

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63

# JMB

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®



## Calculation of Standard Atomic Volumes for RNA Cores and Comparison with Proteins: RNA is packed more tightly than Protein

N. R. Voss\* and M. B. Gerstein

*Molecular Biophysics and Biochemistry, Yale University  
260 Whitney Ave, P.O. Box  
208114, New Haven, CT 06520  
USA*

Traditionally, for biomolecular packing calculations research has focused on proteins. Besides proteins, RNA is the other large biomolecule that has tertiary structure interactions and complex packing. No one has yet quantitatively investigated RNA packing nor compared its packing to that of proteins because, until recently, there were no large RNA structures. Here we address this question in detail, using Voronoi volume calculations on a set of high-resolution RNA crystal structures. We do a careful parameterization, taking into account many factors such as atomic radii, crystal packing, structural complexity, solvent, and associated protein to obtain a self-consistent, universal set of volumes that can be applied to both RNA and protein. We report this set of volumes, which we call the NucProt parameter set. Our measured values are consistent across the many different RNA structures and packing environments. However, our volumes are only defined on well-packed atoms, those with sufficient packing neighbors, that typically occur on the interior of RNA molecular and not the unbounded atoms on the surface. When common atom types are compared between proteins and RNA, nine of 12 types show that RNA has a smaller volume and packs more tightly than protein, suggesting that close-packing may be as important for the folding of RNAs as it for proteins. Moreover, calculated partial specific volumes show that RNA bases pack more densely than corresponding aromatic residues from proteins. Finally, we find that RNA bases have similar packing volumes to DNA bases, despite the absence of tertiary contacts in DNA. Programs, parameter sets and raw data are available online at <http://geometry.molmovdb.org>

© 2004 Published by Elsevier Ltd.

\*Corresponding author

Keywords: packing density; RNA volumes; RNA packing

### Introduction

Numerous methods have been developed to determine atom radii and volumes for proteins<sup>1–11</sup> and have been applied to DNA.<sup>12</sup> These radii and volumes are necessary in understanding protein structure and particularly for uncovering the relationship between packing and stability. Many studies requiring accurate protein radii and volumes have characterized a number of protein properties including: protein energies,<sup>13</sup> protein–protein interactions,<sup>14</sup> standard residue volumes,<sup>5</sup>

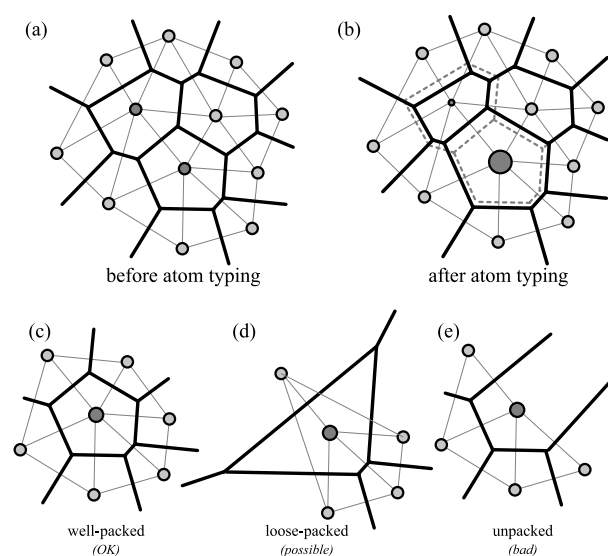
internal core packing,<sup>15,16</sup> packing at the water interface,<sup>17,18</sup> protein cavities,<sup>7,8,19</sup> the quality of crystal structures,<sup>20</sup> analysis of volume by amino acid composition,<sup>21,22</sup> macromolecular motions,<sup>23,24</sup> and even measurements of the fit between enzyme and substrate.<sup>25,26</sup> Standard volumes and radii are also important in an indirect sense for the prediction of side-chain packing.<sup>27–29</sup>

Although a standard protein volume set has been available for years<sup>1</sup> and a DNA volume set was produced recently,<sup>12</sup> no attempts have been made to obtain a standard volume set for RNA molecules. This has been primarily due to a lack of RNA structures other than tRNAs and small oligonucleotides, because their crystallization was once thought to be difficult. Within the past decade, it has been shown that RNA structures can be crystallized in the same way as proteins. This has created a new

Abbreviations used: VDW, van der Waals; PSV, partial specific volume; A-DNA, A-form DNA.

E-mail address of the corresponding author: neil.voss@yale.edu

64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126



**Figure 1.** Voronoi constructs and problems. Effect of atom typing on atom volume. (a) Two-dimensional example of the Voronoi construction. Planes are drawn equidistant between any two atoms. The planes are then intersected to get a volume. (b) For atoms of different sizes the planes are no longer placed equidistant between the atoms, but rather as a ratio function of the van der Waals radius of the atoms. So, large atoms are assigned a larger volume and small atoms are assigned a smaller volume. Three major types of Voronoi packing. (c) Well-packed: polyhedron is closed and surface falls under cutoff value. (d) Loose-packed: polyhedron is closed, but due to lack of neighbors the polyhedron has a large surface area above the cutoff value. (e) Unpacked: Voronoi polyhedron is open and no volume can be calculated. Only well-packed are used to determine the volumes of the atoms.

emphasis on solving RNA structures including: ribosomes, self-splicing introns, and many others. Now that there are several structures available, RNA packing can be addressed and analyzed.

To calculate volumes, we employ the traditional Voronoi polyhedra method.<sup>30</sup> In 1908, Voronoi found a way of partitioning all space amongst a collection of points using specially constructed polyhedra. Here we refer to a collection of “atom centers” rather than “points.” Bernal & Finney<sup>31</sup> first applied this method to molecular systems and Richards<sup>3</sup> first used it with proteins. The methods used in this work have been previously described by others,<sup>3,31</sup> as well as in our earlier work.<sup>9–11</sup> Figure 1 shows how a Voronoi polyhedron is constructed. This construct partitions space such that all points within a polyhedron are closer to the atom defining the polyhedron than to any other atom. The Voronoi planes are shifted from the original equidistant planes (Figure 1(a)) to the modified set (Figure 1(b)) determined by the relative sizes of the van der Waals (VDW) radii of the atoms, i.e., bigger atoms take up more space in the Voronoi construct than smaller ones. Only atoms whose volumes are well-defined (Figure

1(c)) and not loosely packed (Figure 1(d)) or unpacked (Figure 1(e)) are included.<sup>11</sup> Unpacked and loosely packed atoms usually consist of surface atoms or atoms near cavities and therefore do not have enough neighbors to pack tightly. The Voronoi method provides a good estimate of the true volume of an atom and in turn, reliable, self-consistent values for the comparison of atom volumes. Atoms are assigned VDW radii based on their atom type. The typing follows standard united atom conventions and chemical atom typing. A new technique applied in this study is used to test the contributions of crystal symmetry to surface atoms. Since we are interested in large RNA complexes, and the RNA molecules are typically in close contact within the crystal form, it naturally follows to use this additional packing to our advantage as long as there is no effect on the final numbers.

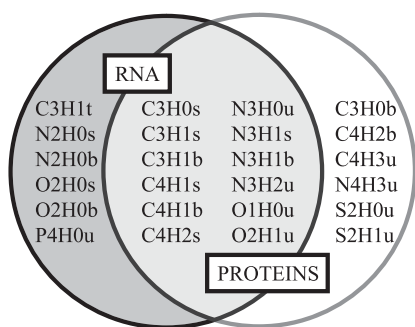
Using Voronoi polyhedra, we report the standard volumes (and many other statistics) for all RNA atoms (49 in total) and all four nucleotides. These atoms are arranged into 18 atom type volumes and radii based on the chemical structure. Further, the benefit of using crystal symmetry to increase the size of the data set is presented. Crystal symmetry had no effect on the final volumes, but increased the population of atoms in our set. Also, we locate less defined atoms within packed RNA structures, such as the backbone has a low percentage of well-packed polyhedra and is, therefore, less defined. We measure the dependence of the volume of the nucleotides on different RNA structural categories (e.g. tRNA, small rRNA, or ribosomes). The final RNA nucleotide volumes are then compared against DNA and organic molecule measurements. In order to evaluate the role of water, ions and proteins in RNA packing, we remove solvent and protein atoms from the calculations and look at the results. This had no effect on the final volumes, but the number of well-packed RNA atoms decreases significantly. We also compare proteins to RNA by comparing the atom types they have in common. The atom types of protein are found to run slightly larger than those of RNA. In addition, from these volumes, we can calculate the partial specific volume of the RNA nucleotides and find that RNA packs more densely than protein.

## Formalism and Results

### Atomic radii calculations

#### Nomenclature

The atomic groups for the RNA atoms are given a nomenclature of the general form “ $X_nH_mS$ ”, where  $X$  indicates the chemical symbol;  $n$ , the number of bonds, which, in most cases, is equivalent to saying  $sp$ ,  $sp^2$ , or  $sp^3$  orbitals;  $H_m$ , the number ( $m$ ) of hydrogen ( $H$ ) atoms attached to the atom where the  $H$  does not change and acts a label for the number ( $m$ ); and  $S$  the subclassification for the atom type,



**Figure 2.** Venn diagram of protein and RNA atom types. Diagramed are all atom types involved in RNA and protein. Types on the right are only involved in RNA while types on the far left are only involved in proteins leaving the central 12 types existing in both RNA and protein.

which is one of the following symbols: *b* (big), *s* (small), *t* (tiny) or *u* (unique). When there are no subclasses for the atom type; *u* (unique) is used. When the atom type needs to be divided into two separate sub-types the type with the larger volume is designated as *b* (big) and the smaller volume *s* (small). In one case (*C3H1*), the atom type requires the addition of a new classification from the previous two subclasses defined previously in proteins.<sup>11</sup> Since the new subclass is smaller than both of the existing *b* (big) and *s* (small) subclasses, its subclass is designated as *t* (tiny). Figure 2 summarizes the 24 different atom types involved in RNA as well as protein and shows which types are common to both.

**Voronoi plane positioning method**

Voronoi polyhedra were originally developed by Voronoi nearly a century ago.<sup>30</sup> While the Voronoi construction is based on partitioning space amongst a collection of “equal” points, all protein atoms are not equal. Some are clearly larger than others. In 1974, a solution was found to this problem,<sup>3</sup> and since then Voronoi polyhedra have been applied to proteins and DNA. Two principal methods of re-positioning the dividing plane have been proposed to make the partition more physically reasonable: method B<sup>3</sup> and the radical plane method.<sup>32</sup> Both methods depend on the radii of the atoms in contact and the distance between the atoms (Figure 1(b)). The simplified method B (or ratio method) divides the plane between the two atoms proportionately according to their covalent radii:

$$d = R + (D - R - r)/2 \tag{1}$$

where *d* is the distance from the atom to the plane, *R*, the VDW radius of the atom, *r*, the VDW of the neighboring atom and *D* is the distance between the two atoms. This method was accepted for a long time, but it was determined that it had a particular flaw. The flaw is vertex error, where the planes

created by neighboring atoms do not perfectly intersect at precise vertices. Vertex becomes a major problem when working with spheres of dramatically different radii. Then the radical plane was introduced which uses a particular quadratic equation to properly divide up the space to obtain precise vertices:

$$d = (D^2 + R^2 - r^2)/2D \tag{2}$$

Because it creates perfect polyhedra, the radical plane method is more pure geometrically than method B. These precise vertices are required for space dividing constructions such as Delaunay triangulations<sup>33</sup> and alpha shapes.<sup>8</sup>

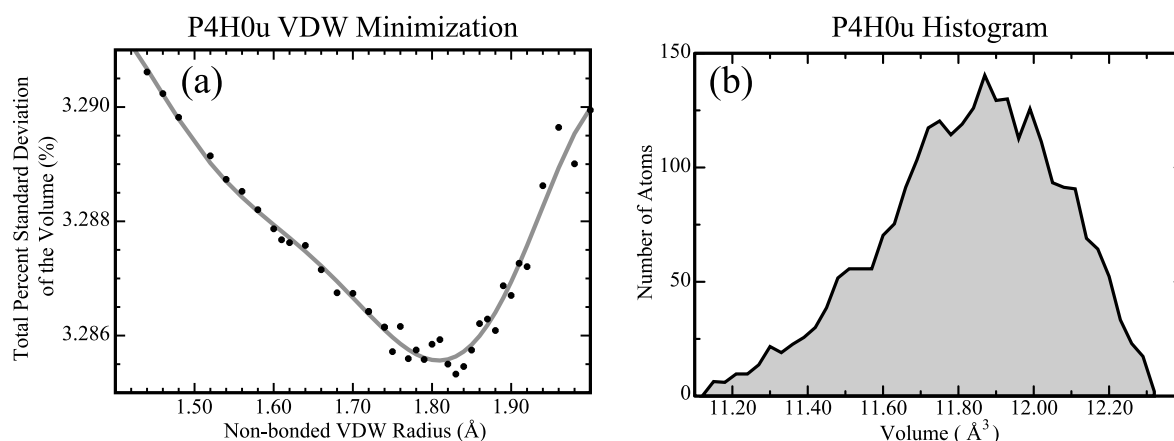
In particular, when comparing the two methods in terms of final volumes there is little difference between the two methods. Even though method B suffers from vertex error, it has been shown to be quite robust for protein calculations, even more robust than the radical plane.<sup>10</sup> In particular, there are two main issues where the methods differ: vertex error and self-consistency. For arbitrary systems with radii of significantly different values, vertex becomes a major issue and the method B is no longer a reasonable approach. However, the radii of proteins atoms do not differ that much and it has been shown that vertex error accounts for one part in 500.<sup>17</sup> In addition, method B has shown to give more self-consistent volumes. It was revealed that the radical plane method actually results in a higher standard deviation than method B, suggesting that it places the plane in a less consistent manner.<sup>10</sup> Further, method B is has been thoroughly tested over the years, while the radical plane is a more recent approach. In addition, the current standard volume set in proteins uses method B for its calculations, so to make direct comparison we will need to have an RNA volume set under the same methodology.

While method B suffers from vertex error, it was reported that this only accounts for one part in of the total volume primarily due to having radii.<sup>17</sup> There are also a two caveats associated with the radical plane method. First, all prior Voronoi research in proteins is based on the method B technique, therefore using radical planes for RNA makes it difficult to draw parallels to protein. Second, volumes calculated by the radical plane result in overall higher standard deviations.<sup>10</sup> Furthermore, in this study, the average standard deviation of the atom types rises from 1.24 to 1.32 for the radical plane method.

Despite this, we report the base volumes for both methods (with little difference) but we use the more traditional method B in all figures, comparisons to protein and radii refinement. The raw data sets and histograms for both methods are also available on the web.

**Importance of atom typing**

Described in more detail by Tsai *et al.*,<sup>11</sup> the distance between the atoms and their intersecting



**Figure 3.** Determining non-bonded VDW radius for the unassigned P4H0u atoms. (a) The normalized standard deviation for the P4H0 atom *versus* its VDW radius. The minimum is found to be 1.82 Å. These values are used for the final Voronoi volume calculations. (b) Histogram of the P4H0u atoms from our final NucProt data set showing one distinct peak.

planes used for Voronoi volume calculation depends on the VDW radius of the atom type. Due to this dependence on atom radius, it becomes increasingly important to obtain accurate atomic classifications and radii. Work done earlier studied the affect of varying the number of atomic classifications and came to the determination that the atom typing system described by *XnHmS* nomenclature was the best balance between over and under fitting for accurate measurements of the volumes for the atoms.<sup>11</sup>

#### VDW radii taken from protein set

The VDW radii for several of the atomic groups involved in RNA structures have analogous atoms in proteins. Several papers have been published on the VDW radius of protein atoms.<sup>5,9</sup> For these overlapping groups, the non-bonded VDW radii of RNA atom groups are simply transferred from their corresponding protein atom groups using the radii defined by the ProtOr set.<sup>9</sup> Whenever there is a small or big designation for the group, the atom group is compared by volume to the protein atoms, e.g., guanine *N1*, of chemical type *N3H1*, is more similar in volume to *N3H1s* than *N3H1b*. Despite vast differences, RNA structure contains only three new atomic groups that completely lack a protein analog, namely *O2H0*, *N2H0* and *P4H0*. *N2H0*, though a new type, is found to be very similar to *N3H1*. Assignment of all the RNA atoms to groups is for the most part straightforward; the only complication came from assignment of the *N2H0* nitrogen atoms and two remaining missing types.

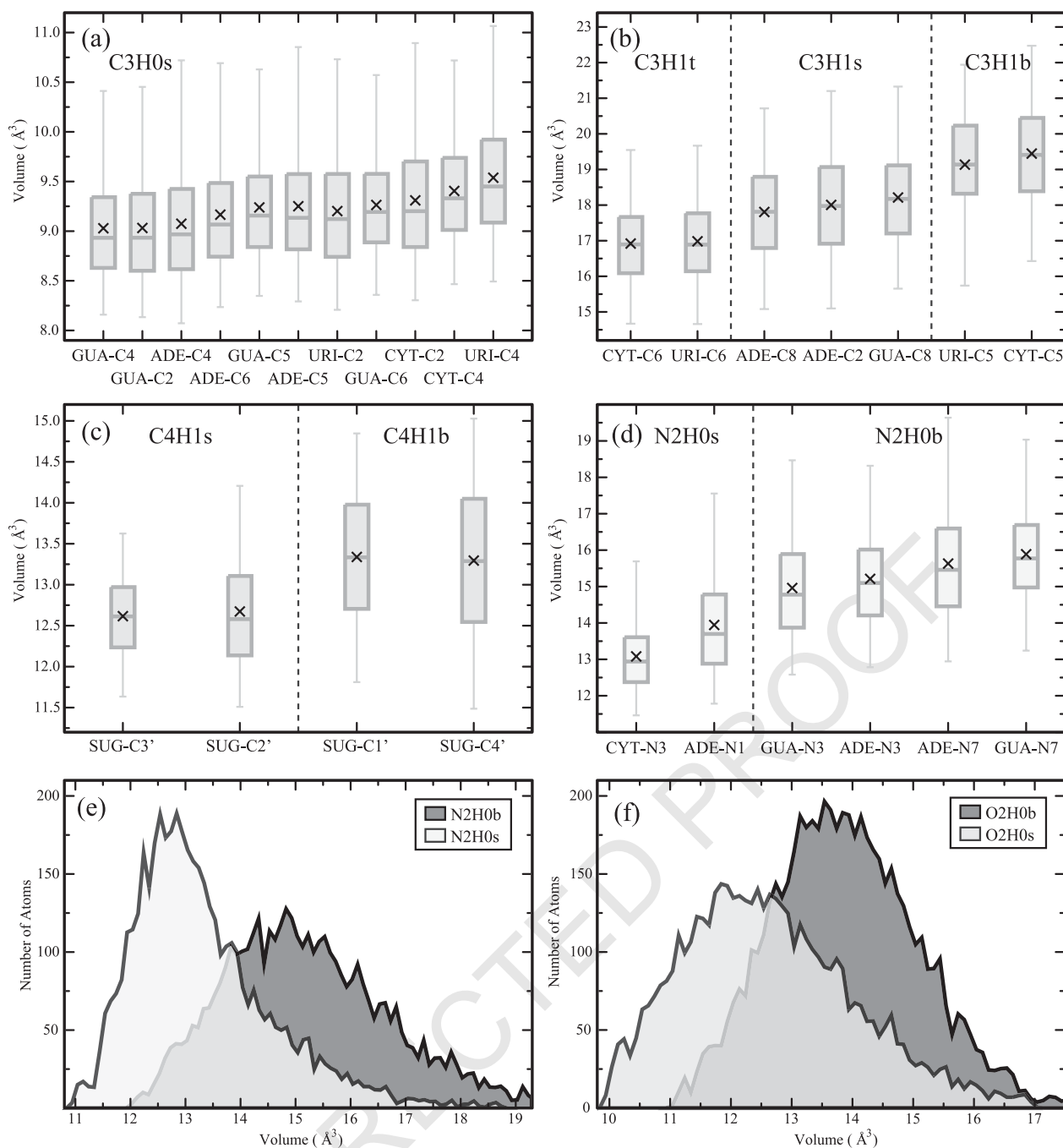
#### Adjusting the bonded VDW radii

Because this new NucProt data set is to include RNA and protein, an investigation into the bonded radii is undertaken to make the values more accurate for both protein and RNA. Using the defined bond length from CNS,<sup>34</sup> the bond radii are

varied for each atom type (grouping small and big subtypes into one type) in order to minimize the sum over all squared bond differences (the bond length—the VDW radius of both atom types bonded) in RNA and protein together. These new bonded radii are not significantly different from the previously published radii,<sup>11</sup> but give a better account of the atom types. For example, the *O1H0u* bonded radius drops the most (from 0.66 to 0.52) reflecting a smaller oxygen atom size due to its double-bonded character, whereas ten of the 24 types change by less than 0.01 Å. These newly adjusted bonded radii should provide a more self-consistent volume data set.

#### New atom types for RNA

Next we need to determine the modal behavior for the unassigned atoms, i.e., do they require small and big subgroups or are they a unique type. *P4H0* only contains one RNA atom and cannot be subdivided further unless the phosphorus atom attached to a guanosine is packed differently than a uridine phosphorus, which is not the case. Therefore, *P4H0* is given a *P4H0u* designation. Further, the *P4H0u* atom type produces a tight histogram (Figure 3(b)), confirming its behavior as a unimodal distribution. The *O2H0* and *N2H0* atom types contain three atoms and six atoms, respectively, of which neither follows a simple distribution. The *O2H0* atom consists of the 3', 4' and 5' sugar oxygen atoms. Individual volume calculations show *O4'* is significantly smaller than both of its type-equivalents *O3'* and *O5'*. The histogram of the *O2H0* atoms (Figure 4(f)) shows a bimodal distribution confirming this assessment. Hence, we design two species of *O2H0* atom: a big class, *O2H0b* (*O3'* and *O5'*), and a small class, *O2H0s* (*O4'*). The *N2H0* atom type is more complex and contains six different types: *ADE-N1*, *ADE-N3*, *ADE-N7*, *GUA-N3*, *GUA-N7*, and *CYT-N3*. From Figure 4(d), the *N3* and *N7* atoms from both purines are grouped into a large



**Figure 4.** Distributions of atom type volumes. (a) Distribution of all atoms composing the *C3H0* group showing one distinct volume. (b) Distribution of all atoms composing the *C3H1* group, suggesting three distinct groups: tiny, small and big. (c) Distribution of all atoms composing the *C4H1* group, suggesting two distinct groups: small and big. (d) Distribution of all atoms composing the *N2H0* group, suggesting two distinct groups: small and big. (e) Volume distribution of the two *N2H0* groups, small and big. (f) Volume distribution of the two types of *O2H0* atoms, small and big.

set and while the *CYT-N3* is significantly smaller than all of the other atoms, it is grouped with the only slightly larger *ADE-N1*. After grouping, a good separation between the small and big subgroups is found (Figure 4(e)).

**Determining non-bonded VDW radius of new types**

For the unassigned atom groups (*O2H0*, *N2H0*, and *P4H0*, from above), a non-bonded VDW radius

needs to be determined. The bonded VDW radii were assigned when the bonded radii were adjusted for all the atom types. All nitrogen-containing atom groups (*N3Hx*, *N4Hx*) in the ProtOr set<sup>10</sup> for proteins are defined as having the same bonded and non-bonded atomic radius, so we felt the *N2H0* should have the same values as its sister atom types because its volume is the same as *N3H1* types. The non-bonded VDW radii for the *P4H0* and *O2H0* do not have existing values and so

their non-bonded VDW radii are determined by varying the non-bonded VDW radius of the atom in question and minimizing the sum of the percent standard deviation of volume (standard deviation of the volume divided by the mean volume) over each atom in RNA. The standard deviation gives an unfair bias to minimizing the error of atoms with larger volumes due to their larger deviations, so by taking the standard deviation divided by the mean this bias is reduced. This method for calculating the missing non-bonded VDW radii results in the most self-consistent set of volumes.

As shown in Figure 3(a), the standard deviation of the volume for phosphorus atom, *P4H0* volume gave a convex curve when its radius was varied. The curve is then fit to a tenth-degree polynomial (only to smooth out the noise without loss of generality) and the *P4H0u* radius is taken to be the minimum of the polynomial fit, which is 1.82 Å. This final value gave an extremely tight unimodal distribution (Figure 3(b)). Likewise, a two-dimensional optimization is employed for the *O2H0* types due to its bimodal distribution (Figure 4(e)), by simultaneously varying the radius of both subtypes. The global minimum of the percent standard deviation is determined to be 1.50 Å for the *O2H0s* and 1.62 Å for *O2H0b* (data not shown).

**Determination of volumes**

*Brief description of Voronoi method*

The volumes of the atoms are determined with

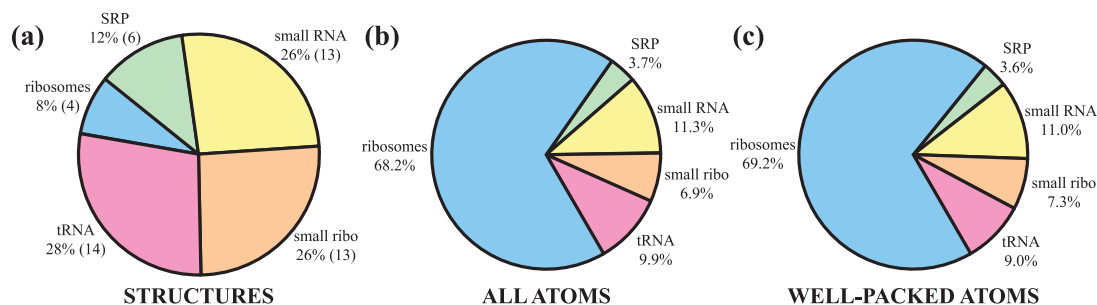
the same Voronoi method as published earlier.<sup>9–11,17</sup> For every pair of atoms, a plane is constructed approximately equidistant from both of the atoms (in actuality the distance is adjusted by the VDW radii of each atom type) and orthogonal to the bond between the atoms (Figure 1). The planes are then intersected, leaving an enclosed polyhedron for each atom. While not every atom has a closed polyhedron, the majority of the polyhedra is closed.

*Assembling the structure set*

Structures were obtained from the NDB<sup>35</sup> by searching for nucleic acids structures containing RNA, with strand lengths greater than 26 nt, to avoid small synthesized RNA molecules, and resolution better than 5 Å. The cutoff value of 5 Å was chosen to include all four ribosomal subunit structures, including both low and high-resolution versions. We compare the high-resolution only data to the entire set and find no difference in the final volumes. After determining our criteria, the search results in a raw structure set of 125 RNA structures. Most of the found structures are redundant (e.g., 50 S ribosomal subunits soaked with various complexes or tRNA with and without synthetases) and some structures contained DNA base-paired with RNA in a complex. After removing the DNA hybrid structures and duplicates, taking care to use the most accurate and detailed structure within each redundant set, a final set is created consisting of 50 unique structures. For comparison purposes, the final sets are broken down into five smaller disjoint

**Table 1.** Summary of structure sets

Set name	Number of PDB files	Number of RNA atoms	%OK no symm	%OK with symm	% of total atoms	% of total "OK" atoms	PDB Ids
<i>Disjoint subsets</i>							
SRP	6	10,137	35.0	37.9	3.7	3.6	1hq1, 1jid, 1lng, 1mfq, 1e8o, 1l9a
Small ribo	13	19,234	36.6	40.3	6.9	7.3	483d, 1msy, 1jbs, 1i6u, 1mms, 1mzp, 1mji, 1dk1, 1g1x, 1qa6, 430d, 364d, 357d
tRNA	14	27,379	33.2	34.8	9.9	9.0	1f7u, 1ehz, 1fir, 1qf6, 1il2, 1h4s, 1b23, 1qtq, 1ser, 1ffy, 1i9v, 1ttt, 1ivs, 2fmt
Small RNA	13	31,438	33.7	37.2	11.3	11.0	1l2x, 437d, 1et4, 1m5o, 1hr2, 1duh, 1cx0, 1kxk, 1f1t, 1l3d, 1hnh, 1kh6, 1nbs
Ribosomes	4	188,911	38.9	39.0	68.2	69.2	1jj2, 1i94, 1n32, 1nkw
All	50	277,099	37.4	38.4	100.0	100.0	1b23, 1cx0, 1dk1, 1duh, 1e8o, 1ehz, 1et4, 1f1t, 1f7u, 1ffy, 1fir, 1g1x, 1h4s, 1hnh, 1hq1, 1hr2, 1i6u, 1i94, 1i9v, 1il2, 1ivs, 1jbs, 1jid, 1jj2, 1kh6, 1kxk, 1l2x, 1l3d, 1l9a, 1lng, 1m5o, 1mfq, 1mji, 1mms, 1msy, 1mzp, 1n32, 1nbs, 1nkw, 1qa6, 1qf6, 1qtq, 1ser, 1ttt, 2fmt, 357d, 364d, 430d, 437d, 483d
<i>Additional sets</i>							
hi-res	9	11,281	49.4	58.7	4.1	6.2	1jbs, 1msy, 483d, 1et4, 437d, 1l2x, 1jid, 1hq1, 1f7u, 1ehz
RNA only	19	33,782	33.7	38.6	12.2	12.2	1ehz, 1fir, 1i9v, 1l2x, 437d, 1et4, 1hr2, 1duh, 1kxk, 1f1t, 1l3d, 1hnh, 1kh6, 1nbs, 483d, 1msy, 430d, 364d, 357d



**Figure 5.** Distribution of atoms within PDB set. Pie charts showing how the structure set breaks down into the five major categories. (a) Number of structures for each subset. (b) Number of atoms for each subset. (c) Number of sufficiently packed or “OK” atoms for each subset. Though the four ribosomal structures account for only 8% of the structures of the pdb set, they account for 69.2% of the atoms used in the final calculations. In the text it is shown that ribosomal and non-ribosomal RNA have the same final base volumes.

subsets: (i) *SRP*—RNA structures involved with the Signal Recognition Particle, (ii) *small-ribo*—small ribosomal RNA fragments, such as the 5 S rRNA structures, (iii) *tRNA*—transfer RNA with and without synthetases, (iv) *small-RNA*—the other remaining small RNA molecules including ribozymes and self-splicing introns, and (v) *ribosomes*—complete ribosomal subunits (Table 1). These structure sets are unfortunately heavily weighted towards ribosomal data (69% of atoms), because of their immense size, despite only being only four of the 50 structures in the set (Figure 5, Table 1). This effect will be addressed later.

**Generation of final volume set**

Surface atoms (as well as loosely packed interior atoms) sometimes lack closed polyhedra or have extended polyhedra and give rise to indeterminate or inflated volumes, respectively. These two special cases of atoms need to be removed from the set of atom volumes in order to obtain a self-consistent data set. The first case of loosely packed atoms occurs when the Voronoi shell is heavily extended (Figure 1(d)). Loosely packed atoms are distinguished from well-packed atoms by their surface area. Atoms above a certain surface area cutoff are characterized as “possible,” meaning that they have a volume, but it is unsure whether it is

relevant. The loosely packed atoms are not used in our final NucProt data set due to their indefinite character. The second case, an atom having insufficient neighbors, leaves the Voronoi shell open ended and produces an indeterminate volume (Figure 1(e)). These unclosed polyhedra are easy to identify for they have no volume and are designated “bad” by the software.<sup>11</sup> Consequently, only atoms with closed polyhedra and small surface areas are then labeled as “ok” atoms. It should also be noted that all atoms (protein, RNA, ions, water, organic molecules and also modified nucleotides and amino acid residues) within the PDB file are taken into account for Voronoi plane positioning. Unfortunately, modified bases volumes are not reported due their small population within our set, thus making it almost impossible to provide any reasonable statistics. There are only 29 pseudouridines within our PDB set and given that at most half the atoms are well-packed, it would be a very unreliable volume for general use.

Despite applying these standard methods for generating the final volume set, extreme atom volumes existed for each RNA atom. Therefore, as an additional measure, the extreme atom volumes are removed from the ends of each RNA atom distribution such that the average range for each distribution drops in half. Dropping the distribution range in half is chosen, because it provides

**Table 2.** Summary of effect of other atoms on the packing calculations

	50 S Ribosomal subunit (1jj2)				
Crystal symm	+	–	–	–	–
Protein	+	+	–	+	–
Ions/water	+	+	+	–	–
Base volumes (Å <sup>3</sup> )					
GUA	145.9	145.9	145.7	146.4	146.2
ADE	140.0	140.0	139.4	138.9	138.2
CYT	115.5	115.5	115.3	115.6	115.0
URI	110.8	110.8	110.6	110.9	110.2
SUG	176.1	176.1	175.4	179.2	177.4
Addition information					
Count	33,245	33,204	28,996	23,007	20,409
%OK	54.0	53.9	47.1	37.3	33.1
%Closed	93.5	93.1	91.2	90.7	84.9
Mean %SD	6.89	6.89	6.78	7.10	6.55

**Table 3.** Base volumes across several different structure sets

Base Volume ( $\text{\AA}^3$ )	Standard Disjoint Sets					Additional Sets			ALL SETS	radical plane	Other Published Sets		
	SRP	ribosome	small ribo	small RNA	tRNA	hi-res	allma	nosymm	complete		organic <sup>a</sup>	A-DNA <sup>c</sup>	B-DNA <sup>c</sup>
<b>GUA</b>	141.3	146.4	144.0	144.7	142.7	144.2	144.0	146.1	145.9	144.5	157.7	145.4	143.8
<b>ADE</b>	134.1	139.7	136.4	137.9	136.2	138.7	137.3	139.4	139.2	137.8	148.8	138.3	136.1
<b>CYT</b>	110.2	115.6	111.8	113.4	111.6	113.2	113.1	115.2	115.0	113.2	122.1	113.9	113.2
<b>URI</b>	102.2	111.3	107.4	106.1	107.9	106.6	106.9	111.0	110.8	109.4	119.2	132.9 <sup>d</sup>	132.6 <sup>d</sup>
Backbone Vol ( $\text{\AA}^3$ )													
<b>SUG</b>	172.9	176.0	175.8	175.2	174.5	176.8	175.7	175.9	176.1	179.3	133.6 <sup>b</sup>	181.8	174.8

<sup>a</sup>Values converted from the work done by Lee & Chalikian.<sup>62</sup> Volumes require conversion of units from cm<sup>3</sup>/mol to  $\text{\AA}^3$ /residue. <sup>b</sup>This value is from a nucleoside, not a nucleotide, and lacks a phosphate group with a volume of approximately 43  $\text{\AA}^3$  (depending on oxidation state). Calculated sugar volume is averaged over three base volumes subtracted from the nucleoside volumes. <sup>c</sup>Values taken from the work done by Nadassy *et al.*<sup>12</sup> <sup>d</sup>Thymine values are used in place of uracil. Thymine should be approximately 27  $\text{\AA}^3$  greater in volume than uracil, based on atom type volumes in thymine.

the best balance between data loss and reduction of the range. This method is very effective because to drop the range in half only 1.25% of the data is removed from the set. In essence, the central 97.5% of the data has half the range of the complete set of data, highlighting some of the extreme values resulting from over-packed atoms due to structural overlap errors or loosely packed atoms missed by the surface area cutoff. After applying these methods, only well-packed atoms are then used in final volume calculations for self-consistency, making it important to maximize the number of well-packed atoms for a good sample size.

**Effect of surface molecules**

The treatment of surface atoms plays an important role in calculating Voronoi volumes because Voronoi volumes rely on neighboring atoms to create polyhedra surrounding each atom. By increasing the number of neighboring atoms it is possible to have more well-packed atoms. To explore these problematic surface atoms, the effects of crystal symmetry, bound proteins, and solvent atoms on structure are examined for their significance. All three factors had little effect on the final volumes, but all make a significant contribution to the number of observations for each atom (Table 2). In the final volume set, atoms from both high and low-resolution structures, both protein containing and protein free RNA structures, and only crystal symmetry generated structures are integrated into our final NucProt data set.

**Final volumes**

We now can provide final volumes. Since we are actually calculating distributions of volumes, i.e. the probability of a volume given an atom type, we provide both histograms and mean values. In particular, we show a sample distributions in Figures 4 and 8†. It is also useful to have explicit mean values for the volumes. The final volumes

† The rest is available on <http://geometry.molmovdb.org/NucProt>

for the RNA bases are 145.9  $\text{\AA}^3$  for guanine, 139.2  $\text{\AA}^3$  for adenine, 115.0  $\text{\AA}^3$  for cytosine, and 110.8  $\text{\AA}^3$  for uracil (Table 3). All four RNA sugar backbones are approximately the same size and so we report only one value of 176.1  $\text{\AA}^3$  (Table 3). The nucleotide volumes are 322.6  $\text{\AA}^3$  for guanosine, 315.0  $\text{\AA}^3$  for adenosine, 290.7  $\text{\AA}^3$  for cytosine, and 285.5  $\text{\AA}^3$  for uridine.

**Discussion**

**Effects on calculations**

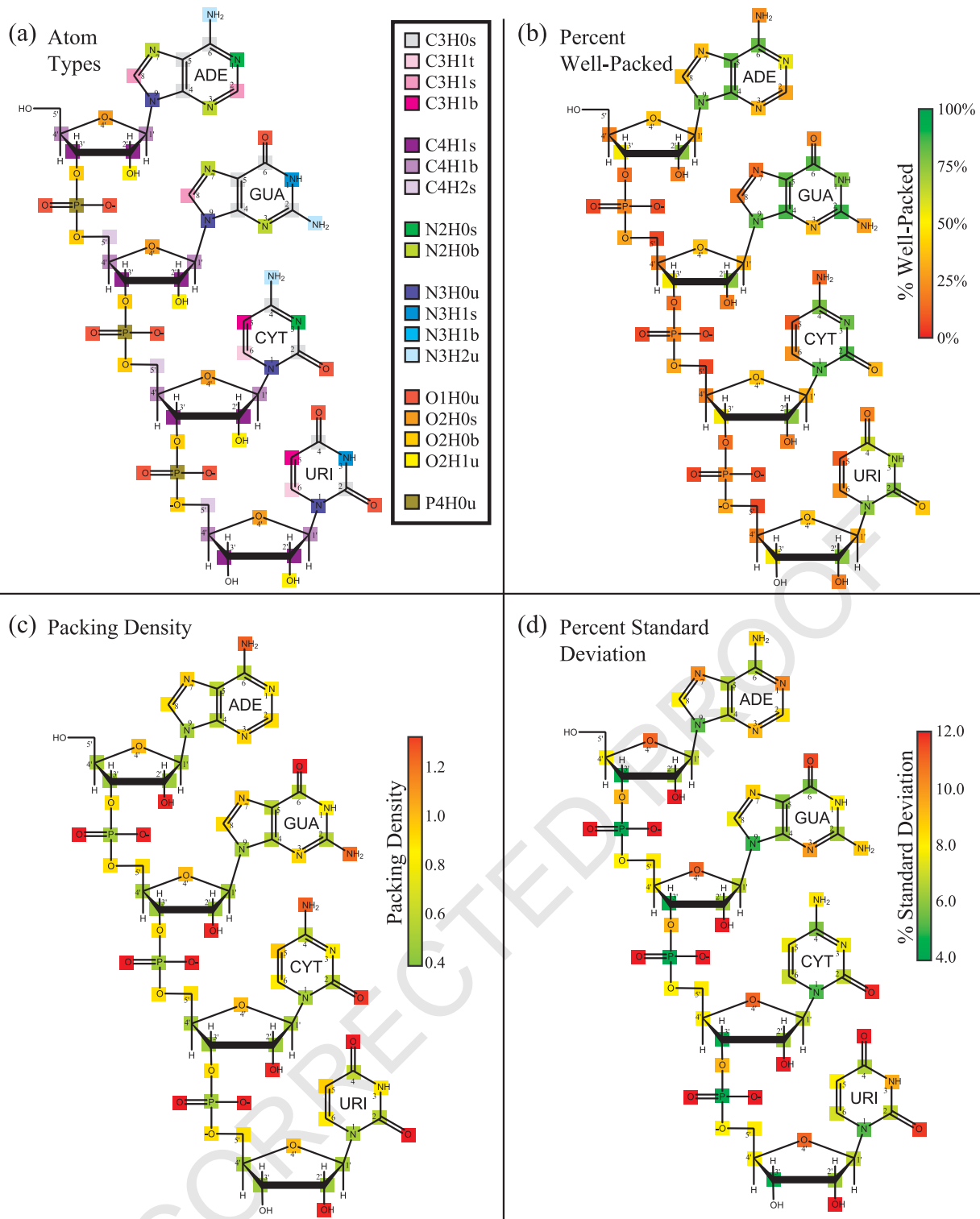
*RNA backbone and base packing*

From the standard deviations and packing percentages, we are able to locate areas within RNA structure that are not well-packed or less defined. In our NucProt data set, only 26.8% of the sugar-phosphate backbone atoms are packed sufficiently to make a volume measurement, which is 20.2% less than the worst base, uracil (at 47.0% well-packed). Further, several backbone atoms have low percentages of well-packed Voronoi polyhedra (Figure 6(b)) and high standard deviations (Figure 6(d)). These results suggest that atoms located in the bases benefit from the tight ring structure of purines and pyrimidines providing inherent packing neighbors as well as base-pairing and base-stacking interactions common in RNA structure. In addition, the atoms located in major groove edge of the RNA bases also have high standard deviations, high packing densities and low percentages of well-packed atoms (Figure 6). This presumably implies a less packed major groove in RNA structures, but it is more likely due to no inherent neighbors in an A-form helix. Therefore, our results suggest that the backbone-sugar regions and major groove atoms are less packed than the interior base atoms.

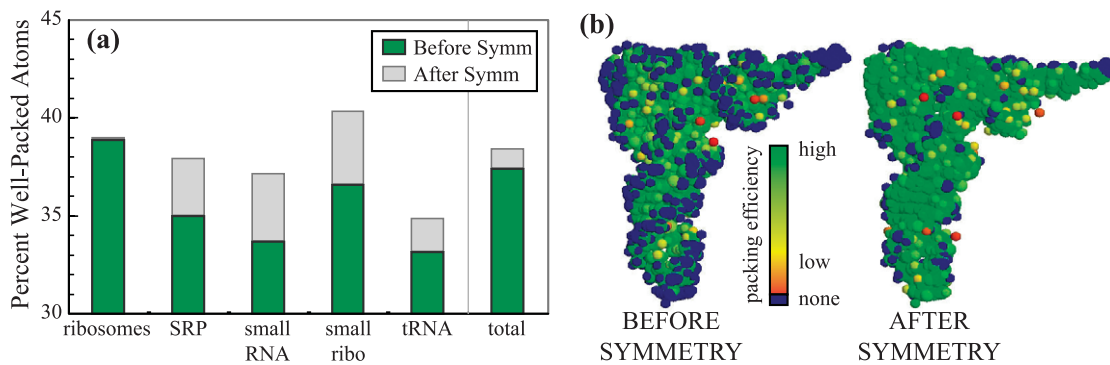
*Role of crystal symmetry in volume size*

Small RNA structures, in general, consist of a single helix and hence lack helix-helix packing, resulting in poorly packed backbones. In an effort to





**Figure 6.** Graphical display of different atomic packing measurements. For each atom in RNA. (a) The atom typing is shown to show what decisions are made to classify the various atoms. (b) The percent well-packed atoms shows that the backbone and extensions off the rings are in general less defined. (c) The packing density (Voronoi volume divided by VDW volume) measures the how tightly each atoms packs. This number tends to be biased by the number of hydrogen atoms bonded to an atom, but still provides insight as another measures of packing. (d) The percent standard deviation of the volume (standard deviation of the volume divided by the mean volume) highlights the unbiased error involved in the volume measurements.



**Figure 7.** Effects of crystal symmetry. (a) Effect of the crystal symmetry on each subset of structures. Ribosomes saw little to no effect, while small RNA molecules, including ribozymes and other small RNAs, see a large jump in their percentage of well-packed atoms. (b) Example of one structure (1ehz), which is by no means the best, where crystal packing helps increase the number of well-packed atoms. Shown is the packing efficiency, i.e., the Voronoi volume of an individual atom divided by the mean volume for the atom. Blue represents atoms that have unclosed Voronoi polyhedra. Packing before crystal symmetry is shown on the left and after is on the right. You can see the dramatic effect of crystal symmetry on obtaining information for surface atoms. The final volumes show no difference between data sets with and without crystal symmetry.

prevent this single helix dilemma, we need to utilize the crystal symmetry contained in the PDB file. Crystal symmetry neighbors have additional relevant packing interactions from their presence within the crystal. Alas, any software found for generating crystal symmetric neighbors is not applicable or does not work well for our purpose. All are incapable of outputting the information to a file or do not have the facility to generate all symmetry neighbors within a given distance of the target structure. Fortunately, matrix information on crystal symmetric rotations and translations are contained within most PDB headers (as well as in the online PDB format description, Appendix 1<sup>36</sup>) and once recognized, it was simple to implement a small script to achieve the additional symmetry neighbors.

As shown in Table 2, crystal symmetry had little to no effect on the final volumes, but does contribute a significantly larger number of acceptable atoms for making calculations (Figure 7). The different disjoint subcategories of structures show that even though the set of ribosomes have little to no effect on the number of well-packed atoms, all other sets containing the smaller structures increase the percentage of well-packed atoms by a dramatic amount. The additional number of well-packed atoms created from the symmetry neighbors not only gives us more data for error analysis, but helps increase the amount of information from atoms involved in backbone packing that would normally have unclosed polyhedra.

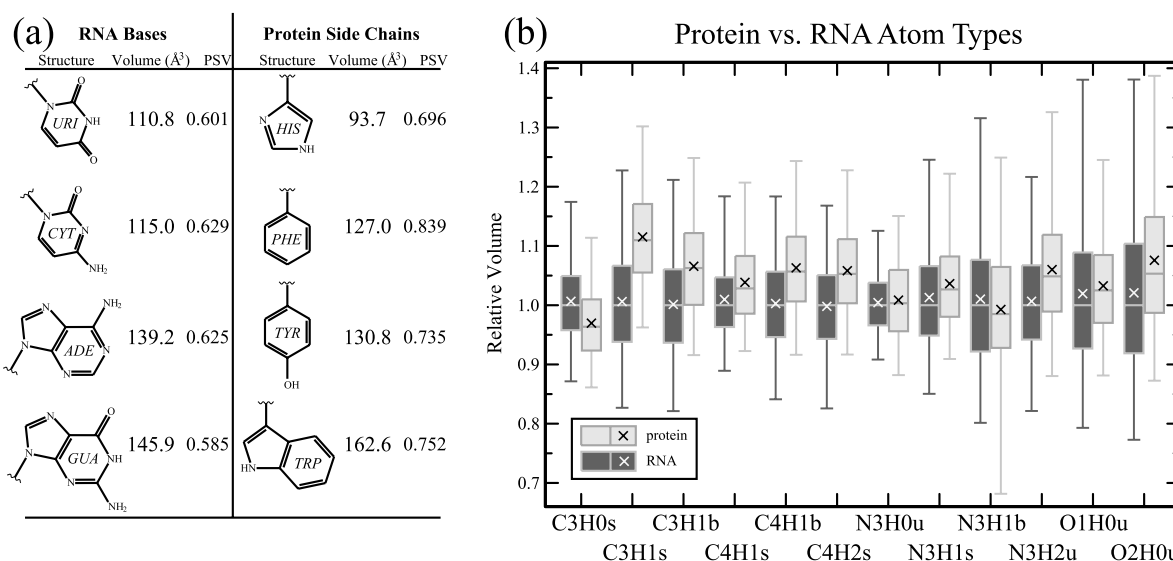
**Roles of different RNA structural categories**

As noted, our data set consists of 69% ribosomal atoms. *A priori* it is unjustified to assume short double-stranded RNAs have the same packing properties as large macromolecular complexes such as the ribosome. The ribosome is large enough in all three dimensions to truly have an interior

while short double-stranded RNA is completely exposed to solvent. Further, we want to confirm that our data set, containing 69% ribosomal atoms, is representative of small RNA molecules as well. There is little difference in nucleotide volume (Table 3) or atom size (data not shown) among the different RNA types. Table 3 clearly shows that the base size, sugar backbone, and entire nucleotide volumes differ by less than 9 Å<sup>3</sup> from the smallest to largest values, which is within the standard deviation. Ribosomal RNA volumes run slightly larger than the other structural categories. This could be due to the larger size of the complexes and their inherent problem of packing helices against other helices. Despite slight variations between the structural categories, the final outcome of the volume calculations suggests that all RNA packs in a universal way.

**Role of water, ions and proteins in RNA structures**

To test the role water, ion, and protein atoms play in RNA structures, we took the largest RNA structure (the refined *Haloarcula marismortui* 50 S ribosomal subunit, 1jj2) and conduct packing tests by systematically removing each kind of atom (Table 3). Crystal symmetry plays a small role in the 50 S subunit, because there is more interior than surface. On the other hand, when the solvent is removed (leaving the RNA and protein), the percentage of well-packed atoms differs by 16.7% from the original value. Similarly, when the protein is removed (leaving the RNA and solvent), the percentage differs by 6.9%. Further, when only the RNA is used the final percentage differs by 20.9%. The final difference of 20.9% is very close to the sum of the other differences of 23.9%, suggesting an independence of the two atom classes. In addition, the solvent has a much larger affect on the loss of well-packed atoms than does the protein. This indicates that RNA atoms pack tightly against the



**Figure 8.** Comparison of RNA volumes to protein. (a) The comparison of aromatic protein residues to RNA bases. On average RNA has a smaller PSV than protein, suggesting RNA packs more densely than protein. (b) Comparing protein to RNA atoms using relative volume. Relative volume is volume divided by the median RNA volume for that atom type. All 12 atoms in the intersection of the protein and RNA atom types are shown for comparison.

solvent. Despite these major differences in the percentages of well-packed atoms, there is no change in the final RNA volumes from the removal of solvent and protein atoms, reinforcing the idea that our set is self-consistent.

### Comparison to proteins and DNA

#### Partial specific volumes

To address the effect of RNA and protein in packing, the partial specific volume (PSV) is computed by taking the calculated volume divided by the atomic mass of the RNA molecule and then changing the units of cubic Ångströms per Dalton (Å<sup>3</sup>/Da) to the classical form of milliliters per gram (ml/g) (using conversion factor of 1 Å<sup>3</sup>/Da =  $N_A \times 10^{-23} = 0.6022$  ml/g). We now provide a new online tool for calculating volumes and PSV for any sequence†.

RNA is found to have an average calculated PSV of 0.569 ml/g, which is significantly more dense than the published protein calculated average of 0.728 ml/g over 13 protein structures.<sup>5</sup> RNA bases are more loosely packed than their complete nucleotide form with an average calculated PSV of 0.610 ml/g and it follows that the sugar and phosphate backbone is more tightly packed with an average calculated PSV of 0.544 ml/g. These values compare well to the published experimental value of 0.540 ml/g.<sup>37</sup> For proteins, it is shown that calculated values are on average 0.5% less than experimental values.<sup>5</sup> Our average PSV is about 5.4% greater than this experimental value, but Durchschlag explains that the experimental RNA PSVs depended heavily on solvent content and the values were difficult to obtain and also may fluctuate greatly.<sup>37</sup> Our results thus conclude that

the calculated PSV using Voronoi volumes for RNA is a good estimate for the experimental PSV.

#### Atom types

There are 18 atom types each in RNA and in protein, but when intersected they only share 12 common atom types. Figure 8(b) analyzes the 12 common types in more detail. Though proteins and RNA for the most part share similar distributions, most protein atom types (nine of 12) run slightly larger (Table 4, Figure 8(b)). The three exceptions to this rule are C3H0s, N3H1b (which run smaller) and N3H0u is almost exactly the same size. The most dramatic effect is shown by C3H1s, which consists of aromatic ring carbon atoms in protein and purine ring carbon atoms in RNA (Table 4). For C3H1s, the 25th percentile of the protein is greater than the 75th percentile of the RNA, suggesting two distinctly different values (Figure 8(b)). One of the reasons why RNA atom types are smaller in volume than equivalent protein types, is their built-in chemical structure. Proteins are chains that have few atoms per residue and pack against one another to achieve tight packing, while RNA contains more than 18 atoms per residue and, therefore has inherent neighboring packing interactions. In fact, the worst packed atoms in RNA are either attached to the phosphorus atom or an extension off the ring structure (e.g. sugar O2', guanine O6, purine O2). In essence, the atom type data shows that RNA is more tightly packed than protein atoms.

#### Similar protein residues

An interesting question to ask is how does the volume of RNA compare to protein amino acid residues of similar chemical structure. Namely, how

1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449

**Table 4.** Summary of atom types in proteins and RNA

Atom type	RNA NucProt Set						Protein NucProt Set (ProtOr)						Comparison, %vol change
	Number	%OK	Count	Mean	Standard deviation	%SD	Number	%OK	Count	Mean	Standard deviation	%SD	
<i>C3H0s</i>	11	81.6	30,479	9.21	0.60	6.5	20	76.1	12,097	8.77	0.63	7.2	-5.0
<i>C3H0b</i>							13	49.4	4418	9.77	0.77	7.9	-
<i>C3H1t</i>	2	24.7	1387	16.95	1.19	7.0							-
<i>C3H1s</i>	3	26.7	2766	17.98	1.44	8.0	8	43.5	1888	20.59	1.81	8.8	12.7
<i>C3H1b</i>	2	10.4	586	19.32	1.51	7.8	8	55.3	2181	21.37	1.90	8.9	9.6
<i>C4H1s</i>	2	64.1	16,554	12.65	0.72	5.7	18	53.2	7227	13.26	1.01	7.6	4.6
<i>C4H1b</i>	2	26.3	6795	13.32	0.97	7.3	6	54.8	3747	14.44	1.33	9.2	7.7
<i>C4H2s</i>	1	6.1	790	21.74	1.77	8.1	20	25.6	4468	23.45	2.34	10.0	7.3
<i>C4H2b</i>							7	27.3	1137	24.42	2.14	8.8	-
<i>C4H3u</i>							9	37.5	3673	36.92	3.25	8.8	-
<i>N2H0s</i>	2	68.7	4461	13.41	1.26	9.4							-
<i>N2H0b</i>	4	29.4	4285	15.33	1.51	9.8							-
<i>N3H0u</i>	4	81.0	10,465	8.79	0.45	5.1	1	74.1	592	8.82	0.66	7.5	0.4
<i>N3H1s</i>	2	74.8	4811	13.63	1.17	8.6	20	62.5	10,356	13.82	1.20	8.7	1.4
<i>N3H1b</i>	4	29.8	4352	15.43	1.67	10.8	4	29.0	500	15.87	2.21	13.9	2.7
<i>N3H2u</i>	3	21.7	2322	22.10	1.94	8.8	4	9.7	286	23.38	2.77	11.8	5.5
<i>N4H3u</i>							1	1.2	12	21.21	1.85	8.7	-
<i>O1H0u</i>	6	13.4	5052	16.29	2.09	12.8	27	36.1	8273	16.17	1.59	9.8	-0.7
<i>O2H0s</i>	1	37.5	4850	12.73	1.40	11.0							-
<i>O2H0b</i>	2	22.3	5753	13.98	1.18	8.4							-
<i>O2H1u</i>	1	19.1	2468	17.39	2.28	13.1	3	20.0	619	18.60	2.45	13.2	6.5
<i>P4H0u</i>	1	20.5	2645	11.86	0.23	1.9							-
<i>S2H0u</i>							2	50.1	280	29.17	2.81	9.6	-
<i>S2H1u</i>							1	51.6	63	34.60	5.73	16.6	-

1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511  
1512

1513 does the volume of an RNA purine compare to  
 1514 tryptophan and how does an RNA pyrimidine  
 1515 compare to phenylalanine, histidine, and tyrosine.  
 1516 Figure 8(a) reports the protein volumes for the side-  
 1517 chains (calculated from the residue volume subtract  
 1518 the volume of glycine) of tryptophan, tyrosine,  
 1519 phenylalanine, and histidine.<sup>17</sup> While these values  
 1520 for the volume are all relatively close to the RNA  
 1521 base volumes (Table 3), the PSV tells a different  
 1522 story (Figure 8(a)). The average PSV for the RNA  
 1523 bases is 0.609 ml/g while the four protein side-  
 1524 chains have an average PSV of 0.755 ml/g,  
 1525 suggesting that the RNA bases are much more  
 1526 dense than protein aromatic side-chains (Figure  
 1527 8(a)). Further, if we divide the total volume by the  
 1528 number of atoms (including the hydrogen atoms),  
 1529 we get an average volume of 9.83 Å<sup>3</sup> per atom for  
 1530 the RNA bases and 11.25 Å<sup>3</sup> per atom for the protein  
 1531 residues. Though the volume per atom numbers are  
 1532 biased by the atom type volumes, this also high-  
 1533 lights that RNA seems to pack more tightly than  
 1534 protein. These results may be due to nucleotide base  
 1535 rings containing more nitrogen atoms than the  
 1536 amino acid aromatic rings. In addition, the RNA  
 1537 rings have more atoms attached to them, creating a  
 1538 large number of inherent neighbors. In addition,  
 1539 RNA duplex base stacking may contribute favor-  
 1540 ably to achieve this tighter packing. Though it is  
 1541 difficult to directly compare these vastly different  
 1542 chemical structures, we find that the RNA bases are  
 1543 more tightly packed than the aromatic protein  
 1544 residues.

#### 1545 DNA volumes

1546 In 2001, Nadassy *et al.* published the standard  
 1547 atomic volumes of double-stranded DNA.<sup>12</sup> Com-  
 1548 paring the volume of RNA in large macromolecular  
 1549 structures to that of A-form DNA (A-DNA) we see a  
 1550 small deviation (Table 3). RNA bases: adenine,  
 1551 guanine and cytosine are larger by only 0.9 Å<sup>3</sup>,  
 1552 0.5 Å<sup>3</sup>, and 1.1 Å<sup>3</sup>, respectively, to that of A-DNA.  
 1553 Since we cannot directly compare uracil to DNA,  
 1554 we compare its volume to the volume for thymidine  
 1555 and they are within the expected difference of 27 Å<sup>3</sup>,  
 1556 due to the extra methyl group. We found that in  
 1557 RNA structures about half of the base volumes are  
 1558 within a standard deviation of the A-DNA base  
 1559 volumes (Table 3), suggesting similar packing of  
 1560 the bases. The sugar-phosphate backbone on the  
 1561 other hand reports a slightly larger difference.  
 1562 The A-DNA sugar plus phosphate reported by  
 1563 Nadassy *et al.* is 5.7 Å<sup>3</sup> larger than the RNA  
 1564 backbone reported here (Table 3). Though one  
 1565 should expect the backbone atoms of A-DNA to  
 1566 be 8.3 Å<sup>3</sup> smaller (based on our atom type volumes)  
 1567 due to the additional volume taken up by the 2'  
 1568 oxygen, this is not the case. Furthermore, in DNA  
 1569 the 2'-carbon volume is reported as 18.0 Å<sup>3</sup>, while  
 1570 we report a volume of 12.67 Å<sup>3</sup>; this drop is  
 1571 expected because of the loss of the hydrogen. But  
 1572 if we take the 2'-oxygen volume of 17.39 Å<sup>3</sup> into  
 1573 account, we now have a total volume of 30.07 Å<sup>3</sup> to

1574 fit into the space of 18.0 Å<sup>3</sup>. RNA structure must  
 1575 accommodate for this additional occupied space. In  
 1576 summary, the published A-DNA volumes are  
 1577 approximately equal for the bases and differ slightly  
 1578 for the backbone where A-DNA is packed less tight  
 1579 than RNA.  
 1580  
 1581  
 1582  
 1583

#### 1584 Implications in RNA packing

1585 Early results for proteins showed protein  
 1586 interiors are more tightly packed than amino acid  
 1587 crystals.<sup>3</sup> These results also indicated that tight  
 1588 packing and detailed interactions are important in  
 1589 protein folding. RNA tends to be seen as a loosely  
 1590 packed molecule, held together primarily by base-  
 1591 pairing and electrostatic interactions through back-  
 1592 bone alterations and metal ion coordination. This is  
 1593 borne out by a survey of a number of prominent  
 1594 papers in RNA structure and folding.<sup>38-55</sup> These  
 1595 papers mention electrostatics and hydrophobic  
 1596 effects as important factors in RNA folding, but  
 1597 none of them mention the importance close pack-  
 1598 ing. For instance, Doudna & Doherty argue that the  
 1599 hydrophobic effect, hydrogen bonding, metal ion  
 1600 coordination and VDW forces all contribute to the  
 1601 formation of compact structures.<sup>51</sup> They say that  
 1602 hydrophobic effects in RNA occur mainly at the  
 1603 level of secondary structure, making a contribution  
 1604 to vertical base stacking. Additionally, they assert  
 1605 that RNA folding is opposed by electrostatic  
 1606 repulsion from the negatively charged phosphate  
 1607 backbone.

1608 However, our results show, surprisingly, that  
 1609 RNA is actually packed more tightly than proteins.  
 1610 In essence, we demonstrate that close packing is as  
 1611 important for RNA folding as for proteins. This  
 1612 suggests a number of interesting energetic calcula-  
 1613 tions that might be worthwhile doing. To empha-  
 1614 size this point, we have modified the text as shown  
 1615 below and changed the title to: "Calculation of  
 1616 Standard Atomic Volumes for RNA Cores and  
 1617 Comparison with Proteins: RNA is packed more  
 1618 tightly than protein".

1619 Another interesting aspect of RNA packing  
 1620 illuminated by our volume calculations concerns  
 1621 the 2'-carbon atom. In DNA the 2'-carbon volume is  
 1622 reported as 18.0 Å<sup>3</sup>,<sup>12</sup> while we report a volume of  
 1623 12.67 Å<sup>3</sup>; this drop is expected because of the loss of  
 1624 the hydrogen. But if we take the 2'-oxygen volume  
 1625 of 17.39 Å<sup>3</sup> into account, we now have a total  
 1626 volume of 30.07 Å<sup>3</sup> to fit into the space of 18.0 Å<sup>3</sup>.  
 1627 Therefore, RNA structure must accommodate this  
 1628 additional occupied space.  
 1629

#### 1630 Practical applications

1631 We now point out in the paper how our  
 1632 parameter set is useful for RNA studies and, in  
 1633 fact, directly increases our understanding of RNA  
 1634 structure. Many applications of our volumes and  
 1635 radii come to mind. We provide three new data sets:  
 1636 a set of atomic RNA volumes, a set of RNA VDW  
 1637 radii and a variety of annotated sets of large,  
 1638

1639 non-redundant, RNA-containing PDB structures.  
 1640 Many programs used for structure solving and  
 1641 model refinement use VDW radii. Any RNA  
 1642 informatics endeavor requires begins with anno-  
 1643 tated sets of PDB structures, which we provide.

1644 **Packing density**

1645  
 1646 Structure–function research can involve the  
 1647 atomic radii and volumes of RNA in order to locate  
 1648 non-standard regions and possibly functional areas.  
 1649 In particular, the volumes can be used to measure  
 1650 the local packing density, the ratio of a given atom  
 1651 to its expected volume within a particular region.  
 1652 The local packing density can be used to determine  
 1653 more and less packed regions within a particular  
 1654 structure. Second, we can use the packing density to  
 1655 locate atoms with extreme volumes. This is useful in  
 1656 evaluating the quality of a crystal structures by  
 1657 locating areas that are packed too loosely or too  
 1658 tightly. Additionally, regions with extreme volumes  
 1659 may pinpoint active sites or other functional  
 1660 features. Finally, packing density is an accepted  
 1661 method for measuring the tightness of fit between  
 1662 RNA and a substrate, such as polymerases, RNases  
 1663 and other RNA-binding molecules.

1664 **Volume and PSV calculation**

1665  
 1666 Before our volume results, calculating molecular  
 1667 volumes for RNA containing macromolecules was  
 1668 limited. Two techniques have existed for determin-  
 1669 ing the volume of unknown particles: electron  
 1670 microscopy and small-angle X-ray scattering.<sup>56</sup>  
 1671 Both methods are problematic. Previous studies of  
 1672 50 S ribosomal subunits to determine their volume  
 1673 did so with a large range of 1.8–4.4 million cubic  
 1674 Angstroms.<sup>56</sup>

1675  
 1676 Using our published volume set, we can estimate  
 1677 the molecular volume of any RNA based solely on  
 1678 its sequence. For example, the Voronoi volume of  
 1679 50 S small subunit structure is 1,374,538 Å<sup>3</sup>. Based  
 1680 on the actual three-dimensional coordinates of the  
 1681 solved ribosome structure, the Richards’ rolling  
 1682 probe method<sup>16</sup> calculates the molecular volume to  
 1683 be 1,400,281 Å<sup>3</sup>. This slight difference of 1.8% is  
 1684 reasonable considering we are only using sequence  
 1685 information. Therefore, in essence, we can get a  
 1686 good estimate for the volume without knowing  
 1687 three-dimension coordinates. Further, if the  
 1688 sequence is known then the mass is readily  
 1689 calculated to obtain a partial specific volume for  
 1690 any unknown structure. Using our volumes, we  
 1691 calculate the partial specific volume of the 50 S  
 1692 subunit to be 0.617 ml/g which compares well to  
 1693 the published value of approximately 0.592 ml/g.<sup>56,  
 1694 57</sup>

1695  
 1696 We have also built a web tool† to perform this  
 1697 calculation of volume and PSV on an arbitrary RNA  
 1698 or protein sequence. For instance, application of  
 1699 the tool to the U65 snoRNA 172 nt consensus  
 1700 sequence,<sup>58</sup> which currently has an unknown  
 1701

structure, shows it to have a volume of 41,700.4 Å<sup>3</sup>  
 and a PSV of 0.569 ml/g.

**Exploration of the ribosome**

One immediate future application of our para-  
 meter sets is the analysis of the ribosome.<sup>59</sup> The  
 ribosome has an extremely complex intertwined  
 folding of protein and RNA that is currently not  
 fully understood. It is an open question how this  
 large macromolecule packs together. We can now  
 use our volume and radii parameters to analyze  
 internal solvent volumes. In a similar sense we can  
 evaluate which helices within the ribosome struc-  
 ture interact with which other helices. This is  
 similar in spirit to work done on membrane  
 proteins.<sup>60,61</sup> Finally, the exit tunnel is a site of  
 antibiotic binding; using our new parameter sets we  
 can trace out the volume and diameter as a function  
 of distance from the active site to better understand  
 how these molecules are functioning to block  
 translation.

**Conclusions**

In this study, we performed a careful parameter-  
 ization of currently available RNA structures to  
 obtain a universal, self-consistent set of volumes,  
 denoted as the NucProt parameter set. This compo-  
 site set can be applied to both RNA and protein. In  
 addition, several factors such as crystal symmetry,  
 structural complexity and protein and solvent  
 interactions are taken into account for their influ-  
 ence on the final results. Using two measures, the  
 percentage of well-packed atoms and final volumes,  
 the impact of each factor was assessed on the data.  
 While all the factors affected the percentage of  
 well-packed atoms, none of them had any affect on  
 the final volumes. From these volume calculations,  
 it is immediately apparent that the RNA backbone  
 is not as tightly packed as its base as determined by  
 its standard deviation and also its percentage of  
 well-packed atoms. For RNA, the calculated partial  
 specific volume corresponded well with its experi-  
 mental value. When compared to proteins, RNA is  
 found to be more dense, because its partial specific  
 volume is smaller. Comparing common atom types  
 between protein and RNA showed that in nine of 12  
 cases, RNA has a smaller volume and is therefore  
 packed tighter. Further, when comparing aromatic  
 protein side-chains to the RNA bases, the partial  
 specific volume for RNA bases was again smaller  
 than the protein side-chains as well as their average  
 volume per atom. Thus, RNA packs more tightly  
 than protein, but based only on well-packed atoms.  
 A-form DNA, on the other hand, has approximately  
 the same base volumes as RNA, though the back-  
 bones differ by more than the bases it is within the  
 standard deviation of the total volume. In con-  
 clusion, RNA packs more tightly than protein and  
 approximately the same as DNA.

*Location of files, programs, scripts and statistics.*

1765 Further details on parameter sets, additional statistics,  
 1766 perl and shell scripts, packaged program files,  
 1767 and the raw volume data are provided online†.

1770 **References**

1771  
 1772 1. Bondi, A. (1964). Van der Waals volumes and radii.  
 1773 *J. Phys. Chem.* **68**, 441–451.  
 1774 2. Chothia, C. (1974). Hydrophobic bonding and acces-  
 1775 sible surface area in proteins. *Nature*, **248**, 338–339.  
 1776 3. Richards, F. M. (1974). The interpretation of protein  
 1777 structures: total volume, group volume distributions  
 1778 and packing density. *J. Mol. Biol.* **82**, 1–14.  
 1779 4. Finney, J. L. (1975). Volume occupation, environment  
 1780 and accessibility in proteins. The problem of the  
 1781 protein surface. *J. Mol. Biol.* **96**, 721–732.  
 1782 5. Harpaz, Y., Gerstein, M. & Chothia, C. (1994). Volume  
 1783 changes on protein folding. *Structure*, **2**, 641–649.  
 1784 6. Li, A. J. & Nussinov, R. (1998). A set of Van der Waals  
 1785 and coulombic radii of protein atoms for molecular  
 1786 and solvent-accessible surface calculation, packing  
 1787 evaluation, and docking. *Proteins: Struct. Funct. Genet.*  
 1788 **32**, 111–127.  
 1789 7. Liang, J., Edelsbrunner, H., Fu, P., Sudhakar, P. V. &  
 1790 Subramaniam, S. (1998). Analytical shape compu-  
 1791 tation of macromolecules: II. Inaccessible cavities in  
 1792 proteins. *Proteins: Struct. Funct. Genet.* **33**, 18–29.  
 1793 8. Liang, J., Edelsbrunner, H., Fu, P., Sudhakar, P. V. &  
 1794 Subramaniam, S. (1998). Analytical shape compu-  
 1795 tation of macromolecules: I. Molecular area and  
 1796 volume through alpha shape. *Proteins: Struct. Funct.*  
 1797 *Genet.* **33**, 1–17.  
 1798 9. Tsai, J., Taylor, R., Chothia, C. & Gerstein, M. (1999).  
 1799 The packing density in proteins: standard radii and  
 1800 volumes. *J. Mol. Biol.* **290**, 253–266.  
 1801 10. Tsai, J. & Gerstein, M. (2002). Calculations of protein  
 1802 volumes: sensitivity analysis and parameter database.  
 1803 *Bioinformatics*, **18**, 985–995.  
 1804 11. Tsai, J., Voss, N. & Gerstein, M. (2001). Determining  
 1805 the minimum number of types necessary to represent  
 1806 the sizes of protein atoms. *Bioinformatics*, **17**, 949–956.  
 1807 12. Nadassy, K., Tomas-Oliveira, I., Alberts, I., Janin, J. &  
 1808 Wodak, S. J. (2001). Standard atomic volumes in  
 1809 double-stranded DNA and packing in protein–DNA  
 1810 interfaces. *Nucl. Acids Res.* **29**, 3362–3376.  
 1811 13. Chothia, C. (1975). Structural invariants in protein  
 1812 folding. *Nature*, **254**, 304–308.  
 1813 14. Janin, J. & Chothia, C. (1990). The structure of protein–  
 1814 protein recognition sites. *J. Biol. Chem.* **265**,  
 1815 16027–16030.  
 1816 15. Janin, J. (1979). Surface and inside volumes in  
 1817 globular proteins. *Nature*, **277**, 491–492.  
 1818 16. Richards, F. M. (1985). Calculation of molecular  
 1819 volumes and areas for structures of known geometry.  
 1820 *Methods Enzymol.* **115**, 440–464.  
 1821 17. Gerstein, M., Tsai, J. & Levitt, M. (1995). The volume  
 1822 of atoms on the protein surface: calculated from  
 1823 simulation, using Voronoi polyhedra. *J. Mol. Biol.* **249**,  
 1824 955–966.  
 1825 18. Gerstein, M. & Chothia, C. (1996). Packing at the  
 1826 protein–water interface. *Proc. Natl Acad. Sci. USA*, **93**,  
 1827 10167–10172.  
 1828 19. Hubbard, S. J. & Argos, P. (1995). Detection of internal  
 1829 cavities in globular proteins. *Protein Eng.* **8**, 1011–1015.  
 1830 20. Pontius, J., Richelle, J. & Wodak, S. J. (1996).

Deviations from standard atomic volumes as a quality  
 1828 measure for protein crystal structures. *J. Mol. Biol.* **264**,  
 1829 121–136.  
 1830 21. Gerstein, M., Sonnhammer, E. L. & Chothia, C. (1994).  
 1831 Volume changes in protein evolution. *J. Mol. Biol.* **236**,  
 1832 1067–1078.  
 1833 22. Gerstein, M. (1998). How representative are the  
 1834 known structures of the proteins in a complete  
 1835 genome? A comprehensive structural census. *Fold.*  
 1836 *Des.* **3**, 497–512.  
 1837 23. Gerstein, M. & Krebs, W. (1998). A database of  
 1838 macromolecular motions. *Nucl. Acids Res.* **26**, 4280–  
 1839 4290.  
 1840 24. Krebs, W. G. & Gerstein, M. (2000). The morph server:  
 1841 a standardized system for analyzing and visualizing  
 1842 macromolecular motions in a database framework.  
 1843 *Nucl. Acids Res.* **28**, 1665–1675.  
 1844 25. David, C. W. (1988). Voronoi polyhedra as structure  
 1845 probes in large molecular systems. *Biopolymers*, **27**,  
 1846 339–344.  
 1847 26. Finney, J. L. (1978). Volume occupation, environment,  
 1848 and accessibility in proteins. Environment and molec-  
 1849 ular area of RNase-S. *J. Mol. Biol.* **119**, 415–441.  
 1850 27. Dunbrack, R. L., Jr (1999). Comparative modeling of  
 1851 CASP3 targets using PSI-BLAST and SCWRL. *Proteins*  
 1852 *Suppl.* **3**, 81–87.  
 1853 28. Koehl, P. & Delarue, M. (1997). The native sequence  
 1854 determines sidechain packing in a protein, but does  
 1855 optimal sidechain packing determine the native  
 1856 sequence? *Pac. Symp. Biocomput.*, 198–209.  
 1857 29. Lee, C. & Levitt, M. (1997). Packing as a structural  
 1858 basis of protein stability: understanding mutant  
 1859 properties from wildtype structure. *Pac. Symp. Bio-*  
 1860 *comput.*, 245–255.  
 1861 30. Voronoi, G. F. (1908). Nouvelles applications des  
 1862 paramètres continus á la théorie de formas quad-  
 1863 ratiques. *J. Reine. Angew. Math.* **134**, 198–287.  
 1864 31. Bernal, J. D. & Finney, J. L. (1967). Random close-  
 1865 packed hard-sphere model II. geometry of random  
 1866 packing of hard spheres. *Disc. Faraday Soc.* **43**, 62–69.  
 1867 32. Gellatly, B. J. & Finney, J. L. (1982). Calculation of  
 1868 protein volumes: an alternative to the Voronoi  
 1869 procedure. *J. Mol. Biol.* **161**, 305–322.  
 1870 33. Tsai, J., Gerstein, M. & Levitt, M. (1997). Simulating  
 1871 the minimum core for hydrophobic collapse in  
 1872 globular proteins. *Protein Sci.* **6**, 2606–2616.  
 1873 34. Brunger, A. T., Adams, P. D., Clore, G. M., DeLano,  
 1874 W. L., Gros, P., Grosse-Kunstleve, R. W. *et al.* (1998).  
 1875 Crystallography & NMR system: a new software suite  
 1876 for macromolecular structure determination. *Acta*  
 1877 *Crystallog. sect. D: Biol. Crystallogr.* **54**, 905–921.  
 1878 35. Berman, H. M., Olson, W. K., Beveridge, D. L.,  
 1879 Westbrook, J., Gelbin, A., Demeny, T. *et al.* (1992).  
 1880 The nucleic acid database. A comprehensive rela-  
 1881 tional database of three-dimensional structures of  
 1882 nucleic acids. *Biophys. J.* **63**, 751–759.  
 1883 36. Berman, H. M., Westbrook, J., Feng, Z., Gilliland, G.,  
 1884 Bhat, T. N., Weissig, H. *et al.* (2000). The Protein Data  
 1885 Bank. *Nucl. Acids Res.* **28**, 235–242.  
 1886 37. Durchschlag, H. (1986). Specific volumes of biological  
 1887 macromolecules and some other molecules of bio-  
 1888 logical interest. In *Thermodynamic Data for Biochemistry*  
 1889 *and Biotechnology* (Hinz, H.-J., ed.), pp. 45–128,  
 1890 Springer, Berlin.  
 1891 38. Draper, D. E. (2004). A guide to ions and RNA  
 1892 structure. *RNA*, **10**, 335–343.  
 1893 39. Su, L. J., Brenowitz, M. & Pyle, A. M. (2003). An

† <http://geometry.molmovdb.org>

1891 alternative route for the folding of large RNAs: 1954  
 1892 apparent two-state folding by a group II intron 1955  
 1893 ribozyme. *J. Mol. Biol.* **334**, 639–652. 1956  
 1894 40. Pyle, A. M. (2002). Metal ions in the structure and 1957  
 1895 function of RNA. *J. Biol. Inorg. Chem.* **7**, 679–690. 1958  
 1896 41. Misra, V. K. & Draper, D. E. (2002). The linkage 1959  
 1897 between magnesium binding and RNA folding. 1960  
 1898 *J. Mol. Biol.* **317**, 507–521. 1961  
 1899 42. Kim, H. D., Nienhaus, G. U., Ha, T., Orr, J. W., 1962  
 1900 Williamson, J. R. & Chu, S. (2002). Mg<sup>2+</sup>-dependent 1963  
 1901 conformational change of RNA studied by fluor- 1964  
 1902 escence correlation and FRET on immobilized single 1965  
 1903 molecules. *Proc. Natl Acad. Sci. USA*, **99**, 4284–4289. 1966  
 1904 43. Silverman, S. K., Deras, M. L., Woodson, S. A., 1967  
 1905 Scaringe, S. A. & Cech, T. R. (2000). Multiple folding 1968  
 1906 pathways for the P4-P6 RNA domain. *Biochemistry*, **39**, 1969  
 1907 12465–12475. 1970  
 1908 44. Hanna, R. & Doudna, J. A. (2000). Metal ions in 1971  
 1909 ribozyme folding and catalysis. *Curr. Opin. Chem. Biol.* 1972  
 1910 **4**, 166–170. 1973  
 1911 45. Ryder, S. P. & Strobel, S. A. (1999). Nucleotide analog 1974  
 1912 interference mapping of the hairpin ribozyme: impli- 1975  
 1913 cations for secondary and tertiary structure forma- 1976  
 1914 tion. *J. Mol. Biol.* **291**, 295–311. 1977  
 1915 46. Rook, M. S., Treiber, D. K. & Williamson, J. R. (1999). 1978  
 1916 An optimal Mg(2+) concentration for kinetic folding 1979  
 1917 of the tetrahymena ribozyme. *Proc. Natl Acad. Sci. USA*, **96**, 1980  
 1918 12471–12476. 1981  
 1919 47. Batey, R. T. & Doudna, J. A. (1998). The parallel 1982  
 1920 universe of RNA folding. *Nature Struct. Biol.* **5**, 1983  
 1921 337–340. 1984  
 1922 48. Strobel, S. A. & Doudna, J. A. (1997). RNA seeing 1985  
 1923 double: close-packing of helices in RNA tertiary 1986  
 1924 structure. *Trends Biochem. Sci.* **22**, 262–266. 1987  
 1925 49. McConnell, T. S., Herschlag, D. & Cech, T. R. (1997). 1988  
 1926 Effects of divalent metal ions on individual steps of 1989  
 1927 the tetrahymena ribozyme reaction. *Biochemistry*, **36**, 1990  
 1928 8293–8303. 1991  
 1929 50. Doudna, J. A. & Cate, J. H. (1997). RNA structure: 1992  
 1930 crystal clear? *Curr. Opin. Struct. Biol.* **7**, 310–316. 1993  
 1931 51. Doudna, J. A. & Doherty, E. A. (1997). Emerging 1994  
 1932 themes in RNA folding. *Fold. Des.* **2**, R65–R70. 1995  
 1933 52. Cate, J. H., Hanna, R. L. & Doudna, J. A. (1997). A 1996  
 1934 magnesium ion core at the heart of a ribozyme 1997  
 1935 domain. *Nature Struct. Biol.* **4**, 553–558. 1998  
 1936 53. Draper, D. E. (1996). Strategies for RNA folding. 1999  
 1937 *Trends Biochem. Sci.* **21**, 145–149. 2000  
 1938 54. Cate, J. H., Gooding, A. R., Podell, E., Zhou, K., 2001  
 1939 Golden, B. L., Szewczak, A. A. *et al.* (1996). RNA 2002  
 1940 tertiary structure mediation by adenosine platforms. 2003  
 1941 *Science*, **273**, 1696–1699. 2004  
 1942 55. Pyle, A. M. & Green, J. B. (1995). RNA folding. *Curr.* 2005  
 1943 *Opin. Struct. Biol.* **5**, 303–310. 2006  
 1944 56. Van Holde, K. E. & Hill, W. E. (1974). General physical 2007  
 1945 properties of ribosomes. In *Ribosomes* (Nomura, M., 2008  
 1946 Tissieres, A. & Lengyel, P., eds), pp. 53–91, Cold 2009  
 1947 Spring Harbor Press, Cold Spring Harbor, NY. 2010  
 1948 57. Hill, W. E., Rossetti, G. P. & Van Holde, K. E. (1969). 2011  
 1949 Physical studies of ribosomes from *Escherichia coli*. 2012  
 1950 *J. Mol. Biol.* **44**, 263–277. 2013  
 1951 58. Ganot, P., Bortolin, M. L. & Kiss, T. (1997). Site-specific 2014  
 1952 pseudouridine formation in preribosomal RNA is 2015  
 1953 guided by small nucleolar RNAs. *Cell*, **89**, 799–809. 2016  
 1954 59. Steitz, T. A., Moore, P. B. (2004). 2017  
 1955 60. Eilers, M., Shekar, S. C., Shieh, T., Smith, S. O. & 2018  
 1956 Fleming, P. J. (2000). Internal packing of helical 2019  
 1957 membrane proteins. *Proc. Natl Acad. Sci. USA*, **97**, 2020  
 1958 5796–5801. 2021  
 1959 61. Gerstein, M. & Chothia, C. (1999). Perspectives: signal 2022  
 1960 transduction. *Proteins in motion. Science*, **285**, 2023  
 1961 1682–1683. 2024  
 1962 62. Lee, A. & Chalikian, T. V. (2001). Volumetric charac- 2025  
 1963 terization of the hydration properties of heterocyclic 2026  
 1964 bases and nucleosides. *Biophys. Chem.* **92**, 209–227. 2027  
 1965

Edited by D. E. Draper

(Received 19 August 2004; received in revised form 24 November 2004; accepted 24 November 2004)