

# Geometrical and Sequence Characteristics of $\alpha$ -Helices in Globular Proteins

Sandeep Kumar and Manju Bansal

Molecular Biophysics Unit, Indian Institute of Science, Bangalore 560 012, India

**ABSTRACT** Understanding the sequence-structure relationships in globular proteins is important for reliable protein structure prediction and de novo design. Using a database of 1131  $\alpha$ -helices with nonidentical sequences from 205 nonhomologous globular protein chains, we have analyzed structural and sequence characteristics of  $\alpha$ -helices. We find that geometries of more than 99% of all the  $\alpha$ -helices can be simply characterized as being linear, curved, or kinked. Only a small number of  $\alpha$ -helices (~4%) show sharp localized bends in their middle regions, and thus are classified as kinked. Approximately three-fourths (~73%) of the  $\alpha$ -helices in globular proteins show varying degrees of smooth curvature, with a mean radius of curvature of  $65 \pm 33$  Å; longer helices are less curved. Computation of helix accessibility to the solvent indicates that nearly two-thirds of the helices (~66%) are largely buried in the protein core, and the length and geometry of the helices are not correlated with their location in the protein globule. However, the amino acid compositions and propensities of individual amino acids to occur in  $\alpha$ -helices vary with their location in the protein globule, their geometries, and their lengths. In particular, Gln, Glu, Lys, and Arg are found more often in helices near the surface of globular proteins. Interestingly, kinks often seem to occur in regions where amino acids with low helix propensities (e.g.,  $\beta$ -branched and aromatic residues) cluster together, in addition to those associated with the occurrence of proline residues. Hence the propensities of individual amino acids to occur in a given secondary structure depend not only on conformation but also on its length, geometry, and location in the protein globule.

## INTRODUCTION

Using basic stereochemical principles and recognizing the planarity of the peptide bond, Pauling et al. (1951) deduced the  $\alpha$ -helical structure as one of the two possible hydrogen-bonded conformations for the polypeptide chains in proteins. Subsequent solution of the three-dimensional structure for myoglobin by x-ray crystallography (Kendrew et al., 1958) proved the validity of this model. Today it is well known that the  $\alpha$ -helix is one of the most common secondary structural motifs in proteins. Because of the regular occurrence of ( $i, i - 4$ ) hydrogen bonds in the main chain, the  $\alpha$ -helices in proteins are generally quite uniform, with their helical parameters unit twist ( $t$ ) and unit height ( $h$ ) lying within well-defined ranges. However,  $\alpha$ -helices have been observed to be distorted for a variety of reasons, for example, occurrence of proline residues (Barlow and Thornton, 1988), solvent-induced distortions (Blundell et al., 1983), and peptide bond distortions (Chakrabarti et al., 1986).

A systematic analysis of helix geometries found in globular proteins was first reported by Barlow and Thornton (1988). They found that most of the  $\alpha$ -helices in globular proteins were smoothly curved. We had previously ana-

lyzed sequences and geometries of 64 long  $\alpha$ -helices with lengths greater than 25 residues taken from 45 globular proteins (Kumar and Bansal, 1996). It was observed that long  $\alpha$ -helices have unique amino acid composition characteristics that distinguish them from the short  $\alpha$ -helices, taken from the data set of Richardson and Richardson (1988). The amino acid distribution in long  $\alpha$ -helices could be correlated with their overall geometry and, in the case of curved helices, with the radius of curvature. We have now extended these studies to study variation in amino acid distribution in  $\alpha$ -helices with respect to several parameters, with a data set that includes helices of all lengths in the June 1995 Protein Data Bank (PDB) (Bernstein et al., 1977). Using HELANAL (Kumar and Bansal, 1996) to characterize  $\alpha$ -helix geometry in a database of 1131  $\alpha$ -helices of nonidentical sequences, from 205 nonhomologous globular protein chains, we found that the following rules are generally observed for all helices: 1) Almost all  $\alpha$ -helices can be unambiguously characterized as being linear, curved, or kinked. 2) The majority of  $\alpha$ -helices are smoothly curved and show varying degrees of curvature, with longer helices preferring a sharp kink to high smooth curvature. 3) Amino acid distribution and propensities of individual amino acids to occur in an  $\alpha$ -helix depend on its location in the protein globule, on its length and geometry, and in the case where the helix is smoothly curved, on its radius of curvature.

## MATERIALS AND METHODS

### Composition of the database

We have used the June 1995 list of nonhomologous (sequence identity  $\leq$  25%) protein chains compiled by Hobohm and Sander (Hobohm et al.,

Received for publication 1 December 1997 and in final form 8 June 1998.

Address reprint requests to Dr. Manju Bansal, Molecular Biophysics Unit, Indian Institute of Science, Bangalore-560 012, India. Tel.: 91-80-3092534; Fax: 91-80-3348535 or 91-80-3341683; E-mail: mb@mbu.iisc.ernet.in.

Dr. Kumar's present address is Bldg. 469, Rm. 147, NCI-FCRDC, P.O. Box B, Frederick, MD 21702.

© 1998 by the Biophysical Society

0006-3495/98/10/1935/10 \$2.00

1992; Hobohm and Sander, 1994) to select a subset database of 205 nonhomologous globular protein chains whose three-dimensional structures have been solved by x-ray crystallography to a resolution of 2.5 Å or better. One thousand one hundred thirty-one (1131)  $\alpha$ -helices of nonidentical sequences and nine or more amino acid residues found in these protein chains constitute the database for the studies presented here. We have only considered helices consisting of nine or more residues because these helices contain at least one residue in which both NH and C=O are hydrogen bonded. Besides this, because HELANAL characterizes the overall geometry of an  $\alpha$ -helix from its local structural features (Kumar and Bansal, 1996, and briefly described below), the use of nine residues as the minimum helix length ensures that we have sufficient data points to deduce its overall geometry.

## Helix boundaries

Initially, protein chain segments of length greater than eight residues and defined as helices by the Dictionary of Protein Secondary Structure (DSSP) (Kabsch and Sander, 1983) were taken for each of the 205 protein chains. Boundaries of these helices were checked and, if necessary, reassigned according to the following additional criteria:

1. The distance  $|O_i - N_{i+4}| \leq 3.5$  Å at helix termini.
2. Angles between successive local helix axes at the helix termini region are less than  $20^\circ$  (Kumar and Bansal, 1996).

## Helix geometry and structural parameters of $\alpha$ -helices

We have used an in-house developed computer program called HELANAL (Kumar and Bansal, 1996) to deduce the overall geometry of an  $\alpha$ -helix from its local structural features, namely, local bending angles and local helix axes computed for every four consecutive  $C^\alpha$  atoms. This window of four  $C^\alpha$  atoms was slid over the length of the helix in steps of one  $C^\alpha$  atom.

The geometry of the helix is assigned from the values of four parameters, namely, the maximum local bending angle in the helix, defining the largest distortion in the helix; the root mean square deviations for the line (rmsL) and the circle (rmsC) fitted to the local helix origins that trace out the three-dimensional path of helix axis; and  $r^2$ , the square of linear correlation coefficient for the line fit. If the maximum bending angle for the helix is  $\geq 20^\circ$ , the helix is classified as kinked. If the maximum bending angle is  $< 20^\circ$ , then we look at the other three parameters. If both rmsC and rmsL are  $> 1$  Å, we do not assign any geometry type to the helix because this implies large and random distortions in the helix. If either one or both of the parameters rmsC and rmsL are less than 1 Å, then we compute the ratio, rmsL/rmsC. A value greater than 1 for the ratio indicates that the circle fits the local helix origins better than the line, and we assign helix geometry as curved. If the ratio is  $\leq 0.7$ , it is implied that the line fits the local helix origins better than the circle. We validate the helix geometry to be genuinely linear if  $r^2 \geq 0.8$ . However, if the validation fails, that is,  $r^2 < 0.8$ , we do not assign any geometrical type to the helix. A value of 0.7–1.0 for the ratio defines a zone where line and the circle fit the local helix origins almost equally well. We have only a few  $\alpha$ -helices in this zone. In this zone, if  $r^2 \geq 0.8$ , i.e., the local helix origins significantly correlate with line, the helix geometry is again classified as linear. But, if  $r^2 \leq 0.5$ , i.e., correlation of the local helix origins with the line is insignificant and could be expected by random chance, the helix geometry is classified as curved. For values of  $r^2$  in the range between 0.5 and 0.8, the helix geometry type is ambiguous. Hence no geometrical type is assigned to the helix.

Helical parameters of  $\alpha$ -helices, viz., unit height and unit twist, are also computed by HELANAL for every four  $C^\alpha$  atoms. For each of the 205 nonhomologous protein chains in our database, there exists a corresponding DSSP file containing secondary structure assignments. Among several other parameters, these DSSP files contain values of virtual torsion angles and  $\phi$ ,  $\psi$  for each residue in the protein chain. We have used the values corresponding to the residues in 1131  $\alpha$ -helices for structural analysis of the helices.

## Estimation of kink angle in the case of kinked helices

As mentioned above, the geometry of an  $\alpha$ -helix was classified as kinked if the maximum local bending angle, computed by HELANAL, exceeds  $20^\circ$ . The kinks are localized in a small region in the middle of the kinked helices. The region of kink is characterized by the large values of the local bending angles ( $> 20^\circ$ ) and hydrogen bond distances ( $> 3.5$  Å) at several consecutive residues, and is identified from plots of the local bending angle at the residue number  $i$  and the hydrogen bond distance ( $O_{i-2} - N_{i+2}$ ) versus the residue number  $i$  in the helix (Kumar and Bansal, 1996). Each kinked  $\alpha$ -helix was then divided into two segments, consisting of fragments before and after the region of kink, and the kink angle was calculated as the angle between the least-squares lines fitted to the local helix origins in these two segments after the kink region was excluded.

## Computation of helix accessibility to the solvent

Solvent accessibility (Vorobjev and Hermans, 1997) of an  $\alpha$ -helix is calculated as the ratio of accessible surface area of the helix within the protein to its accessible surface area upon its isolation from the protein. The accessible surface area of the helix is the sum of the accessible surface areas of all of the atoms in the helix. The accessible surface area for each atom in the helix is calculated by the method of Lee and Richard (1971), using a probe radius of 1.4 Å. For each helix, the whole protein in the crystal asymmetrical unit including all subunits has been considered for computation of its helix accessibility. We have used helix accessibility as a measure of its location in the protein globule.

## Statistical methods

The database of 1131  $\alpha$ -helices is subdivided into various categories according to helix location in tertiary structure (i.e., their solvent accessibility) and geometry (viz., linear, curved, and kinked). Furthermore, the curved helices are classified according to the radius of curvature ranges. Helices in the database are also subdivided into two broad length groups, viz., short helices (9–21 residues) and long helices ( $> 21$  residues), i.e., those that are significantly longer than one full repeat (18 residues) for  $\alpha$ -helices. In each of these classes at least five amino acid residues of each type are present in the database, and propensities for individual amino acid are computed according to the formula given by Williams et al. (1987). The  $\chi^2$  test has been used to determine whether the differences in the amino acid distributions of  $\alpha$ -helices in any two classes are significant (Medhi, 1992). To measure the differences in the sequence compositions of helices in any two classes, Euclidean and Hamming distances are also computed in 20-dimensional amino acid composition space (Chou and Zhang, 1995). Because these two parameters showed identical trends in all calculations, only the values of Hamming distances have been discussed in detail. Change in proportion of an amino acid in a given helix class is considered to be significant if it is greater than  $2\sigma$ , where  $\sigma$  is the estimated standard deviation (Medhi, 1992).

## RESULTS AND DISCUSSION

The commonly assigned structural parameters, unit height and unit twist for the  $\alpha$ -helices, do not show large variations, as also noted in our earlier analysis of long  $\alpha$ -helices (Kumar and Bansal, 1996). Mean unit twist for the 1131  $\alpha$ -helices ranges between  $94^\circ$  and  $103^\circ$ , with an average value of  $(99 \pm 1)^\circ$ . The mean unit height for these helices ranges between 1.3 Å and 1.7 Å, with an average value of  $(1.51 \pm 0.03)$  Å. The average values for the backbone torsion angles ( $\phi$ ,  $\psi$ ) in the middle regions of 1131  $\alpha$ -helices, i.e., the residues other than those in the first and last

turns in the helices, are  $(-63 \pm 7)^\circ$  and  $(-43 \pm 8)^\circ$ , respectively, whereas the average value for the virtual torsion angles computed for every four successive  $C^\alpha$  atoms in the 1131 helices is  $(50 \pm 6)^\circ$ . These observations indicate that structural parameters of  $\alpha$ -helices are quite uniform, and deviations from regularity (linearity) in  $\alpha$ -helices cannot be estimated from these parameters alone.

### Geometrical and length distribution of $\alpha$ -helices in globular proteins

Of 1131  $\alpha$ -helices in our database, with helix lengths ranging between 9 and 37 residues, 1122 (99.2%) can be characterized as being linear, curved, or kinked. In only nine cases (0.8%) the helices are either irregular or their geometries cannot be assigned unambiguously. A majority (821 of 1131, 72.6%) of the  $\alpha$ -helices in globular proteins are smoothly curved, as also found in the earlier analyses by Barlow and Thornton (1988) and Kumar and Bansal (1996). The proportions of linear and kinked helices are 22.6% (256 of 1131) and 4% (45 of 1131), respectively. Fig. 1 shows the distribution of  $\alpha$ -helices with respect to helix length (number of residues in an  $\alpha$ -helix) in each of the three geometric classes, as well as in the whole database of 1131  $\alpha$ -helices. The number of crystallographically observed helices decreases rapidly with the increase in helix length. This observation is in conformity with the observations of Barlow and Thornton (1988), Blundell and Zhu (1995), and Zhu and Blundell (1996). However, our database does not show maxima corresponding to helix lengths of 7, 11, 15.5, and 22 residues, as observed by Srinivasan (1976). The mean helix length in our database of 1131  $\alpha$ -helices is  $(14 \pm 5)$  residues, whereas the mean helix lengths for linear, curved, and kinked helices are  $(12 \pm 3)$ ,  $(14 \pm 5)$ , and  $(20 \pm 6)$ , respectively, indicating that kinked helices are usually longer than the average.

The database of 1131  $\alpha$ -helices is also subdivided into two length classes, namely short  $\alpha$ -helices (containing 9–21 residues) accounting for 1046 (92.5%) of 1131  $\alpha$ -helices and 85 (7.5%) long  $\alpha$ -helices ( $> 21$  residues). The 1046 short  $\alpha$ -helices include 758 (72.5%) curved, 251 (24%) linear, and 28 (2.7%) kinked  $\alpha$ -helices, whereas the 85 long  $\alpha$ -helices include 63 (74.1%) curved, 5 (5.9%) linear, and 17 (20%) kinked  $\alpha$  helices. Thus, among long  $\alpha$ -helices, the proportion of kinked helices increases at the expense of linear helices.

### Longer helices are generally less curved than shorter ones

The mean radius of curvature for 821 curved helices is  $(65 \pm 33)$  Å (range 8–373 Å). Fig. 2 *a* shows the distribution of helices with respect to the radius of curvature. The average radii of curvature in the two length classes of curved helices, viz., short  $\alpha$ -helices with 9–21 residues and long helices with  $>21$  residues, are  $(63 \pm 34)$  Å and  $(96 \pm$

$32)$  Å, respectively. Fig. 2 *b* shows a plot of mean helix length for helices in different curvature ranges versus the radius of curvature. It can be seen that the mean helix length increases linearly with the radius of curvature (square of linear correlation coefficient,  $r^2 = 0.96$ ). It is interesting to find that the minimum radius of curvature for each helix length also increases linearly with the helix length (square of linear correlation coefficient,  $r^2 = 0.93$ ). These observations indicate that longer helices tend to be less curved than shorter ones. Barlow and Thornton (1988) had also noted that long helices in coiled coils are less curved than the short helices in globular proteins. Thus it appears that a high degree of curvature cannot be sustained or propagated over a long polypeptide helix.

### Do $\alpha$ -helices with different lengths or geometries occur in different regions of the protein globule?

Most of the  $\alpha$ -helices in globular proteins are quite short; the mean helix length in the present data set is 14 residues. It is generally believed that the length of  $\alpha$ -helices is limited by the size of the protein (Creighton, 1993). Thus one might expect a correlation between the helix length and the size of the protein as defined by the number of amino acid residues in the protein. In our database, we do not observe any such correlation. For example, a small 56-residue polypeptide chain of ColE1 Rop contains two long  $\alpha$ -helices, forming an antiparallel coiled coil (Banner et al., 1987).

It may also be expected that part or the whole of longer helices may be exposed to solvent, whereas shorter helices can be readily buried in the protein core, i.e., long and short helices may occur in different regions of the protein globule. We have tested this hypothesis by computing solvent accessibility for each of the 1131  $\alpha$ -helices. Accessibility of an  $\alpha$ -helix defines its location in the protein tertiary structure (Blundell and Zhu, 1995; Zhu and Blundell, 1996) and in quaternary structure, in the case of proteins containing more than one subunit. Mean solvent accessibility for the 1131  $\alpha$ -helices is  $(32 \pm 15)\%$ . The majority of the helices (751 of 1131, 66.4%) are “largely buried” in the protein core, i.e., have less than 40% solvent accessibility. 356 (31.5%) helices are “partially buried” (40–60% solvent accessibility), and only 24 (2.1%) are “largely exposed” ( $> 60\%$  solvent accessibility) to the solvent. The mean helix lengths for largely buried, partly buried, and largely exposed helices are  $(14 \pm 4)$ ,  $(14 \pm 5)$ , and  $(16 \pm 6)$ , respectively. Fig. 3 shows a plot of solvent accessibilities of 1131  $\alpha$ -helices versus their lengths. There is no significant correlation between length and solvent accessibility, indicating that an  $\alpha$ -helix of any length can occur anywhere in the protein globule, and there is no segregation between long and short  $\alpha$ -helices into different regions of the protein globule.

A kink in an  $\alpha$ -helix, containing a proline residue after the first four positions, facilitates compensatory hydrogen bond formation between solvent and residues  $i - 3$  and  $i -$

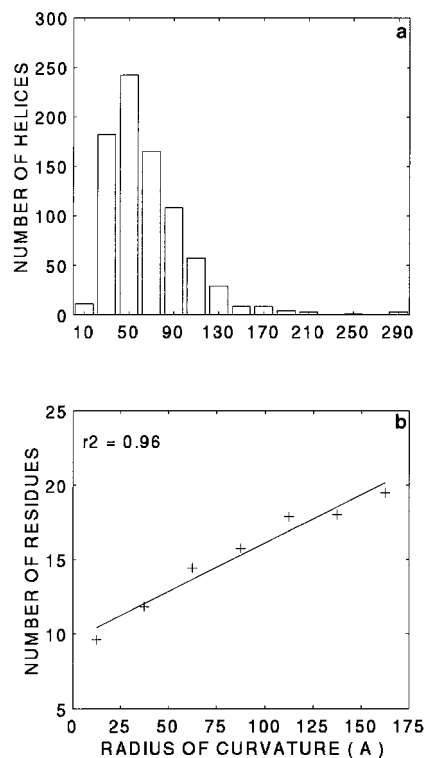
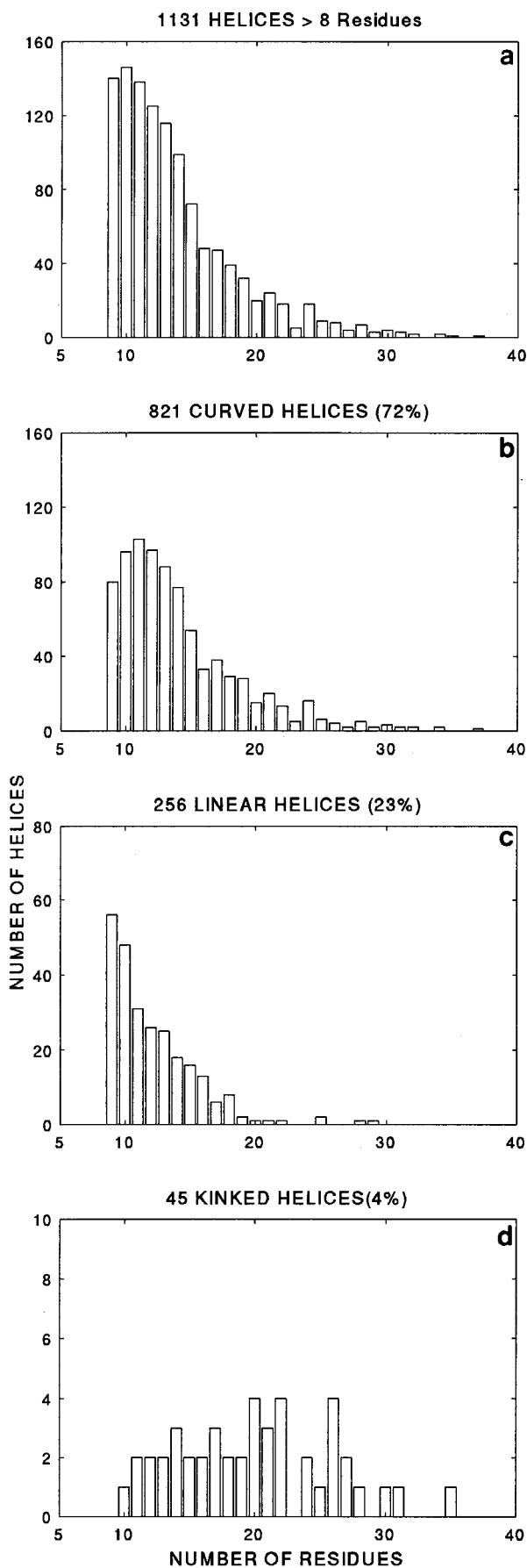


FIGURE 2 (a) Bar diagram showing distribution of 821 curved  $\alpha$ -helices with respect to their radius of curvature. The  $x$  axis denotes the radius of curvature, and the  $y$  axis denotes the number of helices. The mean radius of curvature for all 821 curved  $\alpha$ -helices is  $(65 \pm 33)$  Å. (b) Mean number of residues for  $\alpha$ -helices in each of the seven classes, viz. radius of curvature ranges of 0–25 Å, 25–50 Å, 50–75 Å, 75–100 Å, 100–125 Å, 125–150 Å, and 150–175 Å. The mean helix length (number of residues) ( $y$  axis) is plotted against the midvalues of radii of curvature ranges in the above seven ranges ( $x$  axis).

4 (relative to a proline residue at the  $i$ th position) (Woolfson and Williams, 1990). Besides, water is known to induce bends and kinks in  $\alpha$ -helices (Blundell et al., 1983). According to these concepts, one would expect helix geometry to be correlated with its solvent accessibility and kinked helices to occur predominantly on the protein surface. Fig. 4 shows distributions of linear, curved, and kinked helices with helix accessibility, and it is clear that there is no significant correlation between helix geometry and solvent accessibility. In addition, it is found that solvent accessibilities of Pro and Gly residues at kink regions are similar to the mean residue accessibility of that kinked helix, indicating that there is no noticeable preference for Pro and Gly at kinks to be exposed on the protein surface.

FIGURE 1 Histograms showing the length distribution of all  $\alpha$ -helices (>8 residues in length) in globular proteins. The  $x$  axis denotes the number of residues in the  $\alpha$ -helix (helix length). The  $y$  axis denotes the number of helices. (a) All 1131  $\alpha$ -helices. (b) 821 curved  $\alpha$ -helices. (c) 256 linear  $\alpha$ -helices. (d) 45 kinked  $\alpha$ -helices. In the case of nine  $\alpha$ -helices, no helix geometry is assigned.



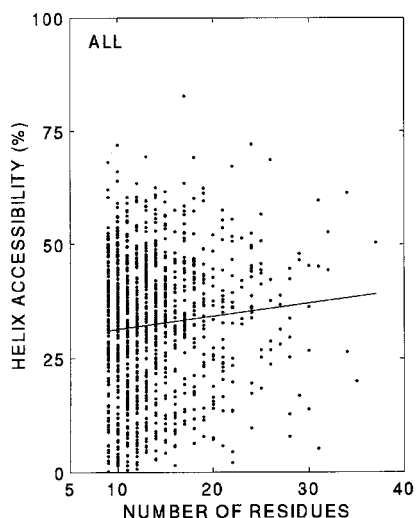


FIGURE 3 Plot showing solvent accessibility of all 1131  $\alpha$ -helices (helix accessibility) with respect to number of residues in the helices. The best fit line indicates that there is no significant correlation between helix accessibility and its length.

### Factors affecting the helix-forming propensities of amino acids

Helix-forming propensities of various amino acids in barnase have been analyzed by Fersht and co-workers (Serrano and Fersht, 1989; Serrano et al., 1992) with site-directed mutagenesis experiments, and they concluded that assigning a unique and generally applicable helix-forming propensity to an amino acid may not be valid. Chakrabarty et al. (1991) also found a strong positional dependence for Ala $\rightarrow$ Gly substitutions in alanine-based peptides. The amino acid distributions in the middle of helices are different from those at helix termini (Argos and Palau, 1982; Richardson and Richardson, 1988; Kumar and Bansal, 1998). Helix propensities of amino acids vary with position of the helix in protein tertiary structure (Blundell and Zhu,

1995; Zhu and Blundell, 1996). In this report, we present further evidence for the variations observed in amino acid propensities in  $\alpha$ -helices with respect to the length and geometries of the helices, as well as their location in the protein globule.

### Amino acid compositions of $\alpha$ -helices vary with their location in the protein globule

In general, we find that all 20 amino acids are less accessible to the solvent when present in  $\alpha$ -helices as compared to their average accessibilities in the whole protein globule. This finding is consistent with the observation in the previous section that most of the  $\alpha$ -helices are buried in the protein core. The question arises whether the amino distributions in the  $\alpha$ -helices vary with their location in the protein globule. Table 1 shows  $\chi^2$  comparison and Hamming distances among amino acid distributions in various helix accessibility classes. It can be seen that differences among the amino acid compositions of helices in different accessibility classes are highly significant. Large and consistently increasing values of Hamming distances reflect that amino acid compositions of these classes are well separated in the 20-dimensional amino acid composition space. Propensities of several amino acids vary with helix accessibility. Fig. 5 shows the variations in propensities of Ala, Leu, Gln, Glu, Lys, and Arg as a function of helix accessibility. Propensities of aliphatic amino acids Ala and Leu decrease linearly with increase in helix accessibility. Proportions of Ala and Leu decrease from 15.4% and 13.0%, respectively, in helices with less than 10% helix accessibility, to 11.6% and 10.4%, respectively, in helices with 50–60% accessibility. The changes in proportions of Ala and Leu are significant at the 95% level of confidence. These observations indicate that Ala and Leu occur more frequently in a buried helix than in a helix near or on the protein surface, as expected on the basis of their hydropho-

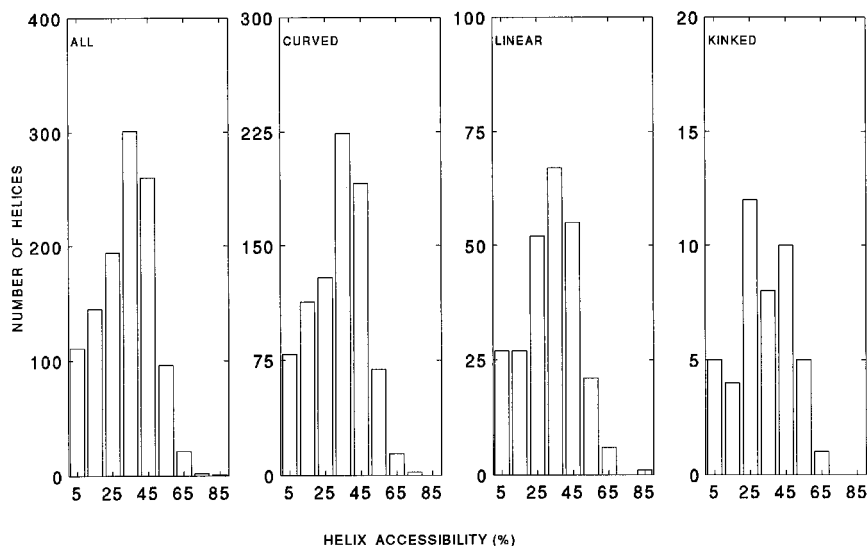


FIGURE 4 Bar diagrams showing the distributions of  $\alpha$ -helices with respect to helix accessibility in various helix geometries. The x axis denotes helix accessibility, and the y axis denotes the number of helices in each case.



acids (Tyr, Trp, or Phe) five times, and His once. The regions of kink in three helices contain both  $\beta$ -branched and aromatic residues.

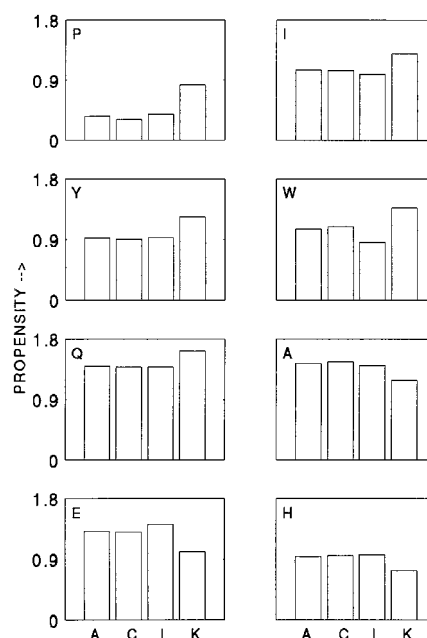
Table 2 shows  $\chi^2$  comparisons and Hamming distances for amino acid distributions in various helix classes. The differences between the amino acid distributions of linear and curved helices are not significant, whereas those between linear and kinked helices as well as between curved and kinked helices are significant, even at the 99.99% level of confidence ( $p < 0.001$ ,  $\chi^2$  values  $> 43.82$ ). These results indicate that amino acid distribution in 45 kinked helices is significantly different from the amino acid distributions in linear and curved helices.

Propensity values for the amino acids that show large differences among the various geometrical classes of helices are shown in Fig. 6. In addition to proline, several other amino acids also show appreciable changes in their propensities to occur in kinked helices. Propensities of Pro, Ile, Trp, Tyr, and Gln show an increase for kinked helices as compared to linear and curved helices, whereas propensities of Ala, His, and Glu show an appreciable decrease. Amino acid distribution in the regions of kink in the 45 kinked  $\alpha$ -helices can be useful in identifying kink-causing residues. Proportions of amino acid residues Ala (13.0% in all helices, 10.7% in kinked helices, 8.7% in the region of kink) and Glu (8.2% in all helices, 6.2% in kinked helices, 3.8% in the region of kink) decrease, whereas the proportion of Pro (1.6% in all helices, 3.7% in kinked helices, 14.1% in the region of kink) increases by large amounts. Changes in proportions of Pro and Glu in the 45 kinked helices as well as in their kinked regions with respect to their proportions in

**TABLE 2** Values of  $\chi^2$  comparison and Hamming distances in the amino acid compositions of various classes of  $\alpha$ -helices, characterized by their geometry and length

Helix classes	$\chi^2$ value	Hamming distance (D. U.)
Geometrical classes in $\alpha$ -helices		
Curved (821) vs. linear (256)	27.01	7.50
Curved (821) vs. kinked (45)	58.85	16.70
Linear (256) vs. kinked (45)	48.71	16.00
821 curved helices classified according to radius of curvature		
0 Å–25 Å (36) vs. 25 Å–50 Å (288)	188.12	19.40
25 Å–50 Å (288) vs. 50 Å–75 Å (245)	43.98	8.90
50 Å–75 Å (245) vs. 75 Å–100 Å (139)	37.42	9.70
75 Å–100 Å (139) vs. 100 Å–125 Å (67)	16.84	9.40
100 Å–125 Å (67) vs. >125 Å (46)	41.47	13.40
Helices classified according to helix length (only middle regions of helices were considered)		
Short (1046) vs. long (85)	44.98	14.80

The number of helices in each class is given in parenthesis, along with the helix class name. For a 19-parameter system such as amino acid distribution, the null hypothesis is rejected at the 95% ( $p < 0.05$ ) level of confidence, if  $\chi^2 > 30.14$ . D. U. stands for distance units. The middle region of an  $\alpha$ -helix corresponds to the residues remaining after the first and last four residues are removed.



**FIGURE 6** Bar diagrams showing variation in propensities with respect to helix geometry for those amino acids that show a large increase or decrease in kinked helices. Amino acids are denoted by a single-letter code in the upper left corner of each box. Bars along the horizontal axis correspond to A: all 1131  $\alpha$ -helices; C: 821 curved  $\alpha$ -helices; L: 256 linear  $\alpha$ -helices; and K: 45 kinked helices, respectively.

all of the 1131  $\alpha$ -helices are significant at the 95% level of confidence. Pro is well known to cause kinks in  $\alpha$ -helical structure (Woofson and Williams, 1990; Von Heijne, 1991; Sankaramakrishnan and Vishveshwara, 1992), but a decrease in the proportion of Glu in kinked helices is interesting and suggests that Glu may have a tendency to resist the occurrence of sharp bends/kinks in  $\alpha$ -helices. Given the conformational flexibility of Gly, it may be expected to have a higher propensity to occur in kinked regions of helices. In our studies, the propensity of Gly to occur in kinked  $\alpha$ -helices remains more or less equal to that in curved or linear helices, but the proportion of Gly increases from 3.3% in the complete kinked helices to 6.1% in the regions of kink, indicating a tendency for Gly to be sequestered in the kinked regions, when it does occur in kinked helices. However, the  $(\phi, \psi)$  values for these Gly residues are similar to those for other amino acids in the kink regions.

### In the case of smoothly curved $\alpha$ -helices, amino acid composition varies with the radius of curvature

Significant differences are observed in amino acid distributions for helices in radii of curvature ranges 0–25 Å and 25–50 Å, 25–50 Å and 50–75 Å, 50–75 Å and 75–100 Å, 100–125 Å, and >125 Å, as seen from Table 2. Differences in amino acid distributions are insignificant between the classes 75–100 Å and 100–125 Å. Variation in propensities

of amino acids that show appreciable differences in five radius of curvature classes, namely, 0–25 Å, 25–50 Å, 50–75 Å, 75–100 Å, and 100–125 Å are shown in Fig. 7. The propensity to occur in the complete dataset of 821 curved  $\alpha$ -helices is also plotted alongside in each figure. The helix-forming propensities of Ala, Asn, and Glu remain nearly constant for helices in all radius of curvature classes,

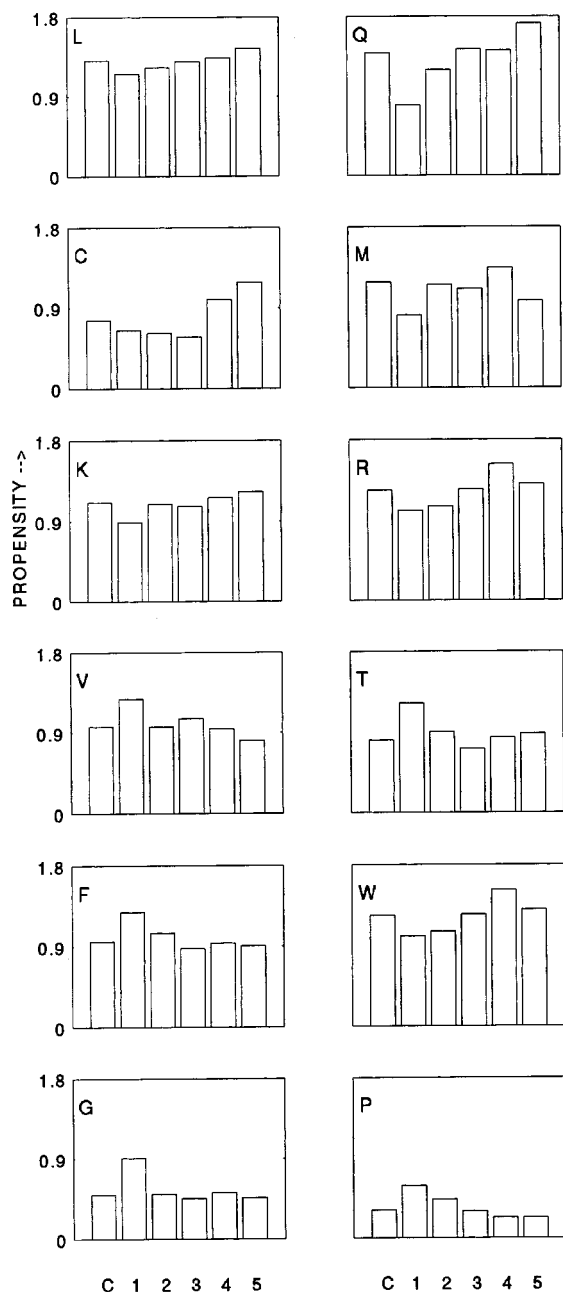


FIGURE 7 Bar diagrams showing the variation in propensities of amino acids in curved  $\alpha$ -helices with respect to radius of curvature. Only amino acids that show interesting features have been shown. They are identified by a single-letter code in the upper left corner of each box. Bars along the horizontal axis correspond to (C) all 821 curved  $\alpha$ -helices, followed by ranges of radius of curvature between (1) 0 Å and 25 Å; (2) 25 Å and 50 Å; (3) 50 Å and 75 Å; (4) 75 Å and 100 Å; and (5) 100 Å and 125 Å, respectively.

whereas Met and Arg show an increase followed by a sharp decrease for large a radius of curvature class (100–125 Å). The amino acid residues Leu, Gln, Lys, and Cys generally show increasing propensities for helices with larger radii of curvature, whereas Val, Thr, Pro, Phe, and Trp show an opposite trend. It is interesting to note that the amino acid residues Pro, Val, Phe, Thr, and Gly, which have all been found to be favored at kink regions, also have considerably higher propensities to occur in the highly curved helices (radius of curvature  $\leq 25$  Å) than in the helices with larger radii of curvature, as seen from Fig. 7. However, Ile, Tyr, and Gln, which are also often found at kinks, have very low propensities for such highly curved helices. In particular, proportions of Thr (4.6% in all 821 curved helices, 6.9% in the 36 highly curved helices,  $\Delta = +2.3\%$ ) and Gly (3.9% in all 821 curved helices, 7.2% in the 36 highly curved helices,  $\Delta = +3.3\%$ ) increase significantly in the case of highly curved helices. The proportion of Thr also decreases significantly in the case of the less curved helices with radius of curvature  $> 125$  Å (4.6% in all 821 curved helices, 2.7% in the case of the 46 less curved helices,  $\Delta = -1.9\%$ ). Given that Gly has greater conformational flexibility than other amino acids, it is possible for Gly residues to cause local bends (less severe than kinks) in  $\alpha$ -helices, thus leading to their being highly curved. It is interesting to note that Pro is also found to occur in highly curved helices, whereas it is rare in helices with large radii of curvature. The above observations clearly demonstrate curvature-dependent sequence variations in curved  $\alpha$ -helices.

#### Amino acid composition of long $\alpha$ -helices differs from that of short $\alpha$ -helices

$\alpha$ -Helices in globular proteins show length-dependent sequence variations, with amino acids having longer side chains and/or a greater number of side-chain functional groups occurring more frequently in longer  $\alpha$ -helices (Kumar and Bansal, 1996). Whereas N-cap and C-cap regions (Presta and Rose, 1988) in  $\alpha$ -helices are expected to be similar for helices of all lengths, length-dependent variations may occur in the middle regions of  $\alpha$ -helices. Table 2 shows  $\chi^2$  comparisons as well as Hamming distances for the middle regions of helices classified into two length classes, viz., short helices (1046 helices with 9–21 residues) and long helices (85 helices with  $>21$  residues). It can be seen that the differences in the distributions of amino acids in the middle regions of helices in these two length classes are significant even at the 99.99% ( $\chi^2 > 43.82$ ) level of confidence. Proportions of amino acid residues Leu ( $\Delta = +1.3\%$ ), Gln ( $\Delta = +1.1\%$ ), and Arg ( $\Delta = +1.4\%$ ) increase substantially, whereas those of Ala ( $\Delta = -2.8\%$ ), Val ( $\Delta = -2.0\%$ ), and Lys ( $\Delta = -0.9\%$ ) show an appreciable decrease for longer helices. Changes in the proportions of Ala, Arg, and Val are significant at the 95% level of confidence. Propensities of individual amino acids to occur in the middle regions of short and long  $\alpha$ -helices are given in Table 3.



**TABLE 3** Statistical propensities for individual amino acids to occur in middle regions (as defined in Table 2) of all 1131  $\alpha$ -helices, along with the propensities in the short and long classes separately

Amino acid	All helices (9–37 residues)	Short $\alpha$ -helices (9–21 residues)	Long $\alpha$ -helices (>21 residues)
Ala (A)	1.56 $\pm$ 0.05	1.63 $\pm$ 0.05	1.32 $\pm$ 0.10
Cys (C)	0.87 $\pm$ 0.09	0.76 $\pm$ 0.10	1.24 $\pm$ 0.24
Asp (D)	0.74 $\pm$ 0.04	0.72 $\pm$ 0.05	0.81 $\pm$ 0.09
Glu (E)	1.16 $\pm$ 0.05	1.18 $\pm$ 0.06	1.08 $\pm$ 0.11
Phe (F)	0.94 $\pm$ 0.06	0.98 $\pm$ 0.07	0.82 $\pm$ 0.11
Gly (G)	0.43 $\pm$ 0.03	0.41 $\pm$ 0.03	0.49 $\pm$ 0.06
His (H)	0.76 $\pm$ 0.07	0.74 $\pm$ 0.08	0.81 $\pm$ 0.15
Ile (I)	1.16 $\pm$ 0.05	1.17 $\pm$ 0.06	1.14 $\pm$ 0.12
Lys (K)	1.12 $\pm$ 0.05	1.15 $\pm$ 0.06	1.01 $\pm$ 0.10
Leu (L)	1.35 $\pm$ 0.05	1.32 $\pm$ 0.05	1.47 $\pm$ 0.11
Met (M)	1.41 $\pm$ 0.09	1.37 $\pm$ 0.10	1.51 $\pm$ 0.21
Asn (N)	0.78 $\pm$ 0.05	0.80 $\pm$ 0.06	0.72 $\pm$ 0.10
Pro (P)	0.10 $\pm$ 0.02	0.08 $\pm$ 0.02	0.13 $\pm$ 0.05
Gln (Q)	1.40 $\pm$ 0.07	1.33 $\pm$ 0.08	1.64 $\pm$ 0.17
Arg (R)	1.33 $\pm$ 0.06	1.26 $\pm$ 0.07	1.57 $\pm$ 0.15
Ser (S)	0.72 $\pm$ 0.04	0.71 $\pm$ 0.05	0.78 $\pm$ 0.09
Thr (T)	0.84 $\pm$ 0.05	0.84 $\pm$ 0.05	0.85 $\pm$ 0.10
Val (V)	1.05 $\pm$ 0.05	1.11 $\pm$ 0.05	0.82 $\pm$ 0.09
Trp (W)	0.98 $\pm$ 0.10	0.94 $\pm$ 0.11	1.10 $\pm$ 0.23
Tyr (Y)	0.94 $\pm$ 0.06	0.92 $\pm$ 0.07	1.03 $\pm$ 0.13

Amino acid propensities and their standard deviations are calculated according to formulae given by Williams et al. (1987). The 1131  $\alpha$ -helices contain a total of 6629 amino acid residues, of which 5125 residues constitute 1046 short helices and 1504 residues are present in the 85 long helices.

These results, in general, confirm the conclusions from our earlier analysis (Kumar and Bansal, 1996), which used a smaller database of short helices, i.e. among similar residues, those with the longer side chains are preferred by long  $\alpha$ -helices over others with shorter side chains (e.g., Leu is preferred over Ala and Val, Gln is preferred over Asn). In addition, residues with a greater number of functional groups have higher propensities for longer helices (e.g., Arg has a higher propensity than Lys, Tyr has a higher propensity than Phe).

### Potential applications of the results

The results presented above are based on a comprehensive analysis of a large database of  $\alpha$ -helices found in nonhomologous globular proteins. Because of the sufficiently large size of the database and the use of more than one statistical measure to characterize various properties, this study is able to avoid pitfalls and small data biases. It provides an improved understanding of the fine features of  $\alpha$ -helical structures and can lead to more accurate secondary structure prediction and better de novo design strategies.

Although several previously known structural properties of  $\alpha$ -helices have been reaffirmed in this analysis, there are some surprises, too. For example, this study establishes for the first time a strong structural relationship between helix length and curvature, as well as amino acid composition. This result is invaluable from the point of view of de novo

design, because it implies that one can regulate the curvature range of an  $\alpha$ -helix by varying its length and composition. Furthermore, one should be able to predict the approximate curvature of an  $\alpha$ -helix based on its length and sequence. This, in turn, requires an accurate estimation of helix length. This can be done by recognizing some “sequence punctuation marks” in the form of characteristic amino acid residues found to occur with high propensities at the helix termini (Kumar and Bansal, 1998). The development of sensitive prediction methods that can incorporate variations in amino acid propensities with respect to various structural parameters described here as well as with respect to the N- and C- terminal regions (Kumar and Bansal, 1998) should result in improvement in secondary as well as higher order structure prediction. For example, a high content of polar amino acid residues can indicate that the helix may lie at or near the protein surface, whereas the presence of an excess of apolar aliphatic residues like Leu can be an indication that the helix is long and is probably located in the core of the protein globule. The presence of  $\beta$ -branched and aromatic residues, in addition to Gly and Pro in the middle of helices, would suggest that the helix either has high curvature or is kinked. Similarly, these results may be very useful in the de novo design of protein motifs based on  $\alpha$ -helical structures. If the designed helix should be long, one may consider using Leu, Arg, and Tyr instead of Ala, Val, Lys, and Phe. To generate  $\alpha$ -helix with high curvature, residues like Thr and Gly should be preferred.

### CONCLUSIONS

The present analysis of 1131  $\alpha$ -helices, with lengths varying from 9 to 37 residues, shows that geometries of nearly all  $\alpha$ -helices in the globular proteins can be simply characterized as being linear, curved, or kinked. Most of the  $\alpha$ -helices (~73%) show varying degrees of smooth curvature, with longer helices generally tending to be less curved. A localized sharp kink appears to be the preferred mode of severe/high bending for longer helices, as indicated by the fact that mean length of the kinked helices is  $20 \pm 6$ , whereas that of linear and curved helices is  $12 \pm 3$  and  $14 \pm 5$  residues, respectively.

$\alpha$ -Helices of various lengths or geometries do not segregate into different regions of a protein globule. Structural features of  $\alpha$ -helices, such as length and geometry, do not correlate with the thermal fluctuations in the coordinates of their constituent atoms. Rather, these features are inherent properties of the helices, encoded in their amino acid sequences. Sequence compositions and propensities of individual amino acids to occur in  $\alpha$ -helices vary with helix length, geometry, curvature, and location in the protein globule. Residues with polar side chains, viz., Gln, Glu, Lys, and Arg, are all significantly more frequent in helices that occur at or near the protein surface. Very significant differences are seen between the amino acid composition of 45 kinked  $\alpha$ -helices and those of linear or smoothly curved

helices. Helix-breaking residues (such as  $\beta$ -branched and aromatic residues, in addition to Gly and Pro) are often found in the region of kink, suggesting a thermodynamic correlation. Gly, Pro, Val, Thr, and Phe are also found to occur, with more than their expected frequencies, in the highly curved helices. Among amino acids with similar side chains, those with longer and/or a greater number of functional groups are found much more frequently in longer helices (viz. Leu > Ala/Val, Gln > Asn, and Arg > Lys). Hence propensities of individual amino acids to occur in a given secondary structure depend not only on its conformation but also on the length, geometry, curvature, and location of the secondary structure in the protein globule. The sequence-structure correlation observed in this analysis, along with those observed in a previous paper (Kumar and Bansal, 1998), if incorporated in protein structure prediction algorithms, should lead to better prediction of secondary structure as well as de novo design of novel motifs containing  $\alpha$ -helices with predefined length and geometry.

The authors are grateful to Prof. N. V. Joshi for advice on statistical analysis and to Dr. M. RaviKiran for several useful suggestions during the course of this investigation. The authors also thank Mr. Anirban Ghosh for help in preparation of the manuscript.

SK acknowledges the Council for Scientific and Industrial Research of India for a fellowship.

## REFERENCES

- Argos, P., and J. Palau. 1982. Amino acid distribution in protein secondary structures. *Int. J. Protein Pept. Res.* 19:380–393.
- Banner, D. W., M. Kokkinidis, and D. Tsernoglou. 1987. Structure of the ColE1 Rop protein at 1.7 Å resolution. *J. Mol. Biol.* 196:657–675.
- Barlow, D., and J. M. Thornton. 1988. Helix geometry in proteins. *J. Mol. Biol.* 201:601–619.
- Bernstein, F. C., T. F. Koetzle, G. J. B. Williams, E. F. Meyer, Jr., M. D. Brice, J. R. Rodgers, O. Kennard, T. Shimanouchi, and Tasumi, M. 1977. The Protein Data Bank: a computer based archival file for macromolecular structures. *J. Mol. Biol.* 112:535–542.
- Blundell, T. L., D. Barlow, N. Barkakoti, and J. M. Thornton. 1983. Solvent induces distortions and curvature of  $\alpha$ -helices. *Nature.* 306:281–283.
- Blundell, T. L., and Z.-Y. Zhu. 1995. The  $\alpha$  helix as seen from the protein tertiary structure: a 3-D structural classification. *Biophys. Chem.* 55:167–184.
- Chakrabarti, P., M. Bernard, and D. C. Rees. 1986. Peptide bond distortion and curvature of  $\alpha$  helices. *Biopolymers.* 25:1087–1093.
- Chakrabarty, A., J. A. Schellman, and R. L. Baldwin. 1991. Large differences in the helix propensities of alanine and glycine. *Nature.* 351:586–588.
- Chou, K. C., and C. T. Zhang. 1995. Prediction of protein structural classes. *Crit. Rev. Biochem. Mol. Biol.* 30:275–349.
- Creighton, T. E. 1993. *Proteins: Structure and Molecular Properties*, 2nd Ed. W. H. Freeman and Company, New York.
- Hobohm, U., and C. Sander. 1994. Enlarged representative set of protein structures. *Protein Sci.* 3:522–524.
- Hobohm, U., M. Scharf, R. Schneider, and C. Sander. 1992. Selection of representative protein data sets. *Protein Sci.* 1:409–417.
- Kabsch, W., and C. Sander. 1983. Dictionary of protein secondary structure: pattern recognition of hydrogen bonded and geometrical features. *Biopolymers.* 22:2577–2637.
- Kendrew, J. C., G. Bodo, H. M. Dintzis, R. G. Parrish, H. Wycoff, and D. C. Phillips. 1958. A three dimensional model of the myoglobin molecule obtained by x-ray analysis. *Nature.* 181:662–666.
- Kumar, S., and M. Bansal. 1996. Structure and sequences characteristics of long  $\alpha$  helices in globular proteins. *Biophys. J.* 71:1574–1586.
- Kumar, S., and M. Bansal. 1998. Dissecting  $\alpha$  helices: position specific analysis of  $\alpha$  helices in globular proteins. *Proteins Struct. Funct. Genet.* 31:460–476.
- Lee, B., and F. M. Richards. 1971. The interpretation of protein structures: estimation of static accessibility. *J. Mol. Biol.* 55:379–400.
- Medhi, J. 1992. *Statistical Methods: An Introductory Text*. Wiley Eastern Limited, New Delhi.
- Pauling, L., R. B. Corey, and H. R. Branson. 1951. The structure of proteins: two hydrogen bonded helical configurations of the polypeptide chains. *Proc. Natl. Acad. Sci. USA.* 37:205–211.
- Presta, L. G., and G. D. Rose. 1988. Helix signals in proteins. *Science.* 240:1632–1641.
- Richardson, J. S., and D. C. Richardson. 1988. Amino acid preferences of specific locations at the ends of  $\alpha$  helices. *Science.* 240:1648–1652.
- Sankaramakrishnan, R., and S. Vishveshwara. 1992. Geometry of proline-containing alpha-helices. *Int. J. Pept. Protein Res.* 39:356–363.
- Serrano, L., and A. R. Fersht. 1989. Capping and  $\alpha$  helix stability. *Nature.* 342:296–299.
- Serrano, L., J. L. Neira, J. Sancho, and A. R. Fersht. 1992. Effect of alanine versus glycine in  $\alpha$  helices on protein stability. *Nature.* 256:453–455.
- Srinivasan, R. 1976. Helical length distribution from protein crystallographic data. *Ind. J. Biochem. Biophys.* 13:192–193.
- Von Heijne, G. 1991. Proline kinks in transmembrane alpha-helices. *J. Mol. Biol.* 218:499–503.
- Vorobjev, Y. N., and J. Hermans. 1997. SIMS: computation of a smooth invariant molecular surface. *Biophys. J.* 73:722–732.
- Williams, R. W., A. Chang, D. Juretic, and S. Loughran. 1987. Secondary structure predictions and medium range interactions. *Biochim. Biophys. Acta.* 916:200–204.
- Wolfson, D. N., and D. H. Williams. 1990. The influence of proline residues on alpha-helical structure. *FEBS Lett.* 277:185–188.
- Zhu, Z.-Y., and T. L. Blundell. 1996. The use of amino acid patterns of classified helices and strands in secondary structure prediction. *J. Mol. Biol.* 260:261–276.