

ture; instead, they are reiterated in hierarchic fashion ranging from the whole protein monomer through supersecondary structures down to individual helices and strands.

The concept of a domain is subject to reinterpretation when considered in the broader context of its molecular hierarchy. The conventional large domains cited in the literature (near the top of the hierarchy) are seen to be structural composites, with subparts that are domains in their own right. Thus, the large, spatially distinct chain segments that can be resolved within a protein may not represent strictly autonomous units, but may arise instead in consequence of a concluding step in the sequential folding process.

Acknowledgments

X-Ray coordinates were provided by the Brookhaven Data Bank.²⁹ The work was supported by PHD Grant 29458 and by an NIH Research Career Development Award.

²⁹ F. C. Bernstein, T. G. Koetzle, G. J. B. Williams, E. F. Meyer, Jr., M. D. Brice, J. R. Rogers, O. Kennard, T. Shimanouchi, and M. Tasumi, *J. Mol. Biol.* **112**, 535 (1977).

[30] Calculation of Molecular Volumes and Areas for Structures of Known Geometry

By FREDERIC M. RICHARDS

Introduction

At the macroscopic level the concepts of area and volume are quite clear, and the methods of measurement straightforward in principle. At the level of individual molecules the definitions become less obvious. When different methods of measurement are used, it is not clear that area and volume are single valued characteristics. For a one-component system of known molecular weight, the mass density is sufficient to unambiguously define the mean volume per molecule. For a system of two or more components, thermodynamic parameters, the partial molar volumes, can be defined uniquely, but these values bear no necessary simple relation to the actual physical volumes of the individual components which, in turn, may not be uniform throughout the mixture. The concepts of total molecular area are even more elusive (see below).

For the purpose of this discussion we shall concentrate on the definition and calculation of certain geometrical volumes and areas that can be derived from high-resolution structural data. No attempt will be made to give a critical discussion of the use (or potential misuse) of these values.

The following basic input data are required:

1. The list of Cartesian coordinates of the atom centers.
2. One or more lists of attributes for each atom:
 - a. Assigned van der Waals radii for each atom (required in most procedures).
 - b. Assigned covalent radius for each atom (use depends on the volume algorithm selected).
 - c. Covalent connectivity (needed if b is used, and also to assemble atoms into reasonable packing groups)
3. Choice of radius for a spherical probe.

In the presently used procedures no allowance is made for errors or fluctuations in the coordinates. The structure is simply a collection of points in three-dimensional space. Sensitivity to error, if required, is estimated by repeating the calculation with appropriately altered coordinate lists and comparing the results. None of the algorithms place any restrictions on the position of the atom centers.

The van der Waals Envelope

Because of the diffuse radial distribution of the electron density surrounding any atomic center, the apparent position of the surface of a molecule will depend on the technique used to examine it. For chemically bonded atoms the distribution is not spherically symmetric nor are the properties of such atoms isotropic. In spite of all this the use of the hard sphere model has a venerable history and an enviable record in explaining a variety of different observable properties. As applied specifically to proteins, the work of G. N. Ramachandran and colleagues has provided much of our present thinking about permissible peptide chain conformation.¹ Different approaches using more realistic models, complex mathematics, and even quantum mechanical approximations have improved the details but have not altered the basic outline provided by the hard sphere approximation. The steepness of the repulsive term in the potential function for nonbonded interactions is responsible for the success of "hard" in the hard sphere.

In spite of the general success of the hard sphere approximation, the van der Waals envelope of a molecule is not unique and is defined differ-

¹ G. N. Ramachandran and B. Sasisekharan, *Adv. Protein Chem.* **23**, 284 (1968).

ently for different purposes. The bases on which the radii of the individual atoms are derived and the uses to which they are put differ. There may be no simple set of "correct" values. The radii are closely connected to the nonbonded potential energy function. Given this function for a pair of atoms the sum of the radii may be equated either to the value of the interatomic separation at the minimum or to the smaller value where the potential is zero. In the biochemical literature the former has been more commonly used. For the Lennard-Jones 6-12 potential the two values differ by about 12%.

The parameters of the nonbonded potential functions can be derived by fitting the functions to the observed packing in molecular crystals and making use of the structural and thermodynamic data available for a large number of such crystals.² The derived functions, however, will depend on how the lattice energy is partitioned. A full expression for the lattice energy may include the electrostatic interaction of all partial charges, bond lengths and angles, special treatment of hydrogen bonds, etc. The nonbonded terms will then account only for the residual energy. On the other hand the partial charges can be omitted and the total energy partitioned among the various pair-additive, nonbonded terms. The expressions, and thus the derived radii, will, of course, be different even for the same input data. The proper use of these numbers requires that the basis for their derivation be known and appropriately accounted for.

When working with macromolecules it is frequently convenient not to deal with individual atoms but with small groups of atoms which are considered to be adequately represented by a sphere and characterized by a single radius. These groups usually consist of a single heavy atom and one or more hydrogen atoms. (The relation between the group radii and the individual atom radii are not always obvious.) Such groups have been referred to as "unified atoms" by Dunfield *et al.*³ and as "extended atoms" by Karplus and colleagues.⁴ Some group radii that have been used in various studies are listed in the table.

The complex surface that results from the intersection of a number of spheres is referred to as the van der Waals surface. This surface has a defined area and it encloses a defined volume. Although the construction is easy to visualize and is logically consistent, it should be recognized that no chemical procedure ever directly measures this particular area or volume. However, the various areas and volumes computed by algorithms

² F. A. Momany, L. M. Carruthers, R. F. McGuire, and H. A. Scheraga, *J. Phys. Chem.* **78**, 1595 (1974).

³ L. G. Dunfield, A. W. Burgess, and H. A. Scheraga, *J. Phys. Chem.* **82**, 2609 (1978).

⁴ B. R. Gelin and M. Karplus, *Biochemistry* **18**, 1256 (1979).

SOME LISTS OF VAN DER WAALS RADII FOR SELECTED GROUPS OF ATOMS

Symbol	Designation	Bondi ^a	Lee and Richards ^b	Shrake and Rupley ^c	Richards ^d	Chothia ^e	Richmond and Richards ^f	Gelin and Karplus ^g	Dunfield <i>et al.</i> ^h and Nemethy <i>et al.</i> ⁱ
—CH ₃	Aliphatic, methyl	2.0	1.80	2.0	2.0	1.87	1.9	1.95	2.13
—CH ₂ —	Aliphatic, methyl	2.0	1.80	2.0	2.0	1.87	1.9	1.90	2.23
>CH—	Aliphatic, CH	—	1.70	2.0	2.0	1.87	1.9	1.85	2.38
≡CH	Aromatic, CH	—	1.80	1.85	<i>j</i>	1.76	1.7	1.90	2.10
>C=	Trigonal or aromatic	1.74	1.80	1.5/1.85	1.7	1.76	1.7	1.80	1.85
—NH ₃ ⁺	Amino, protonated	—	1.80	1.5	2.0	1.50	.7	1.75	—
—NH ₂	Amino or amide	1.75	1.80	1.5	—	1.65	1.7	1.70	—
X (O/N)	Amide (N or O unknown)	1.75	1.55	1.5	1.6	1.65	1.7	—	—
>NH	Peptide, NH or N	1.65	1.52	1.4	1.7	1.65	1.7	1.65	1.75 (N)
=O	Carbonyl oxygen	1.5	1.80	1.4	1.4	1.40	1.4	1.60	1.56
—OH	Alcoholic hydroxyl	—	1.80	1.4	1.6	1.40	1.4	1.70	—
—OM	Carboxyl oxygen	—	1.80	1.89	1.5	1.40	1.4	1.60	1.62
—SH	Sulphydryl	—	1.80	1.85	—	1.85	1.8	1.90	—
—S—	Thioether or —S—S—	1.80	—	—	1.8	1.85	1.8	1.90	2.08

^a A. Bondi, "Molecular Crystals, Liquids and Glasses." Wiley, New York, 1968. Radii assigned on the basis of observed packing in condensed phases.

^b Lee and Richards.²⁵ Values adapted from A. Bondi, *J. Phys. Chem.* **68**, 441 (1964).

^c Shrake and Rupley.²⁶ Values taken from L. C. Pauling, "The Nature of the Chemical Bond," 3rd ed. Cornell Univ. Press, Ithaca, New York, 1960.

^d Richards.⁶ Minor modification and extension of Bondi (1968) set (see footnote a. above). Rationale not given.

^e Chothia.¹⁵ From packing in amino acid crystal structures. Personal communication from T. Koetzle quoted.

^f Richmond and Richards.¹⁴ No rationale given for values used.

^g Gelin and Karplus.⁴ Origin of values not specified.

^h Dunfield *et al.*³ Detailed description of deconvolution of molecular crystal energies. Values represent one-half of the heavy-atom separation at the minimum of the Lennard-Jones 6-12 potential functions for symmetrical interactions.

ⁱ G. Nemethy, M. S. Pottle, and H. A. Scheraga, *J. Phys. Chem.* **87**, 1883 (1983).

^j See original paper.

that are closer to any other atom center than they are to i . The reduced list provides the vertices of the limiting polyhedron from which, in turn, the faces and edges can be derived.

Approach Based on the Convex Nature of the Polyhedron. Since the limiting polyhedron surrounding atom center i is convex, all vertices are either in a given face plane or on the same side of that plane as i .⁶ Any potential vertex on the opposite side of any of the faces is not a member of the final set of vertices. In an iterative procedure vertices are then retained or rejected on the basis of position with respect to each plane in the current set of faces.

The algorithm for implementing this procedure starts by setting up the equations for all planes in the set $\{A\}$ of atoms around i . (See below for the selection of $\{A\}$ and for the various equations that may be used to define these planes.) An arbitrary but very large tetrahedron (i.e., four vertices and four face planes) is set up around i . The position of each of the four vertices with respect to the planes in $\{A\}$ is then examined. A vertex (vertices) not on the same side as i of a given plane is (are) eliminated and replaced by new vertices produced by the selecting plane and the planes contributing to the eliminated vertex or vertices. The index designation for planes, lines, and vertices discussed above make it easy to do the bookkeeping at this stage. The original file of planes for $\{A\}$ is arranged and searched in order of increasing distance from i to the plane. After each plane is checked, the vertex list is changed as required. One pass through the list of planes is then sufficient to yield the limiting polyhedron. The file now contains the position of each vertex and the equation for each face. This procedure is general and is independent of the formula by which the equations for the planes are developed. While the limiting polyhedrons are uniquely defined, the full set of polyhedrons may or may not accurately account for all space depending on the definition used for the planes associated with $\{A\}$.

Approach Based on Distance Selection. For the Voronoi and Radical Plane procedures (see below) each potential vertex can be located at a defined distance from each of four atom centers, one of which is i .^{7,8,10} All potential vertices from the atom set $\{A\}$ are calculated. For a vertex to be part of the limiting polyhedron, it must be no closer to any other atom center than it is to i . The distance of each vertex to all atoms in $\{A\}$, other than the four defining atoms, is tested against the distance to i , and is accepted or rejected on this basis. From the indices identifying each vertex in the final list the faces of the polyhedron can be established by searching the list for all sets of triplets having a common index. Thus a

⁶ B. J. Gellatly and J. L. Finney, *J. Mol. Biol.* 161, 305 (1982).

final list is obtained giving each face, the number of vertices in each face, and the total number of faces and vertices as in the first procedure. This approach requires that the distance relations be specified as equalities, but the final set of Voronoi or radical plane polyhedrons do account accurately for all space.

The Atom Set $\{A\}$. For either procedure the efficiency of the calculation depends on making the set of atoms, $\{A\}$, as small as possible. For a completely arbitrary set of points, the total list would have to be surveyed, in principle, since there would be no way in advance of knowing how asymmetric any limiting polyhedron might be. In practice with macromolecules this is not a problem since the points are reasonably uniformly distributed and the resulting polyhedra are quite compact. In the program used by Richards⁶ it was found by trial that all necessary positions were included with a comfortable margin of safety if only atoms less than 6.5 Å from i were selected for $\{A\}$. The time for atom selection from the coordinate list is minimized if the list is loaded into a coarse cubic lattice permitting a grid search.

A more systematic and less arbitrary procedure for the set selection was given by Brostow *et al.*⁹ The algorithm used for the construction of the limiting polyhedron is comparable to the convex polyhedron approach, but somewhat different criteria are used in selecting and limiting the atom list. The authors suggest that their algorithm is more efficient than that used by Richards or Finney, although no one seems to have made benchmark runs with all three programs on the same data set.

Volume of the Polyhedron

Once the vertex list is complete, the volume of the resulting limiting polyhedron is readily computed. The area of each face can be calculated from the component triangles. The length of the face normal to the center i is already known. The cone volumes associated with each face are summed. The individual atom volumes can be examined or more commonly combined into packing units such as side chains or whole residues.

Although the principal use of this procedure to date has been to assign individual atom volumes, Finney¹¹ has pointed out that the polyhedrons represent a wealth of information about the surroundings of each atom. The extent of interactions with each neighbor is reflected in a clearly defined way by the area of the shared face. The direction of the interaction is defined by the normal to the face. The neighbors may be other protein atoms or potential solvent molecules.

¹¹ J. L. Finney, *J. Mol. Biol.* 119, 415 (1978).

Selection of Position of the Plane along the Interatomic Vector

The Voronoi Construction. In the original Voronoi procedure the planes are drawn as bisectors of the lines between the points, as in Fig. 3 (see Finney⁷ and Richards,⁶ method A). The points are considered intrinsically equal and no special characteristics are assigned. If d_{ij} is the interatomic distance and p_{ij} the distance from atom i to the intersection of the plane and the vector, then

$$p_{ij}/d_{ij} = 1/2 \quad (1)$$

The procedure gives a single unique face between each atom pair. This face is part of the limiting polyhedrons about both atoms. The method is exact in that all of the space is precisely accounted for without error. The appropriate distance relations are

$$(x_i - x)^2 + (y_i - y)^2 + (z_i - z)^2 = L_c^2, \quad i = i, j, k, l \quad (2)$$

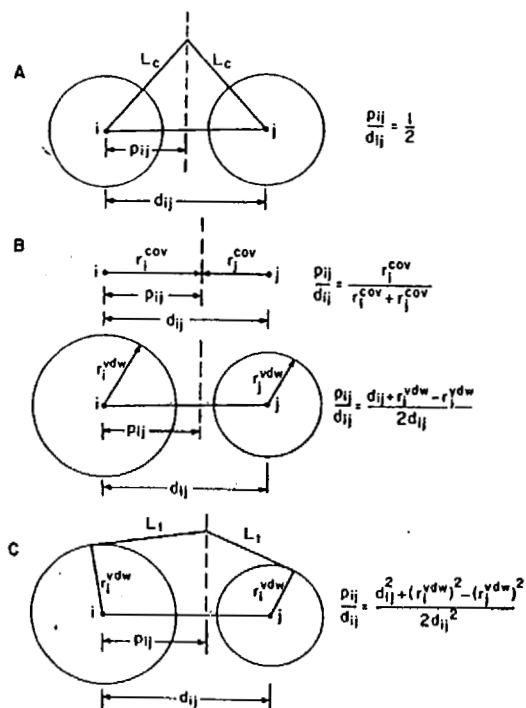


FIG. 3. Definitions of the planes dividing the interatomic vectors in the different partitioning procedures: (A) Voronoi; (B) Richards' method B; (C) radical plane.

where (x_i, y_i, z_i) and (xyz) are the coordinates of the four atoms and corresponding vertex, and L_c is the common vertex to atom center distance.

Richards' Construction B: While mathematically rigorous, the Voronoi procedure does not make much physical sense since different atoms do have different intrinsic sizes and clearly are not equal. In an attempt to overcome this difficulty, Richards²⁵ (method B) suggested a modified procedure in which the planes defining the polyhedron do not bisect the interatomic vector but cut it in a ratio which depends on the van der Waals or covalent radii of the two atoms involved (Fig. 2):

Covalent interaction:

$$p_{ij}/d_{ij} = r_i^{\text{cov}} / (r_i^{\text{cov}} + r_j^{\text{cov}}) \quad (3)$$

Noncovalent interaction:

$$p_{ij}/d_{ij} = (d_{ij} + r_i^{\text{vdw}} - r_j^{\text{vdw}}) / 2d_{ij} \quad (4)$$

This procedure assigns more space to those atoms which are intrinsically larger and less to the smaller atoms. While physically reasonable, the method has lost the mathematical rigor of the strict Voronoi procedure. All of the space is not accounted for. A shared face between two adjacent atoms will not necessarily contain the same vertices for the two polyhedrons surrounding the atoms, thereby leading to the vertex error problem. The little error polyhedrons are variable in size and position, but may, in aggregate, represent considerable volume. A careful comparison has recently been made by Gellatly and Finney¹⁰ who conclude that for ribonuclease S the modified procedure underestimates the total volume of the molecule by about 4%. The use of the procedure will thus depend very much on the purpose for which the numbers are to be used. While it is capable of including both covalent and noncovalent characteristics, the absolute volume errors are variable and may be substantial. Caution is indicated.

The Radical Plane Construction. A different space partitioning method has been developed by Fischer and Koch¹² and applied by Gellatly and Finney¹⁰ to the protein volume problem. The procedure is mathematically accurate and provides a rational basis for handling unequal spheres (Fig. 3). The radical plane is the locus of points from which the tangent lengths L_i to the two spheres are equal. The distance equations are now

$$(x - x_i)^2 + (y - y_i)^2 + (z - z_i)^2 - (r_i^{\text{vdw}})^2 = L_i^2, \quad i = i, j, k, l \quad (5)$$

¹² W. Fischer and E. Koch, *Z. Kristallogr.* 150, 245 (1979).

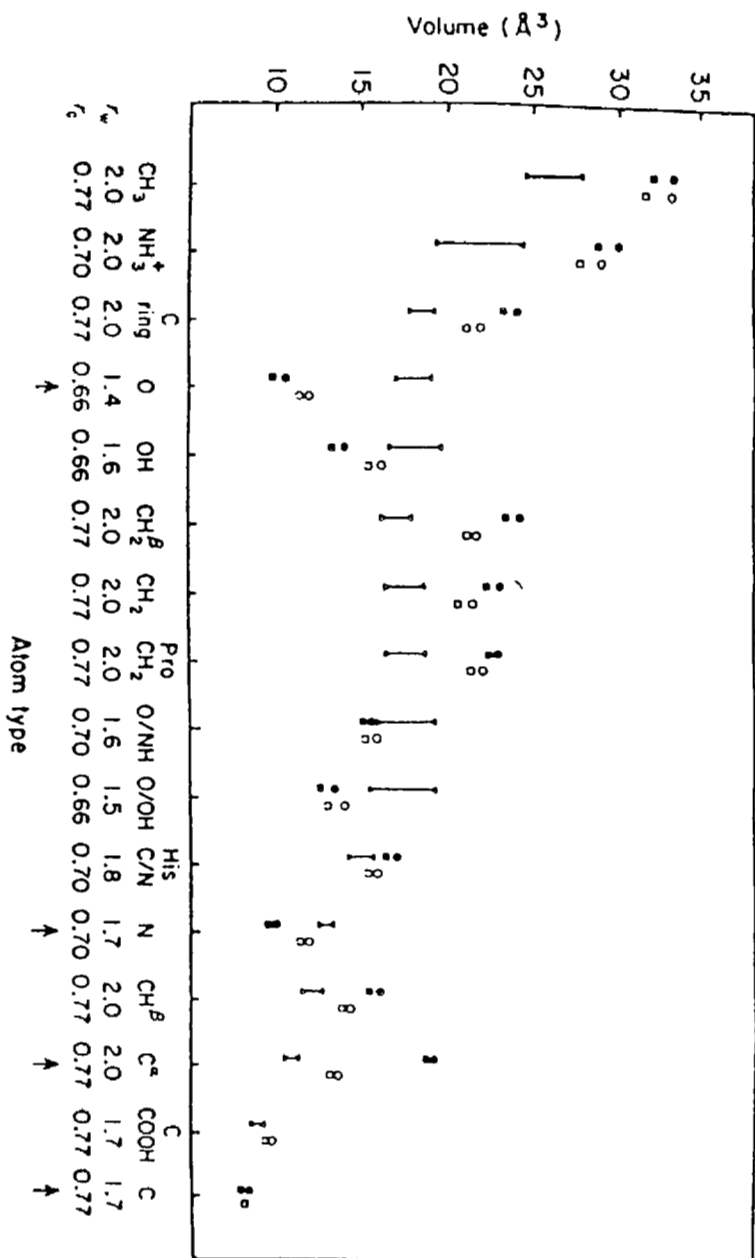


Fig. 4. Mean volumes for atom groups as calculated by the three procedures, and using the same data set for ribonuclease S. The van der Waals radii and the covalent radii used are given below each atom type. Arrows indicate main-chain atoms. The surface was defined by the uniform density probe procedure employing probes with either a 1.7 or 1.4 Å radius. (V) Voronoi, 1.7 Å; (▲) Voronoi, 1.4 Å; (●) Radical, 1.7 Å; (■) Radical, 1.4 Å; (○) Richards B, 1.7 Å; (□) Richards B, 1.4 Å. Radical and Voronoi values for C carbonyl groups are superimposed. (Reproduced with permission from Cellatly and Finney.¹⁰)

the vertex errors are individually calculated and corrected for (including sub-allocating the surface-involved error tetrahedra between protein and solvent) we argue in favour of the use of the radical method.

When considering volume calculations over atom groups, however, the choice is less clear: absolute volumes are of less interest than deviations of occupied volumes from the mean, and therefore as long as the vertex error is reasonably uniformly distributed this problem is less significant. Again we would reject the use of Voronoi's method because of the non-physical partitioning of non-bonded interactions, even though in some cases the consequent volume spread is masked by the local packing variations. Provided groups are chosen, with constant covalent environments (e.g. main chain atoms except glycine, whole side-chains) the differences in covalent treatment between radical and Richards' (and also Voronoi's methods) are completely cancelled, and as both Richards' and radical methods partition non-bonded interactions reasonably, there is little to choose between them.

For calculations of occupied volumes and volume distributions for single atoms (or atom groups such as CH₂), we would argue that no procedure is satisfactory unless the atoms are grouped together with a constant covalent environment, in which case the same considerations apply as for the larger groups such as main chain and side-chains. For a variable covalent environment, the spread of the resulting volume distributions will be significantly influenced by the placing of the covalent partitioning planes. The effect will be present for both Richards' and radical methods, though the form of the equation is such that the effect will be greatest for radical.

Clearly, it is impossible to devise a volume partitioning procedure that is both rigorous and consistent with the different chemical constraints in proteins. We can handle a system of interacting van der Waals' atoms rigorously, using radical planes, but as soon as we have to deal with covalent interactions, we must *either* use van der Waals' criteria to partition a covalent bond *or* abandon geometrical rigour.

We argue that a discussion of the preferability of using radical or Richards' method for examining packing efficiency of an atom or group of atoms with variable covalent environment would be largely academic and of little value. If we ask questions about packing efficiency, then covalent-bond partitioning is physically irrelevant, the identity of the repulsive electron shell between the two atoms having been lost in the covalent interaction. Therefore, any discussion of packing efficiency and variations for atoms or groups with a variable covalent environment must necessarily consider data that are perturbed by volume variations that are *not* due to the packing constraints being investigated. The perturbations will be smaller for Richards' than for the radical method, so if such comparisons are required, then Richards' method B is to be preferred over the radical planes method.

Other discussions and applications of the volume calculations are given by Richards¹⁴ and Chothia.¹⁵

Addendum on Cavities

Large Grid Approximation. In the above discussion all of the space inside the hypothetical solvent shell is assigned to the protein atoms. If

¹⁴ T. J. Richmond and F. M. Richards, *J. Mol. Biol.* **119**, 537 (1978).

¹⁵ C. Chothia, *Nature (London)* **254**, 304 (1975).

there is a hole in the structure, the volume that it represents is assigned to the surrounding atoms as specified by their limiting polyhedrons. Such a hole only appears as a lowering in the packing density for this group of atoms. If the hole is modest in size and the number of protein atoms large, the variation from the mean of the packing density may not be obvious and location of the hole not easy to derive. There is no unique and mathematically satisfactory solution to this problem so far, but several approximate procedures have been suggested.

In an attempt to focus on the cavity structure Richards¹⁶ made use of the large grid that he had set up in defining the solvent shell (see above). The grid positions outside of the van der Waals envelope of the protein but inside the solvent shell were used to define the cavities. The centers of these empty cubes were used as pseudoatoms in a modified Voronoi calculation (Fig. 5). The limiting polyhedrons were defined by neighboring grid positions and the van der Waals envelope of neighboring protein atoms (see Fig. 1b). Each empty grid position thus had a volume associated with it. The connectivity of these positions could be evaluated to get an idea of the volume and shape of the cavities. The procedure is crude and unlikely to give more than a very rough idea of the cavity distribution. Nonetheless, the volumes assigned to the protein atoms are lowered a little and the standard deviation of their volume distributions markedly reduced. The general cavity distribution can be visualized easily. The purpose of defining and examining the cavities was to consider volume fluctuations in the dynamics of proteins. Such fluctuations must reflect changes in cavity volume, since the van der Waals envelope is essentially incompressible under normal conditions.

The large grid (2.8 Å) used in the above calculation missed many of the smaller cavities which did not happen to include a grid position. If completely empty, a single grid cube is almost big enough to hold an entire water molecule.

Small-Grid Approximation. A. Perlo and F. M. Richards (unpublished) tried to improve on the estimate of total cavity volume by using a much finer grid (0.5 Å). The algorithm does not use the Voronoi procedure at all. When the van der Waals envelope of the protein is inserted into the lattice of this small grid, a single atom may cover 100 or more positions. The lattice is checked sequentially in each of the three principal lattice directions. Each empty position is given a number representing its minimum distance in grid units to the van der Waals envelope (Fig. 6a). The problem now is to decide which of these positions are truly external, and thus bulk solvent, and which are cavities either internal or non-solvent-

¹⁶ F. M. Richards, *Carlsberg Res. Commun.* 44, 47 (1979).

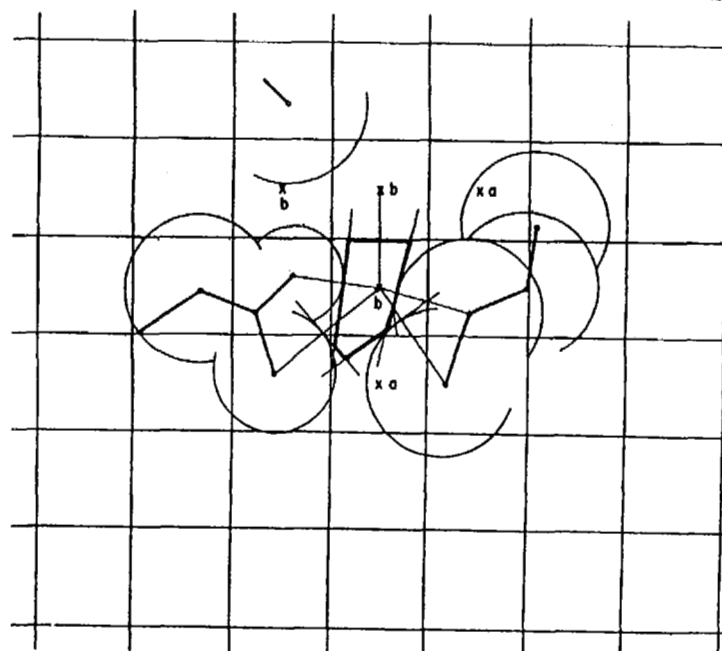


FIG. 5. Volume construction for cavities using the coarse grid (edge = 2.8 Å) and Richards Method B. The van der Waals envelopes of a few atoms are shown. Grid positions whose centers are inside the van der Waals envelope of the protein are labeled a, those outside b. The b locations are used as pseudoatom positions in the volume calculation (see text for discussion). (Reproduced with permission from Richards.¹⁶)

accessible on the surface. The probe is a string of digits against which the lattice positions are checked. With a 0.5 Å grid a water molecule is a sphere with a diameter of about seven grid positions. The closest grid approximation of a sphere would be represented as 1, 3, 3, 4, 3, 3, 1. For ease in subscript manipulation the actual approximation used is 1, 2, 3, 4, 3, 2, 1.

Starting from the edge of the lattice box, known to be "outside" the protein, each group of seven consecutive lattice positions is tested against the probe. If each lattice position is not a protein position and is characterized by a number equal to or greater than that of the probe, then all of these positions could be part of a solvent molecule and are so designated. Otherwise they are classed as potential cavity positions (Fig. 6b). As the probe moves along an axial direction the number of protein surfaces passed is counted in order to assign cavities as external or internal. The

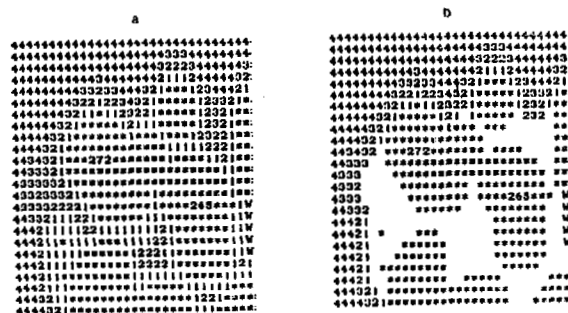


FIG. 6. Cavity definition by the small grid procedure of Perlo and Richards. Grid spacing is 0.5 Å. A small section of a single plane through pancreatic trypsin inhibitor is shown. The asterisks identify those grid positions which are inside the van der Waals envelope of the protein. The numbers in the asterisk area are the serial numbers from the coordinate list of atoms whose centers happen to lie in or close to this plane. In (a) each grid position outside of the asterisk area contains a number which represents the distance in grid units to the nearest part of the van der Waals envelope in any of the three axial directions. Distances equal to or greater than 4 are also listed as 4. The test of this filled lattice for potential cavity positions is described in the text. Such cavity positions located by the algorithm are shown as clear areas in (b). The "W"s are part of the van der Waals envelope of one of the four interior water molecules identified in this structure.

initial cavity list is large. This is reduced as the probe check is carried out in the other two axial directions. From the final list the sum of the number of cavity positions gives directly the total cavity volume. The total protein volume is given by the sum of this volume and the volume (number of grid positions) inside the van der Waals envelope. A cavity as small as a single grid position (0.125 Å³) will be recognized as will internal cavities large enough to hold a solvent molecule. A very similar procedure has been used by Kossiakoff.^{17,18}

A complete test of the approximations in this algorithm has not been made, but it has been used for estimating the volume fluctuations during a molecular dynamics simulation of pancreatic trypsin inhibitor (A. Perlo, F. M. Richards, N. Swaminathan, and M. Karplus, unpublished). It has also been used by Pickover and Engelman¹⁹ in their study of the extended low angle X-ray scattering curves of solutions of several proteins.

The procedures of Connolly²⁰ for depicting protein surfaces (referred to below) can also be used to identify cavities and to estimate their vol-

¹⁷ A. Kossiakoff, *Nature (London)* 296, 713 (1982).

¹⁸ A. Kossiakoff, *Brookhaven Symp. Biol.* 32, 281 (1983).

¹⁹ C. A. Pickover and D. M. Engelman, *Biopolymers* 21, 817 (1982).

²⁰ M. Connolly, Ph.D. Dissertation, University of California, Berkeley (1981).

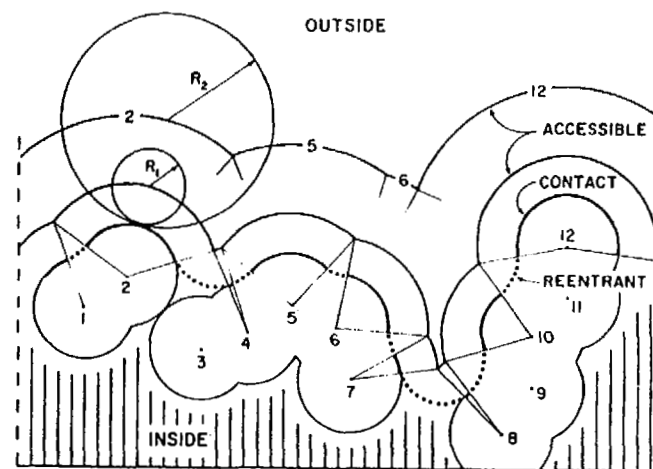


FIG. 7. Schematic representation of possible molecular surface definitions. A section through part of the van der Waals envelope of a hypothetical protein is shown with the atom centers numbered. The accessible surfaces generated by two probes of different size, R_1 and R_2 , and the geometrical definition of contact and reentrant surfaces are shown. (Reproduced with permission from Richards.²¹)

ume. In this procedure cavities appear only if they are large enough to contain at least one water molecule. In his survey of a number of proteins with this algorithm, Connolly found that the total cavity volume was of the order of 3% or less of the total volume of the protein.

A different approach developed for glass structures but not yet applied to proteins has been described by Finney and Wallace.²¹

Area

Definitions

On the molecular scale any conceivable probe has dimensions comparable to the features of the surface being examined. Consider the cross section of part of the surface of the hypothetical macromolecule shown in Fig. 7. The trace of the van der Waals envelope of some of the atoms of the structure is shown. A spherical probe of radius R_1 is allowed to roll on the outside while maintaining contact with the van der Waals surface. It will never contact atoms 3, 9, or 11. Such atoms are considered not to be part of the surface of the molecule and are referred to as interior atoms.

²¹ J. L. Finney and J. Wallace, *J. Non-Cryst. Solids* 43, 167 (1981).

The question of how to define and quantitate the surface is a matter of convenience. One straightforward procedure is simply to use the continuous sheet defined by the locus of the center of the probe, the "accessible surface." Another alternative would be to consider the "contact surface," those parts of the molecular van der Waals surface that can actually be in contact with the surface of the probe. This would provide a series of disconnected patches. The "reentrant surface" is also a series of patches defined by the interior-facing parts of the probe when it is simultaneously in contact with more than one atom. Considered together the contact and reentrant surfaces represent a continuous sheet, which might be called the "molecular surface."²²

By the nature of the geometrical construction there are no reentrant sections of the accessibility surface, i.e., viewed from the molecule each spherical segment is convex. This does entail a possible loss of information as the ratio of contact-to-reentrant surface may be a useful measure of molecular surface roughness. This can be seen qualitatively by inspecting Fig. 7. The molecular surface also has the advantage that the area approaches a finite limiting value as the size of the probe increases. To date most reports have calculated and discussed the accessible area.

With any of the surface definitions the actual numbers derived will depend on the radius chosen for the probe. An example of the change that is produced by probe size is shown in Fig. 7. In going from R_1 to R_2 the number of noncontact or interior atoms increases from three to eight. The accessible surface becomes much smoother (as does the molecular surface, not shown); there is only a slight dimple replacing the deep crevice revealed by the R_1 probe. The appearance of deeply convoluted features or actual holes in the interior of the protein becomes very sensitive to the choice of probe radius. The smaller the probe the larger the number of feature that will be revealed. About the smallest physically reasonable probe is a water molecule considered as a sphere of radius of 1.4 or 1.5 Å. The ratio of this number to the van der Waals radii assumed for the individual atom or atom groups will markedly affect the calculated areas for individual atoms.

The accurate calculation of the surface area is a complex geometrical problem. This problem has in fact been solved rigorously and a closed form analytical expression has been derived.^{23,24} Various approximate methods, whose accuracy varies but is adequate for many purposes, have been more widely used.

²² F. M. Richards, *Annu. Rev. Biophys. Bioeng.* **6**, 151 (1977).

²³ T. J. Richmond, *J. Mol. Biol.* **178**, 63 (1984).

²⁴ M. Connolly, *J. Appl. Crystallogr.* **16**, 548 (1983).

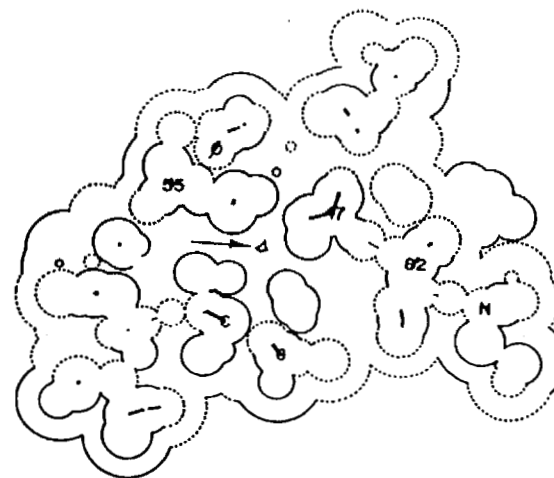


FIG. 8. Superposition of sections through the van der Waals envelope and the accessible surface of ribonuclease S. The arrow indicates a cavity inside the molecule large enough to accommodate a solvent molecule with a radius of 1.4 Å, although it appears to be unfilled in the electron density map. In places the accessible surface is controlled by atoms above or below the section shown. The dashed outline is the surface of N or O atoms, the solid outline C or S atoms. (Reprinted with permission from Lee and Richards.²⁵)

Procedure of Lee and Richards

The procedure reported by Lee and Richards²⁵ developed from a program used to graphically portray the van der Waals surface of a protein. For the area calculation the radius chosen for each atom in the structure was the van der Waals radius for that particular atom plus the radius of the hypothetical probe, most often set at 1.4 Å. The structure was then sectioned by a series of planes perpendicular to one of the principal axes. The intersections of the enlarged atom spheres with this plane gave a set of circles of varying size. The outer arcs defined by the intersections of these circles represented the trace of the accessible surface of the protein on that plane (Fig. 8). Some internal surface appeared on occasion, and represented cavities in the structure that were large enough to hold one or more probe spheres. Such cavities were recognized by hand inspection of the lists of surface arcs. The total length of the trace of the accessible surface multiplied by the spacing of the planes gave an approximation to the area of the surface associated with that plane. The sum of such surface

²⁵ B. Lee and F. M. Richards, *J. Mol. Biol.* **55**, 379 (1971).

increments over the whole set of planes provided the total accessible area of the molecule. This number approached a limiting value as the spacing of the planes was decreased. A practical balance of computing time against numerical accuracy suggested a spacing of between 0.25 and 0.5 Å as being appropriate.

Richmond¹⁴ modified this program so that accessible areas are calculated for each atom separately. Any atom list could thus be processed without calculation over the whole protein each time. An example of a section through one atom and its overlapping neighbors is shown in Fig. 9. Normally, only the accessible area of the atom is recorded as a single scalar quantity. However, the segment lists have much more information potentially available. The shape of each accessible patch on the atom can be visualized along with its vector directions with respect to the coordi-

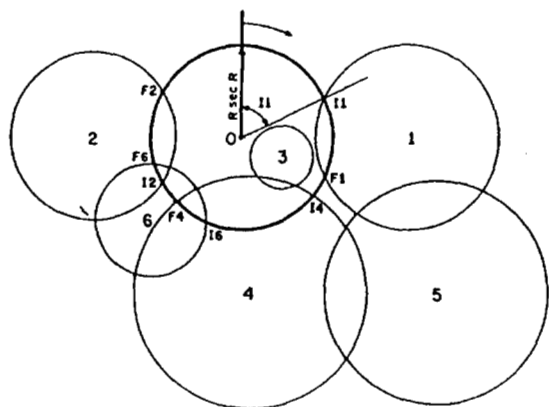


FIG. 9. The Lee-Richmond procedure for calculating the accessible area of an individual atom. The heavy circle is the intersection of the sphere of the target atom (0) with a particular plane. The intersection of the spheres of the neighboring atom, 1 through 6, are shown as lighter circles of differing sizes which are determined by their expanded radii and vertical position with respect to the plane. The intersection check is started at the top as shown by the arrow. For each circle the initial (I) and final (F) intersections (if any) of each circle are accumulated in two lists as their angular position on the 0 circle measured from the starting position. Any atom from the full coordinate list that is possibly close enough to intersect the 0 circle is checked, but the order of checking is arbitrary. The full list of intersections is then ordered on the basis of increasing values of the I intersections. The I and F lists are then tested against each other to get all continuous occluded segments of 0 (in this example I1 to F1 and I4 to F2). Any remaining segments of 0 are part of the accessible surface by definition (i.e., F1 to I4 and F2 to I1). Multiplication of the summed segment lengths by the radius RSECR and by the interplanar spacing gives the approximate accessible area for atom 0 in this plane. Such numbers are summed for the full set of planes to get the total accessible area of 0.

nate axes. An example of such a presentation is shown in Fig. 10. While this format may be useful for visualization, it is likely that the solvent-adjacent polyhedral faces and face normals of the limiting polyhedrons are easier to deal with computationally as suggested by Finney.⁷

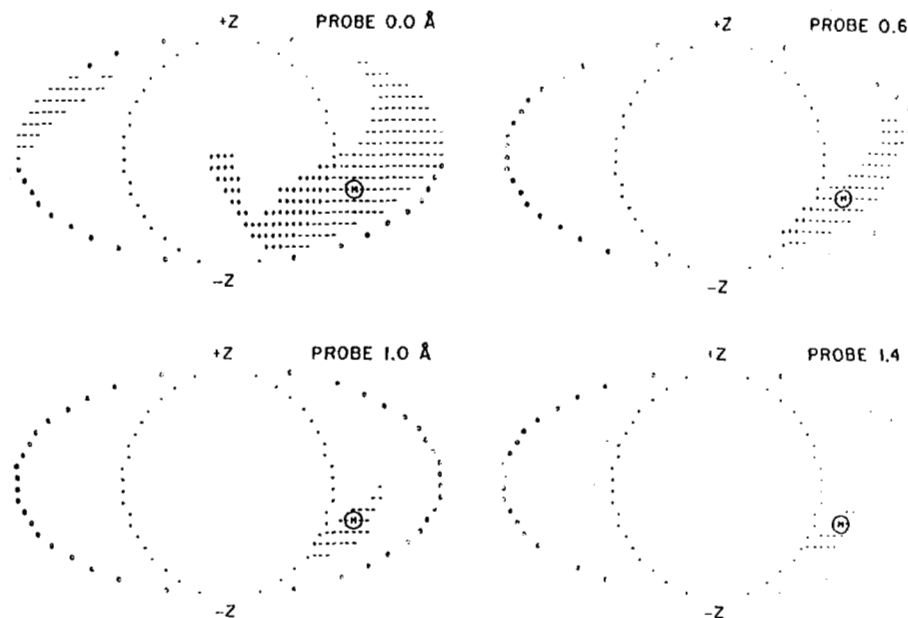


FIG. 10. A modified stereographic projection of the accessible areas on the surface of an atom drawn from segment lists of the type described in Fig. 9. The inner circle of points represents the equator of the sphere, and surrounds the hemisphere above the plane of the paper. The hidden hemisphere is imagined cut along the furthest meridian and then opened up on each side to give the two areas inside the ellipse of asterisks. The right and left lines of asterisks represent the cut meridian. Accessible areas are shown as + on the front hemisphere and - on the back hemisphere. In this example the Z axis of the coordinates is vertical. The +X axis comes toward the viewer. The +Y axis is to the right in the plane of the paper. For illustration the amide nitrogen atom of Phe 120 in ribonuclease S is shown. The H in the small circle gives the N-H vector direction for the amide group (the same direction for each panel). The area patterns are shown for different values of the probe radius. For $r = 0.0 \text{ \AA}$ the indicated area is the actual van der Waals surface of the NH group. The clear space representing nonaccessible area is that part of the NH sphere occluded by covalent attachment to the C and C_α atoms of the main chain and by any van der Waals overlap caused by misplacement of nonbonded atoms in the X-ray structure. The decrease in accessible area with increase in probe size is seen in the other panels. At $r = 1.40 \text{ \AA}$ there is still a small accessible region around the N-H vector which suggests that this particular proton could exchange with solvent with little or no alteration in the structure of the protein.

Procedure of Shrake and Rupley

A different algorithm for calculating the solvent-exposed areas of atoms was developed independently by Shrake and Rupley.²⁶ Again, a sphere of expanded radius equal to the van der Waals radius plus the probe radius (taken as 1.4 Å) is set up around each atom. The central atom, whose area is to be calculated, is represented by a set of 92 points distributed nearly uniformly over the surface of the sphere. Each point is then checked against surrounding atoms to find out if it is within any of the spheres. Points outside the spheres of all surrounding atoms lie on the accessible surface and their number is a direct measure of the accessible area of the central atom. For the occluded points of the central atom, the test atom closest to any particular point is credited with occluding the point. Thus the neighboring atoms collectively provide the environment of the central atom and can be scored quantitatively for their influence on the central atom.

These programs provided actual areas in Å², and these are the values frequently reported. In the original paper Lee and Richards²⁵ also defined the term *accessibility* which is a dimensionless quantity varying between 0 and 1. It represents the ratio of the accessible area in a particular structure to the accessible area of the same group in a reference compound. The latter is normally taken as gly-X-gly, where the group of interest is in the residue X. Accessibility, so defined, is being used at this time, particularly in electrostatic calculations where interactions are modified by these dimensional factors.²⁷

Procedure of Wodak and Janin

Wodak and Janin²⁸ have proposed an approximate analytical expression for the accessible area rather than the numerical calculation described above. The equations are differentiable and can be used directly as a factor to incorporate solvent influences in energy minimization procedures. The derivation assumes a random distribution of spheres surrounding the target atom and includes a correction for excluded volumes. Although the expression is not accurate for a specific atom, averages taken over all, or large parts, of a structure become very good approximations of total surface area. The original paper should be consulted for the derivation of the following equations, where r_o = van der Waals radius of the

²⁶ A. Shrake and J. A. Rupley, *J. Mol. Biol.* **79**, 351 (1973).

²⁷ J. B. Matthew, G. I. H. Hanania, and F. R. N. Gurd, *Biochemistry* **18**, 1919 (1979).

²⁸ S. J. Wodak and J. Janin, *Proc. Natl. Acad. Sci. U.S.A.* **77**, 1736 (1980).

target atom; r_i = van der Waals radius of a neighboring atom; r_w = radius of the spherical probe; d_i = interatomic distance between atoms o and i ; b_i' = maximum area of target atom covered by atom i ; b_i = minimum area of target atom covered by atom i ; n = number of atoms which occlude any area surrounding the target atom; S = area of expanded target atom = $4\pi(r_o + r_w)^2$; A = approximate value of accessible area of target atom.

$$b = \pi(r_o + r_w)(r_o + r_i + 2r_w - d_i)[1 + (r_i - r_o)/d_i] \quad (7)$$

$$b' = \pi(r_o + r_w)(r_o + r_i - d_i)[1 + (r_i - r_o - 2r_w)/d_i] \quad (8)$$

Define

$$A' = S \sum_{i=1}^n [1 - (b_i - b_i')/S] \quad (9)$$

$$B' = \sum_{i=1}^n b_i' \quad (10)$$

Then

$$A_c = A' - B' \quad (A_c = 0 \text{ if } A' < B') \quad (11)$$

$$A_c/d_i = A'/d_i - B'/d_i \quad (12)$$

where

$$A'/d_i = \frac{-A'(b_i/d_i - b_i'/d_i)}{(S - b_i + b_i')} \quad (13)$$

$$B'/d_i = b_i'/d_i \quad (14)$$

The original paper and references therein should be consulted for possible further approximations using single spheres for entire residues and for use of these functions in defining domain structures.

(Note that a mathematically accurate description of accessible areas for any collection of spheres has been developed by T. J. Richmond²³ and M. Connolly.²⁴)

Surface Representations

The computer presentation of van der Waals surfaces in the form of packing models has been highly developed by R. Feldman at the National Institutes of Health. However, the algorithms have not been reported to generate numbers related to the area or volume of these figures.

A particularly effective presentation of the continuous molecular surface of molecules, including both the contact and reentrant sections, has been developed by M. L. Connolly. Interior cavities can be recognized

and enumerated. A brief overview of the computer presentations has been given by Langridge *et al.*²⁹ and Bash *et al.*³⁰ Both areas and volumes are provided in newer programs.^{24,31} Unfortunately, no details of any of these algorithms are available in published form. The latter are particularly important in providing analytical expressions for the area which can be differentiated and built into energy minimization procedures, as described in detail by Richmond.²³

²⁹ R. Langridge, T. E. Ferrin, I. D. Kuntz, and M. L. Connolly, *Science* **211**, 661 (1981).

³⁰ P. A. Bash, N. Pattabiraman, C. Huang, T. E. Ferrin, and R. Langridge, *Science* **222**, 1325 (1983).

³¹ M. L. Connolly, *Science* **221**, 709 (1983).