70	3.81	3.80
73	3.79	3.94
76	3.66	3.49
79	3.54	2.96
82	3.55	2.58
85	3.46	2.35
88	3.48	2.21
91	3.58	2.17
94	3.72	2.15
Mean	3.6	3.2

Table 4. Information amount of simultaneous distribution of $I(l)$

The values of l, l' such as 2, 5, 8, ..., $l=(3n-1)$, ... correspond to 94, 91, 88, ..., $(96-l)$, ... in other tables, respectively.

		On the big set		On the small set	
l	l'	$H(l, l')$	$H(l) + H(l') - H(l, l')$	$H(l, l')$	$H(l) + H(l') - H(l, l')$
1	2	7.36	0.05		
2	4	7.26	0.05	4.25	0.06
4	8	7.02	0.04	4.36	0.01
8	16	6.15	0.06		
16	32	5.81	0.58	5.27	0.38
5	10	7.02	0.04		
10	20	7.10	0.05		
20	40	7.37	0.04		
40	80	6.34	0.15		
7	14	6.97	0.03		
14	28	7.38	0.04		
19	41	7.22	0.05	7.20	0.02
34	37	7.88	0.05	7.70	0.03

(b) The case when intensity has been divided into only two degrees, strong and weak

In the previous section $H(I)$ was found to be 3.6 bits on average for the big set. Similar calculations gave $H(I)=2.7, 1.7$ and 0.8 bits when the number of degrees of intensity was reduced to 8, 4 and 2 respectively. As is to be expected, less redundancy is obtained when the intensity is specified with less accuracy. In other words, the better the accuracy of measurements, the easier the structure analysis because of the more redundancy.

If the intensity is specified with less and less accuracy, the total information given by I^l s for all l finally does not reach the necessary information amount required to solve the structure or at least to determine the structure uniquely apart from homometric ones. In the present case where the intensity is specified by two degrees, a sum of $H(I)$ on l is 25.34 and 18.04 bits for the big set and the small set respectively. These figures are probably less than those figures, 24.5 and 14.65, to be recovered if redundancy is subtracted. To speak in terms of averages, the structure cannot be determined in this case.

Actually, however, it is possible to determine the structure even in such a case, provided that the intensity distribution of reflexions has marked characteristics as seen in the example of 96R-SiC. This real structure (Tokonami, 1966) consists mostly of the 6H type with some 21R. The observed data classified into strong and weak show more or less disagreement with I^l values calculated for every member of the small set. The number of reflexions which showed such a discrepancy was counted by checking each member of the small set. The statistics of all members as regards the number of reflexions which showed a discrepancy is listed in Table 5. Fortunately only one structure was found to show no discrepancy, while the three next nearly similar structures already showed disagreement for two reflexions. The majority of members showed discrepancies for 10 to 14 reflexions. Although no survey was

Table 9. *The information amount contained in $V_{z/96}$ (on the small set)*

z	$H(z)$
1	0.00
2	1.81
3	1.81
4	3.73
5	2.11
6	3.83
7	3.92
8	3.65
9	3.73
10	3.85
11	3.76
12	3.80
13	3.74
14	3.77
15	3.76
16	3.30
17	3.81
18	3.74
19	3.72
20	3.86
21	3.77
22	3.78
23	3.80
24	3.81
25	3.72
26	3.83
27	3.87
28	3.58
29	1.81
30	3.58
31	2.09
32	0.00
33	2.09
34	3.58
35	1.81
36	3.59
37	3.86
38	3.84
39	3.72
40	3.82
41	3.82
42	3.75
43	3.86
44	3.74
45	3.90
46	3.84
47	3.74
48	3.29

realizable crystal. The set of all meaningful points in \mathbf{I}^h space, which will be denoted by M , cannot be outside a super-cube with length 1 for each edge in this space. If a point in \mathbf{x}_j space is shifted along each of the J coordinate axes by an infinitesimal distance, a corresponding point in \mathbf{I}^h space shifts along each of J directions also by an infinitesimal distance. Thus any movement in \mathbf{x}_j space can be mapped in a sub-space of J dimensions in \mathbf{I}^h space, and the whole \mathbf{x}_j space is mapped in \mathbf{I}^h space as a sub-space M of J dimensions.

The extent to which the structure analysis has been well carried out is usually expressed by a parameter such as the R index or R factor, which corresponds to the distance between the observed point I_{obs} and the calculated point I_{cal} in \mathbf{I}^h space. When the point

I_{obs} happens to belong to M , it is possible to make this distance zero. However, in general this would not happen.

The sub-space M does not distribute uniformly in \mathbf{I}^h space, being dense in some regions and disperse elsewhere. Two non-equivalent points more or less distant from each other in \mathbf{x}_j space sometimes happen to have their mapping points very close to each other in \mathbf{I}^h space; these structures are pseudohomometric, and it is difficult to tell which point represents the true structure, when the observed point falls near these two points. In such a case, more accurate observed intensities are required to determine the structure. On the other hand, when the relevant point falls in the sparsely populated region, even less accurate measurements will enable us to distinguish the true structure from a nearly similar but false one.

Suppose we plot in \mathbf{x}_j space the value of the distance of I_{cal} from the fixed point I_{obs} in \mathbf{I}^h space. If observed data are good enough, the true structure will agree with the point which has the minimum distance. This multi-dimensional map may be called 'an R factor map' as used in the review articles (Hosoya, 1961, 1964). This idea seems to have been used in various investigations (for instance, Hosoya, 1958) and has extensively been studied by Milledge (1962). The R factor map with contour lines drawn may have numerous peaks and troughs, and thus the \mathbf{x}_j space may be divided into many regions, as it were, by a multi-dimensional network of watersheds. The number of these regions may be large but will certainly be far smaller than the number of points to be considered in the \mathbf{x}_j space. Improvement of the Monte Carlo method of Vand & Niggli (1961) so as to make it practicable was effected only by introducing the optimal shift method (Niggli, Vand & Pepinsky, 1961), which is nothing but the Monte Carlo method for the above-mentioned regions.

Karle & Hauptman (1964) used H' terms which are beyond the range of the observed H terms, when they devised a method of refining the Patterson map. It may not be self-evident that it is allowable to use unobserved reflexion terms outside of the limiting sphere. However, this may be understood by the following considerations. When the projection onto a partial space with H dimensions is given, the sub-space M in \mathbf{I}^h space with H' dimensions more or less restricts the coordinates in certain regions, and, therefore, the values of $H' - H$ unobserved terms can be inferred to some extent.

In the above discussion it has been assumed that the distribution of the points in \mathbf{x}_j space is uniform. But some regions in \mathbf{x}_j space are entirely prohibited, for instance, because of steric hindrance due to the finite size of atoms. This fact may be clearly seen in a short note by Goedkoop, MacGillavry & Pepinsky (1951). Moreover, such a non-uniformity comes from coordination of atoms in covalent crystals, from Pauling's rules in ionic crystals or from the form of molecules in molecular crystals. In the example in the previous

chapter, 32^96 points in x_j space were reduced to several tens of millions because of symmetry and steric hindrance. A set of the observed I_{obs} values is not arbitrary but, as was pointed out in the above, more or less restricted. It is therefore favourable to take the weight of the points in x_j space into consideration. The results on SiC in the previous chapter showed various expectation values for reflexion intensities. This comes from the non-uniformity in x_j space, because both the big set and the small set include only those structures which satisfy certain limitations.

In relation to this discussion, it should be mentioned here that Hauptman (1964) showed how to infer the shape of the molecule directly from the intensity distribution without the information about phases. Conversely, when the shape of a molecule or a chemical unit of a crystal is known, a kind of unitary intensity, or the ratio of the observed intensities and the square of an absolute value of the structure factor for the above-mentioned unit, is very helpful to reduce the number of parameters J .

Let M' denote the mapped points of x_j space limited by steric hindrance and other possible information; then M' is a subset of M and has more sparse distribution than M has in I^h space, and therefore the requirement of experimental accuracy of I_{obs} for determining the structure uniquely becomes less severe. The extreme case is seen in the structure analysis of 96R-SiC (Tokonami, 1966), in which only the qualitative observed values were used; on modification so as to make the structure point fall on M' , utilizing certain algebraic characters found in the Patterson function,

the unique solution was obtained. In usual crystal analyses, it often happens that a few reflexions are considered to be subject to extinction effects and their intensity data are taken into account with small weights. This may be allowed from the viewpoint that only the points on M' have meaning in I^h space.

The authors express their thanks to Prof. S. Miyake for his advice in improving the method of presentation of these concepts.

References

- DIAMOND, R. (1963). *Acta Cryst.* **16**, 627.
 GOEDKOOP, J. A., MACGILLAVRY, C. H. & PEPINSKY, R. (1951). *Acta Cryst.* **4**, 491.
 HAUPTMAN, H. (1964). *Acta Cryst.* **17**, 1421.
 HOSOYA, S. (1958). Dissertation, Univ. of Wales.
 HOSOYA, S. (1961). *J. Cryst. Soc. Japan*, **3**, no. 1, 2 (in Japanese).
 HOSOYA, S. (1964). *J. Cryst. Soc. Japan*, **6**, 56 (in Japanese).
 JAGODZINSKI, H. (1949). *Acta Cryst.* **2**, 201.
 KARLE, J. & HAUPTMAN, H. (1964). *Acta Cryst.* **17**, 392.
 KRISHNA, P. & VERMA, A. R. (1965). *Z. Kristallogr.* **121**, 36.
 MILLEDGE, H. J. (1962). *Proc. Roy. Soc. A*, **267**, 566.
 NIGGLI, A., VAND, V. & PEPINSKY, R. (1961). In *Computing Methods and the Phase Problem in X-ray Crystal Analysis*, p. 161. Oxford: Pergamon Press.
 SAYRE, D. (1952). *Acta Cryst.* **5**, 843.
 SHANNON, C. E. (1948). *Bell Syst. T. J.* **27**, 379.
 TOKONAMI, M. (1966). *Miner. J. Japan*, **4**, 401.
 TOKONAMI, M. & HOSOYA, S. (1965). *Acta Cryst.* **18**, 908.
 VAND, V. & NIGGLI, A. (1961). In *Computing Methods and the Phase Problem in X-ray Crystal Analysis*, p. 266. Oxford: Pergamon Press.

Acta Cryst. (1967). **23**, 25

Prévision de quelques Images de Dislocations par Transmission des Rayons X (Cas de Laue symétrique)

PAR DANIEL TAUPIN

*Centre de Calcul Numérique, Laboratoire de Physique Théorique et Hautes Energies,
Faculté des Sciences, 91-Orsay (Seine et Oise), France*

(Reçu le 12 octobre 1966)

The author's dynamical theory of the diffraction of X rays by distorted crystals is applied to the computation of the images obtained in the Laue case (transmission) when a single dislocation is contained in an otherwise perfect crystal. Some results which have been obtained for various orientations of the Burgers vector are presented: dislocation far from the upper surface (entry), the width of the incident beam being infinite; dislocation in the vicinity of the upper surface, the width of the incident beam being infinite; dislocation in the vicinity of the upper surface, the incident beam being very narrow. The results are in good agreement with topographs published by different authors.

Dans un précédent travail (Taupin, 1964a, b) nous avons montré comment, à partir des équations de Maxwell, on pouvait étendre la théorie dynamique, maintenant devenue classique, de la diffraction des ray-

ons X par les cristaux parfaits au cas où les cristaux comportaient des défauts ou déformations élastiques, sans toutefois être limité à l'approximation en colonne comme dans la théorie de Howie & Whelan (1961,