# Analysis on sliding helices and strands in protein structural comparisons: A case study with protein kinases

V S Gowri, K Anamika, S Gore[1] and N Srinivasan[*]

*Molecular Biophysics Unit, Indian Institute of Science, Bangalore 560 012, India*
[1]*Present address: Department of Biochemistry, University of Cambridge, 80, Tennis Court Road, Cambridge CB2 1GA, UK*

*[*]Corresponding author (Fax, 91-80-2360 0535; Email, ns@mbu.iisc.ernet.in)*

Protein structural alignments are generally considered as 'golden standard' for the alignment at the level of amino acid residues. In this study we have compared the quality of pairwise and multiple structural alignments of about 5900 homologous proteins from 718 families of known 3-D structures. We observe shifts in the alignment of regular secondary structural elements (helices and strands) between pairwise and multiple structural alignments. The differences between pairwise and multiple structural alignments within helical and $\beta$-strand regions often correspond to 4 and 2 residue positions respectively. Such shifts correspond approximately to "one turn" of these regular secondary structures. We have performed manual analysis explicitly on the family of protein kinases. We note shifts of one or two turns in helix-helix alignments obtained using pairwise and multiple structural alignments. Investigations on the quality of the equivalent helix-helix, strand-strand pairs in terms of their residue side-chain accessibilities have been made. Our results indicate that the quality of the pairwise alignments is comparable to that of the multiple structural alignments and, in fact, is often better. We propose that pairwise alignment of protein structures should also be used in formulation of methods for structure prediction and evolutionary analysis.

## 1. Introduction

It is well-known that the 3-D structures of homologous proteins are conserved better than the sequence of amino acids (Chothia and Lesk 1986; Hubbard and Blundell 1987). Even if the sequence similarity between the homologues is extremely poor (less than about 25% of sequence identity) they share a common 3-D fold and, often, gross functional properties are similar (Holm and Sander 1994; Murzin *et al* 1995; Orengo *et al* 1997). Clearly, for the homologues of very distant relationship, sequence-based alignment is unreliable and alignment obtained on the basis of 3-D structural similarities is more reliable. Hence, in general, protein structure-based alignments are more accurate than sequence similarity-based alignments.

Protein structural alignments are in extensive use over a long time in fold recognition, comparative modelling and in various analyses on protein evolution. The quality of the structural alignments plays an important role in such studies. Several structural comparison algorithms have been developed over the last few decades and some of these developments are commonly used. For example, DALI is a pairwise structural alignment program which compares protein structures based on the alignment of distance matrices generated based on the C$\alpha$-C$\alpha$ distance (Holm and Sander 1993). One of the widely used programs

for superposition of multiple structures is STAMP (Russell and Barton 1992) which employs rigid-body superposition algorithm. MUSTANG (Konaguruthu *et al* 2006) is another multiple structural superposition algorithm that is based on progressive pairwise heuristic algorithm. In the COMPARER approach, Sali and Blundell (1990) use structural features such as solvent accessibility and secondary structure and relationships such as hydrogen bonding and employed dynamic programming in aligning protein structures. Orengo and Taylor (1990) use iterative double dynamic programming approach progressively. Gerstein and Levitt (1996) use an approach that alternates between dynamic programming and rigid-body superposition. TOPS is a pattern search algorithm for superposition of distantly related multiple protein structures (Williams *et al* 2003). However, this algorithm is limited to all $\beta$ protein structures and performs less effectively for the other classes of protein structures namely all alpha, alpha/beta and alpha+beta in comparison to STAMP.

It is well understood that the multiple sequence alignments are informative compared to pairwise sequence alignments. Residue variations in positions in multiple sequence alignment serve as excellent guideline in arriving at the accurate alignment. While multiple structural alignments enable identification of conserved core of a protein fold it is not clear if the quality of the alignment is better than the pairwise structural alignment. Here, we present an analysis on the comparison of the quality of the pairwise structural alignments and multiple structural alignments following the indications obtained in our earlier broader analysis of alignment of protein structures (Balaji and Srinivasan 2001). As loops are known to be variable among homologous proteins we considered the alignment involving regular secondary structural elements namely helices and strands in our analysis. Similarities in solvent accessibilities of residues aligned in pairwise and multiple structural alignments have been used as diagnostic features in assessing the alignment quality. A detailed analysis on structural alignment comparison has been presented for one of the most divergent family of proteins which is the protein kinase family.

## 2.    Materials and methods

### 2.1    Dataset

The dataset for the comparison of structural alignments include 718 multi-member protein family alignments from the PALI database (Balaji *et al* 2001; Gowri *et al* 2003) spread across four SCOP (Murzin *et al* 1995) classes namely, all alpha, all beta, alpha/beta and alpha+beta. For the analysis on protein kinase structural comparisons, the 3-D structures of kinases in their active conformation are considered, as kinases are known to adopt characteristic conformation when they are active (Krupa *et al* 2004) and

alter the conformation significantly when a kinase is in inactive state (Huse and Kuriyan 2002). The six protein kinase structures considered for the detailed analysis are, cAMP-dependent protein kinase [1ATP, (Zheng *et al* 1993)], cyclin-dependent protein kinase [1JST, (Russo *et al* 1996)], phosphorylase kinase [1PHK, (Owen *et al* 1995)], PKB/Akt protein kinase [1O6L, (Yang *et al* 2002)], protein kinase C $\theta$ [1XJD,(Xu *et al* 2004)] and mitogen activated protein kinase [2ERK, (Canagarajah *et al* 1997)].

### 2.2    Softwares

Pairwise and multiple structural alignments are generated using STAMP (structural alignment of multiple proteins) program (Russell and Barton 1992), which encodes rigid-body superposition procedure. SSTRUC (Smith 1989) program has been used to map the secondary structural elements in these protein structures. Residue accessibility calculations were performed using NACCESS (Hubbard and Thornton 1993).

### 2.3    Approach

For every possible pairs within a family the direct pairwise alignment (DPA) is compared with the pairwise alignment extracted from the multiple structural alignments (PMA). In the former case (DPA) alignment is made by considering only two proteins at a time and the latter case (PMA) corresponds to extraction of alignment of two proteins from a simultaneous alignment of all the homologues in the family. In every DPA and PMA, the regular secondary structural regions namely $\alpha$-helices and $\beta$-strands are particularly investigated in detail to analyse shifts in helix-helix and strand-strand alignments. Root mean square deviations (RMSD) of the relative side-chain accessibilities of the residues in the equivalent secondary structures have been used as a measure of accuracy of the alignment. Relative sidechain accessibility of a residue in a protein is expressed in percentage and it is the ratio of accessible surface area of the residue ($X$) in the protein to its accessible surface area in the Ala-$X$-Ala tripeptide in the extended conformation. Highly corresponding side chain accessibilities of aligned residues is considered to indicate an alignment of good quality. Root mean square deviations are calculated for every helix-helix pair and strand-strand pair from both the DPA and PMA using the following formula:

$$\mathrm{RMSD_{acc.}} = \sqrt{\frac{\sum\limits_{res=1,N} (\mathrm{RSA_{res,str1}} - \mathrm{RSA_{res,str2}})^2}{N}}$$

where, $\mathrm{RMSD_{acc.}}$ is the root mean square deviation of the amino acid side-chain accessibility of the secondary

structural elements (helices and strands), $RSA_{res,str1}$ and $RSA_{res,str2}$ refer to the amino acid side-chain accessibility values of the equivalent residues in structure 1 and structure 2 respectively. $N$ is the total number of residues in a particular secondary structural element.

Normalized positional shift (NPS) between DPA and PMA for a pair of proteins is calculated using

$$NPS = \frac{\sum_i S_i}{\sum_i R_i}$$

where the summation $i$ is made for the number of helix-helix or strand-strand pairs involved in the alignment, $S_i$ represents the sum of number of positions shifted in PMA compared to DPA for every residue pair present in the ith pair of regular secondary structure and $R_i$ represents number of residues in the regular secondary structure, in the $i$th pair, corresponding to the shorter of the two proteins involved in the alignment.

## 3.    Results and discussion

### 3.1    *Comparison of pairwise and multiple structural alignments*

We have compared the quality of direct DPA with the pairwise alignments extracted from the PMA of the homologous protein structures. As loops are known to adopt different conformations in homologous proteins we focused on possible differences in alignment within helical or $\beta$-strand regions. As expected in many of the helix-helix and strand-strand pairs in our dataset, the alignments from DPA and PMA matched completely. However a substantial proportion of pairs of homologous proteins showed differences between DPA and PMA. Whenever DPA and PMA differed for a given pair of protein structures we used root mean square deviation of the side-chain solvent accessibility of the aligned residues in equivalent helices and strands to assess the quality of DPA and PMA. RMSD of the accessibility values ($RMSD_{acc}$) of all the helices and strands are compared between DPA and PMA for every pair of proteins involved in the analysis.

Percentage of pairs of helices with the $RMSD_{acc}$ from DPA better than PMA [i.e. $RMSD_{acc}$ (DPA) < $RMSD_{acc}$ (PMA)] and those with $RMSD_{acc}$ from PMA better than DPA [$RMSD_{acc}$ (DPA) > $RMSD_{acc}$ (PMA)] are plotted for the three classes namely all alpha, alpha/beta and alpha+beta as shown in figure 1a. Similar bar diagram is plotted for the beta strands for the three classes namely all beta, alpha/beta and alpha+beta (figure 1b). The plots suggest that, in general, there is a higher proportion of pairs of regular secondary structural elements (helices and strands) with $RMSD_{acc}$ for DPA better (lower $RMSD_{acc}$) than that of PMA.

This comparison clearly demonstrates that, in general, the pairwise structural alignment of the protein structures is more optimal compared to the multiple structural alignment.

This result is in contrast to the conclusions from sequence analysis wherein the multiple sequence alignment is certainly far more accurate and robust than pairwise alignment.

### 3.2    *Analysis of difference between DPA and PMA*

Difference between PMA and DPA is unlikely to be pronounced if the extent of evolutionary divergence in a protein family is low which is indicated by high sequence identities between pairs of homologous proteins in the family. In order to explore this proposition we analysed the relationship between sequence identities between homologous proteins in a pair and extent of difference between DPA and PMA as indicated by normalized positional shift (NPS). Figure 2 shows a plot between sequence identity and NPS with both the parameters averaged for every 5% of sequence identity between homologues in a pair. It is clear that with increase in the sequence identity between homologues the difference between DPA and PMA decreases. Hence, it is particularly important to use pairwise alignments between homologues compared to PMAs when the evolutionary divergence in the family is high.

We have also analysed the nature of difference between DPA and PMA for all the pairs showing difference between the two kinds of alignment. In particular, we analysed the number of residue positions the PMA is shifted compared to the corresponding DPA. Figure 3 shows the frequency of helix-helix and strand-strand pairs with non-zero shift in the alignment positions between DPA and PMA. Helix-helix alignments often differ between DPA and PMA by 4 residue positions which corresponds approximately to a turn of the helix. In the alignment between two homologues with a topologically equivalent pair of helices differing in helix lengths by 3 to 6 residues, the shorter helix may "slide" over the longer helix by one turn resulting in comparable quality of root mean square deviation of C$\alpha$ atoms. Such a "slide" seems to be the case often between DPA and PMA. From figure 3 it can also be noted that the difference between DPA and PMA for the strand regions is often corresponding to two or one alignment positions. While the ideal fully extended strand may be viewed as a "helix" with two residues per "turn" the extended strands in proteins deviate markedly from the ideal nature and number of "turns" in such structures is usually one to two. The differences between DPA and PMA for the strand regions seem to be consistent with this notion and hence correspond to sliding "turns" in extended strands.

In order to understand the origin of such shifts in the aligned regular secondary structures, the difference between the lengths of the equivalent secondary structures

**Figure 1.**    Percentage of pairs of secondary structural elements for various SCOP classes (a) alpha helices (b) beta strands. Number of pairs with better DPA compared to PMA (and *viceversa*) as indicated by the RMSD of the solvent accessibilities of the side chains of aligned residues are indicated.

has been analysed. A plot of the frequency of helix-helix and strand-strand pairs with non-zero shifts as a function of the difference between the lengths of the equivalent secondary structures suggests that there is considerable

number of misaligned helix-helix (~3000) and strand-strand pairs (~7000 pairs) with zero length difference between them (figure 4). Further, a plot of the frequency of equivalent helix-helix pairs with zero length differences

**Figure 2.** Plot of normalized positional shift between DPA and PMA against sequence identity. Both X and Y-axes parameters are averaged for every 5% of sequence identity.



**Figure 4.** Plot of frequency of helix-helix and strand-strand alignment segments against the length difference between the segments.



**Figure 3.** Plot of frequency of helix-helix and strand-strand alignment segments against shift in the alignment positions between DPA and PMA.



**Figure 5.** Plot of frequency of helix-helix alignment segments with zero length difference against the shift in the alignment positions between DPA and PMA.

and their shift values suggests that there are approximately 1800 misaligned pairs with a shift of 4 residues (figure 5). However, this proportion is considerably less compared to the proportion of helix-helix pairs with a shift of 4 residues (~11000 pairs) as shown in figure 3.

### 3.3 *Structural comparison of protein kinases*

Comparison of the pairwise structural alignments with those of multiple structural alignments for the homologous protein families from the PALI database has been extended to one of the most diverse family of proteins which is the family of protein kinases. In four out of the 15 pairs of structural alignments,

shift in the secondary structural elements have been observed. Two such examples are discussed in detail below.

3.3.1 *Structural comparison of cyclin-dependent protein kinase and PKB/Akt protein kinase:* Cyclin-dependent protein kinase (1JST) and PKB/Akt protein kinase (1O6L) share gross structural similarity with $C_\alpha$ RMSD of 1.87Å between the structures. Pairwise structural comparisons from the DPA with that of the PMA however differs drastically. Comparison of the pairwise structural alignments from DPA and PMA show shifts in the alignments of regular secondary structural elements. These differences are reflected in the

**Figure 6.** Structural superposition of helices from cyclin-dependent protein kinase (red) with PKB/Akt protein kinase (blue). Root mean square deviation of the side-chain accessibilities of the helices (RMSD$_{acc}$) and the structure dependent sequence alignment blocks from the direct pairwise represented as DPA and pairwise extracted from the multiple structural alignment represented as PMA. This figure and figure 5 have been produced using Setor (Evans, 1993).



**Figure 7.** **(a)** Structural superposition of PKB/Akt protein kinase (blue) and Phosphorylase kinase (orange) with the four helices ($\alpha$B, $\alpha$D, 3$_{10}$ B, C and $\alpha$F) are marked. **(b)** The structural superposition of the four helices ($\alpha$B, $\alpha$D, 3$_{10}$ B, C and $\alpha$F) from the PKB/Akt protein kinase structure with the Phosphorylase kinase structure with the root mean square deviation of the side-chain accessibilities of the respective helices (RMSD$_{acc}$) for the direct pairwise alignment (DPA) and pairwise extracted from the multiple structural alignment (PMA).

pairwise structure-based sequence identities calculated for DPA is about 25% and for PMA is about 11%. The structural superposition of the two helices along with the alignment blocks and the $RMSD_{acc}$ calculated for the alignments extracted from DPA and PMA are shown in figure 6. Shift of 4-7 residues has been observed in two topologically equivalent helices in this structural alignment comparison. A 3 and 7 residue shift in the helices correspond approximately to one and two turns of the helices. Repetitive nature of local structural motifs in a regular secondary structure, such as a helix, is a contributing factor to the difference in the alignments.

Comparison of the side-chain accessibilities of the aligned residues in these two helices from the DPA and PMA suggest that the alignment of the two helices from CDK structure with the equivalent helices from the PKB/Akt structure is more optimal in the DPA compared to the PMA. This is numerically represented by the root mean square deviation of the side-chain accessibility ($RMSD_{acc}$) of the amino acids in these helices shown in figure 6.

3.3.2 *Structural comparison of PKB/Akt protein kinase and phosphorylase kinase:* PKB/Akt protein kinase and phosphorylase kinase share sequence identity of 31% and their structures are highly similar with RMSD of 1.17Å for optimal match of $C_\alpha$ positions. Pairwise structural superposition of these two protein kinase structures is shown in figure 7a with the four helices labelled according to the convention used in cAMP-dependent protein kinase structure (Knighton *et al* 1991). Structural superposition from DPA and PMA of these four helices showing the shift of 4–7 residues are shown in figure 7b.

Comparison of the RMSD of the side-chain accessibilities of the four helices between DPA and PMA are 18.2 (DPA), 22.74 (PMA) for $\alpha$B helix; 24.89 (DPA), 26.53 (PMA) for $\alpha$D helix; 1.47 (DPA), 6.99 (PMA) for $3_{10}$ B, C helix and 6.41 (DPA), 9.96 (PMA) for $\alpha$F helix suggests that the $RMSD_{acc}$ of the alignment of the topologically equivalent secondary structures from DPA is less than that of PMA. Hence, the structural alignment extracted from the DPA is more optimal compared to that from PMA.

## 4. Conclusions

Structure-based alignments are generally considered as the 'golden standard' as they are more accurate than sequence-based alignments especially between distant homologues. Structure-based alignments are used frequently in fold recognition, secondary structure prediction, profile generation etc. It is well known that the multiple sequence alignments are more accurate compared to the pairwise sequence alignments. Our comparison of the structural alignments suggests that, in general, pairwise structural

alignment is better than the multiple structural alignments. Certainly DPA and PMA are of comparative quality but with critical differences. Variations in the alignments generated using different algorithms can be expected especially in the case of distantly related proteins. However the part of present analysis on protein kinase alignments has been performed manually and investigated case-by-case basis in order to ensure the robustness of structural alignment made by STAMP program.

Comparison of DPA and PMA for the protein kinase structures in their active conformation showed shift in the alignments of regular secondary structural elements in spite of gross structural similarity between the kinase structures. Such differences in the alignment quality have a profound influence in our understanding of conformational differences of active and inactive forms of protein kinases. Detailed analysis of kinase structural comparisons showed that, in general, DPA is found to be better than PMA.

The repetitive nature (regularity) in the helical and extended strand structures is a reason for the differences between different structure-based alignments (PMA and DPA) obtained using rigid-body superposition. In this sense, repetitive nature of the regular secondary structural regions may be considered analogous to the low complexity regions in protein sequences wherein one or more residue types occur in more than usual frequency over a stretch. Environmental information encoded in solvent accessibility measure could be used to assess the results of different alignments obtained using rigid-body superposition.

Our analysis suggests that, pairwise structural alignments should be considered, in addition to multiple structural alignment, to improve the purpose of use of structure-based alignments. Use of pairwise alignments of proteins can thus be expected to improve prediction protocols such as secondary structure prediction and fold recognition.

## References

Balaji S and Srinivasan N 2001 Use of a database of structural alignments and phylogenetic trees in investigating the relationship between sequence and structural variability among homologous proteins; *Prot. Eng.* **14** 219–226

Balaji S, Sujatha S, Kumar S S C and Srinivasan N 2001 PALI: A database of Phylogeny and ALIgnment of homologous protein structures; *Nucleic Acids Res.* **29** 61–65

Canagarajah B J, Khokhlatchev A, Cobb M H and Goldsmith E J 1997 Activation mechanism of the MAP kinase ERK2 by dual phosphorylation; *Cell* **90** 859–869

Chothia C and Lesk M 1986 The relation between the divergence of sequence and structure in protein; *EMBO J.* **5** 823–826

Evans S V 1993 SETOR: hardware lighted three-dimensional solid model representations of macromolecules; *J. Mol. Graph.* **11** 134–138

Gerstein M and Levitt M 1996 Using iterative dynamic programming to obtain accurate pairwise and multiple alignments of protein structures; in *Proceedings of the Fourth International Conference on Intelligent Systems for Molecular Biology.* (AAAI Press, USA) pp 59–67

Gowri V S, Pandit S B, Karthik P S, Srinivasan N and Balaji S 2003 Integration of related sequences with protein three-dimensional structural families in an updated version of PALI database; *Nucleic Acids Res.* **31** 486–488

Holm L and Sander C 1993 Protein structure comparison by alignment of distance matrices; *J. Mol. Biol.* **233** 123–138

Holm L and Sander C 1994 The FSSP database of structurally aligned protein fold families; *Nucleic Acids Res.* **22** 3600–3609

Hubbard S J and Thornton J M 1993 *'NACCESS' Computer Program* (Department of Biochemistry and Molecular Biology, University College London, UK)

Hubbard T J and Blundell T L 1987 Comparison of solvent-inaccessible cores of homologous proteins: definitions useful for protein modeling; *Protein Eng.* **1** 59–71

Huse M and Kuriyan J 2002 The conformational plasticity of protein kinases; *Cell* **109** 275–282

Knighton D R, Zheng J H, Ten Eyck L F, Xuong N H, Taylor S S and Sowadski J M 1991 Structure of a peptide inhibitor bound to the catalytic subunit of cyclic adenosine monophosphate-dependent protein kinase; *Science* **253** 414–420

Krupa A, Preethi G and Srinivasan N 2004 Structural modes of stabilization of permissive phosphorylation sites in protein kinases: distinct strategies in Ser/Thr and Tyr kinases; *J Mol Biol.* **339** 1025–1039

Konagurthu A S, Whisstock J C, Stuckey P J and Lesk A M 2006 MUSTANG: A multiple structural alignment algorithm; *Proteins Struct. Funct. Bioinformatics* **64** 559–574

Murzin A G, Brenner S E, Hubbard T and Chothia C 1995 SCOP: a structural classification of proteins database for the investigation of sequences and structures; *J. Mol. Biol.* **247** 536–540

Orengo C A and Taylor W R 1990 A rapid method for protein structure alignment; J. *Theor. Biol.* **147** 517–551

Orengo C A, Michie A D, Jone S, Jones D T, Swindells M B and Thornton J M 1997 CATH-a hierarchic classification of protein domain structures; *Structure* **5** 1093–1108

Owen D J, Noble M E, Garman E F, Papageorgiou A C and Johnson L N 1995 Two structures of the catalytic domain of phosphorylase kinase: an active protein kinase complexed with substrate analogue and product; *Structure* 3 467–482

Russell R B and Barton G B 1992 Multiple protein sequence alignment from tertiary structure comparison: assignment of global and residue confidence levels; *Proteins Struct. Funct. Genet.* **14** 309–323

Russo A A, Jeffrey P D and Pavletich N P 1996 Structural basis of cyclin-dependent kinase activation by phosphorylation; *Nat. Struct. Biol.* **3** 696–700

Sali A and Blundell T L 1990 Definition of general topological equivalence in protein structures. A procedure involving comparison of properties and relationships through simulated annealing and dynamic programming; *J. Mol. Biol.* **212** 403–428

Smith D 1989 *SSTRUC A Program to Calculate Secondary Structural Summary* (Department of Crystallography, Birkbeck College, University of London, UK)

Williams A, Gilbert D R, Westhead D R. 2003 Multiple structural alignment for distantly related all beta structures using TOPS pattern discovery and simulated annealing. *Protein Eng.* **16** 913–923

Xu Z B, Chaudhary D, Olland S, Wolfrom S, Czerwinski R, Malakian K, Lin L, Stahl M L *et al* 2004 Catalytic domain crystal structure of protein kinase C-theta (PKCtheta); *J.Biol.Chem.* **279** 50401–50409

Yang J, Cron P, Good V M, Thompson V, Hemmings B A and Barford D 2002 Crystal Structure of an Activated Akt/Protein Kinase B Ternary Complex with Gsk-3 Peptide and AMP-Pnp; *Nat. Struct. Biol.* **9** 940–944

Zheng, J H, Trafny E A, Knighton, D R, Xuong N H, Taylor S S, Teneyck F and Sowadski. J M 1993 2.2 Angstrom refined crystal-structure of the catalytic subunit cAMP-dependent protein kinase complexed with Mn ATP and a peptide inhibitor; *Acta Crystallogr. D. Biol. Crystallogr.* **49** 362–365