



BIOLOGICAL  
CRYSTALLOGRAPHY

**Volume 71 (2015)**

**Supporting information for article:**

**Identification of local variations within secondary structures of proteins**

**Prasun Kumar and Manju Bansal**

## S1. Brief description of ASSP algorithm

ASSP first calculates the local structure parameters like twist, rise, virtual torsion angle and radius by considering four contiguous C<sup>α</sup> atoms (repeating unit or step) with a sliding window of one C<sup>α</sup> atom at a time along the length of the protein structure. Subsequently, the continuity in the protein structure is checked based on these parameters and continuous stretches are further divided into the SSEs.

### S1.1. Calculation of the parameters

Let CA<sub>1</sub>(x<sub>1</sub>, y<sub>1</sub>, z<sub>1</sub>), CA<sub>2</sub>(x<sub>2</sub>, y<sub>2</sub>, z<sub>2</sub>), CA<sub>3</sub>(x<sub>3</sub>, y<sub>3</sub>, z<sub>3</sub>) and CA<sub>4</sub>(x<sub>4</sub>, y<sub>4</sub>, z<sub>4</sub>) are the Cartesian coordinates of four contiguous C<sup>α</sup> atoms (Supplementary Fig. S9). Two successive points are joined to each other by pseudo bonds and can be given as:

$$\textcircled{\text{O}} \mathbf{B1} = \mathbf{CA1CA2} = (x_2 - x_1)\mathbf{i} + (y_2 - y_1)\mathbf{j} + (z_2 - z_1)\mathbf{k}$$

$$\textcircled{\text{O}} \mathbf{B2} = \mathbf{CA2CA3} = (x_3 - x_2)\mathbf{i} + (y_3 - y_2)\mathbf{j} + (z_3 - z_2)\mathbf{k}$$

$$\textcircled{\text{O}} \mathbf{B3} = \mathbf{CA3CA4} = (x_4 - x_3)\mathbf{i} + (y_4 - y_3)\mathbf{j} + (z_4 - z_3)\mathbf{k}$$

Where, **i**, **j**, **k** are unit vectors along X, Y and Z axes respectively. The vectors obtained by subtracting the pseudo bond vectors are:

$$\textcircled{\text{O}} \mathbf{V1} = \mathbf{B1} - \mathbf{B2}$$

$$\textcircled{\text{O}} \mathbf{V2} = \mathbf{B2} - \mathbf{B3}$$

**V1** and **V2** lie in a plane perpendicular to the axis of a helix described by these four atoms. The direction cosines (*l*, *m*, *n*) of the helix axis **U** were obtained from the cross product of vectors **V1** and **V2**:

$$\textcircled{\text{O}} \mathbf{U}(l, m, n) = (\mathbf{V1} \times \mathbf{V2}) / (|\mathbf{V1} \times \mathbf{V2}|)$$

Finally, various geometric parameters were derived using following equations:

$$\textcircled{\text{O}} \text{Twist } (\mathbf{t}) = \cos^{-1} ((\mathbf{V1} \cdot \mathbf{V2}) / (|\mathbf{V1}| |\mathbf{V2}|))$$

$$\textcircled{\text{O}} \text{Vtor } (\boldsymbol{\theta}) = \text{Calculation is similar to that of } (\varphi, \psi). \text{ The only difference is that the four atoms involved are C}^\alpha. \text{ Since, the bonds between them are virtual, the torsional angle is called Vtor (Virtual torsional angle)}$$

$$\textcircled{\text{O}} \text{Rise } (\mathbf{h}) = ((\mathbf{B2} \cdot \mathbf{U}) / |\mathbf{U}|)$$

$$\textcircled{\text{O}} \text{Radius } (\mathbf{r}) = (\text{sqrt } (|\mathbf{V1}| |\mathbf{V2}|) / (2(1 - \cos^{-1} \mathbf{t})))$$

$$\textcircled{\text{O}} \text{Bending Angle } (\mathbf{BA}) = \text{It is defined as an angle between successive local helical axes corresponding to two sets of four C}^\alpha \text{ atoms CA}_1, \text{ CA}_2, \text{ CA}_3, \text{ CA}_4 \text{ and CA}_4, \text{ CA}_5, \text{ CA}_6, \text{ CA}_7. \text{ The bending angle obtained will be at CA}_4$$

### S1.2. Identification of continuous stretches

Two contiguous steps will be said to be a part of a continuous stretch, if and only if the absolute value of (Twist difference ( $\Delta\text{Twist}$ ), Rise difference ( $\Delta h$ ), Vtor difference ( $\Delta V_{\text{tor}}$ ))  $\leq (35^\circ, 1.1\text{\AA}, 50^\circ)$ , where

$$\textcircled{O} \Delta\text{Twist} = |\text{Twist1} - \text{Twist2}|$$

$$\textcircled{O} \Delta h = |h1 - h2|$$

$$\textcircled{O} \Delta V_{\text{tor}} = |V_{\text{tor1}} - V_{\text{tor2}}|$$

Here (Twist1, h1, Vtor1) and (Twist2, h2, Vtor2) are the structural parameters of repeating units 1 and 2 respectively.

The twist vs rise and twist difference vs rise difference plots for  $\alpha$ -helices identified by *DSSP*, *STRIDE* and *ASSP* are shown in Supplementary Fig. S10.

### S1.3. **Classifying the continuous stretches**

Continuous stretches as a whole or part of it are further classified into different type of SSEs by first assigning characters (**A/ a, G/ g, I/ i, P, S** and **U**) to individual step based on the step parameter values.

### S1.4. **Final arrangement**

No two secondary structures should be overlapping and in order to address this issue, final check of the same is done. Here we make use of minimum length criteria of different Secondary structure elements (SSEs) identified by *ASSP*. The minimum possible length for  $\pi$ ,  $\alpha$ , 310 and PPII helices are 5, 4, 3 and 3 residues respectively. The overlap between two SSEs (assuming 1st SSE is S1 and 2nd is S2) can only be of one residue at the junction. Broadly two cases are possible:

#### **Case 1: When the two SSEs are of different SSE types**

Based on the length of two SSEs, it is further divided into four different categories.

- i. **Length of S1 is the minimum possible and that of S2 is more than the minimum possible.**  
In this case, the termini of S1 remain the same, while N2 of S2 becomes N1, with the length of S2 reduced by one.
- ii. **Length of S1 is more than the minimum possible and that of S2 is the minimum possible.**  
In this case, the termini of S2 remain the same, while C2 of S1 becomes the C1 with the length of S1 reduced by one.
- iii. **Length of S1 and S2 both are more than the minimum possible.** The treatment to this case will be similar to the case (ii).
- iv. **Length of S1 and S2 both are minimum possible.** If the SSEs are of same type, both are merged to give a single SSE of the same type and in this case the bending angle is not

checked. In case, when the SSE types are different, the SSEs will be merged and the new SSE will be of the longer one.

- v. **When the length of either of S1 or S2 is less than minimum possible and other SSE is longer than the minimum possible:** The SSE with the length < minimum possible, will be merged to the longer one

### Case 2: When the two SSEs are of same types

When two secondary structure elements (SSE) are of same type and the bending angle at the overlapping residue  $\leq 60^\circ$ , the two SSEs are combined with the N1 coming from the 1st SSE and C1 from 2nd. In other case, where the bending angle at the overlapping residue  $> 60^\circ$ , we treat them as two separate SSEs and further arrangement is same as of case 1.

**Table S1** Brief description about different Secondary Structure assignment algorithms. The algorithms are divided according to the categories mentioned in the main text.

Sl. No.	Algorithm	Description	Reference
<b>Category (i)</b>			
1	<i>DSSP</i>	Detects the hydrogen-bond patterns using bond energy criterion	(Kabsch & Sander, 1983)
2	<i>STRIDE</i>	Uses ( $\phi$ , $\psi$ ) along with hydrogen bond pattern	(Frishman & Argos, 1995)
3	<i>PROSS*</i>	Uses only on the backbone dihedral angles ( $\phi$ , $\psi$ )	(Srinivasan & Rose, 1999)
4	<i>SECSTR</i>	Uses DSSP like hydrogen bond definition and was developed to identify and analyze $\pi$ -helices	(Fodje & Al-Karadaghi, 2002)
5	<i>DSSP-PPII*</i>	Identifies the PPII helices in the region not assigned as a major SSE by DSSP and gives the output in the DSSP format	(Mansiaux <i>et al.</i> , 2011)
<b>Category (ii)</b>			
6	Levitt <i>et.al.</i>	Uses distance and virtual torsion angle made by the CA atoms over a sliding window of four residues	(Levitt & Greer, 1977)
7	<i>DEFINE-S</i>	Uses only CA coordinates and compares the distance between various CAs with the distances	(Richards & Kundrot, 1988)

		in ideal SSEs	
8	<i>P-CURVE</i>	To start with, it chooses the successive repeating unit and does the analysis of mathematical analysis of protein curvature	(Sklenar <i>et al.</i> , 1989)
9	<i>P-SEA</i>	Solely based on the CA atoms. Uses three distance, one angle and one dihedral angle	(Labesse <i>et al.</i> , 1997)
10	<i>XTLSSTR*</i>	Calculates two angles and three distances for assigning SSEs. The algorithm is driven by the concept of circular dichroism (CD) of a protein in the far ultraviolet range.	(King & Johnson, 1999)
11	<i>STICK</i>	Finds a set of best fit axes and later takes the average rise of the residues along each axis	(Taylor, 2001)
12	<i>SST</i>	Uses minimum message length inference for SSEs assignment	(Konagurthu <i>et al.</i> , 2012)
<b>Category (iii)</b>			
13	<i>KAKSI</i>	Uses CA distances and backbone dihedral angles to show the concordance with the assignments found in the Protein Data Bank	(Martin <i>et al.</i> , 2005)
14	<i>PALSSE</i>	Mainly uses distance and torsion angle constraints to identify core elements and later extends them to longer segments	(Majumdar <i>et al.</i> , 2005)
15	<i>SEGNO*</i>	The CA atoms along with the backbone dihedral angles ( $\varphi$ , $\psi$ ) and the angle-distance hydrogen bond	(Cubellis <i>et al.</i> , 2005)

\* Algorithms identify PPII-helices also

**Table S2** Percentage content in-terms of amino acid residues assigned as  $\alpha$ ,  $3_{10}$ ,  $\pi$  helix and strand content by different algorithms in four mentioned datasets (HRes, MRes, LRes and NMR). KAKSI and PALSSE do not differentiate between  $\alpha$ ,  $3_{10}$  and  $\pi$  and assign them as helix only.  $\pi$  value for many of the algorithms is 0 because values are rounded to first decimal place.

	HRes				MRes				LRes				NMR			
	$\alpha$	$3_{10}$	$\pi$	$\beta$	$\alpha$	$3_{10}$	$\pi$	$\beta$	$\alpha$	$3_{10}$	$\pi$	$\beta$	$\alpha$	$3_{10}$	$\pi$	$\beta$
ASSP	35.2	3.5	0.9	22	35.7	3.4	0.8	22.6	32.1	3.6	0.8	19.4	32.5	3.7	0.8	16.1
DSSP	35.3	4.8	0	22.5	36.1	4.2	0	22.9	33.7	3.3	0	20.4	33.7	1.6	0	17.3
KAKSI	36.4			22	38			22.5	35.1			19	32.2			15.2
PALSSE	57.6			23.2	57.3			24.1	54.6			22.8	54.6			20
SST	34.8	1.5	0.5	22.3	35.3	1.5	0.4	23	33.1	1.7	0	20.7	30.8	2.2	0.5	18.2
STRIDE	36.5	5	0	22.6	37.3	4.4	0	23.3	34.3	3.5	0.3	21.2	35.1	2	0	18.8

**Table S3** Mean (std. dev.) values of twist, rise and radius for the PPII helices assigned by different algorithms

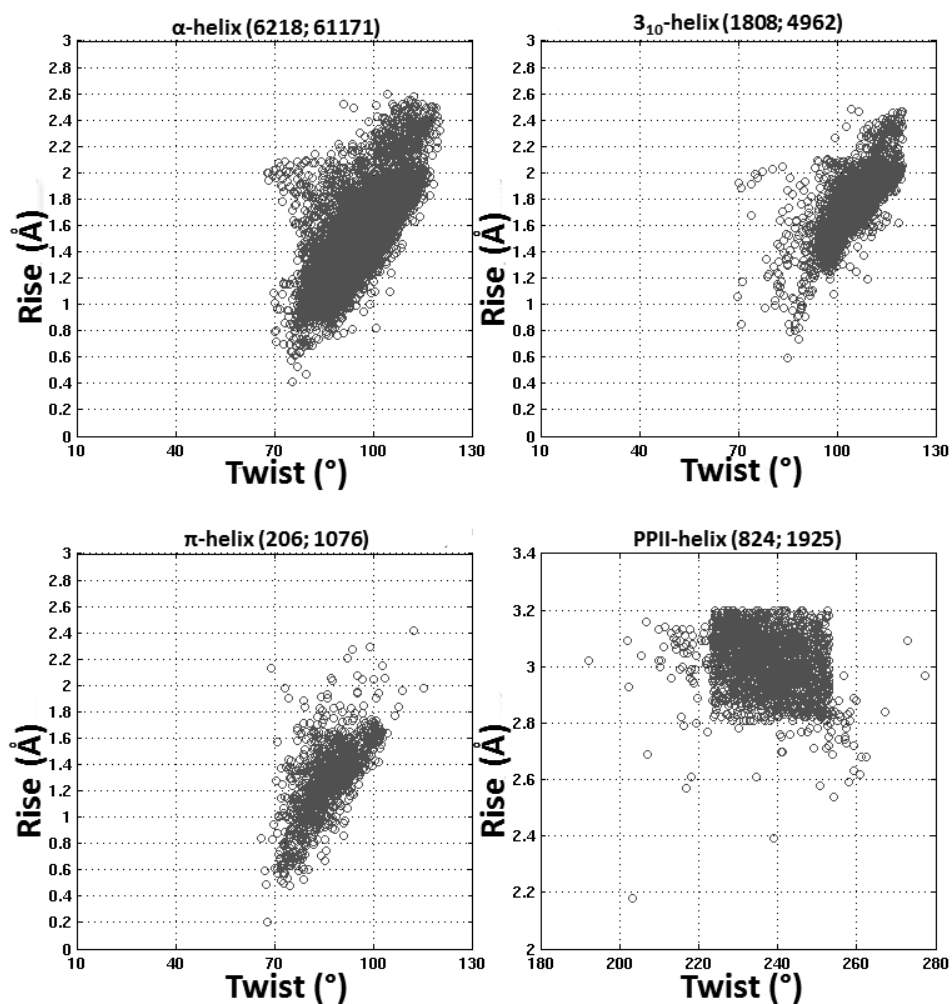
	ASSP	DSSP-PPII	PROSS	SEGNO	XTLSSTR
Twist ( $^{\circ}$ )	237.6 (9.2)	224.4 (44)	231.3 (35.9)	234.6 (27.8)	234.4 (33.7)
Rise ( $\text{\AA}$ )	3.0 (0.1)	2.8 (0.5)	2.9 (0.4)	2.9 (0.4)	2.9 (0.5)
Radius( $\text{\AA}$ )	1.3 (0.1)	1.4 (0.4)	1.4 (0.3)	1.4 (0.3)	1.3 (0.4)

**Table S4** Comparison of left handed helices as identified by ASSP and (Novotny & Kleywegt, 2005). The corresponding assignments by other algorithms, which have been discussed in the referenced paper, are not tabulated here.

Sl. No.	PDB ID: Chain ID	Protein Name	ASSP	(Novotny & Kleywegt, 2005)
1	1BD0:A	Alanine racemase	40-44 ( $\alpha$ )	40-44 ( $\alpha$ )
2	1AUT:L	Activated protein C	101-104 ( $3_{10}$ )	101-104 ( $3_{10}$ )
3*	1B9W:A	Merozoite surface protein 1 (P. cynomolgi)	52-55 ( $3_{10}$ )	52-55 ( $3_{10}$ )
4	1G2L:B	Coagulation factor X	258-261 ( $3_{10}$ )	258-261 ( $3_{10}$ )
5	1KLI:L	Coagulation factor VII	94-97 ( $3_{10}$ )	94-97 ( $3_{10}$ )
6	1N1I:A	Merozoite surface protein 1 (P. knowlesi)	57-60 ( $3_{10}$ )	57-60 ( $3_{10}$ )

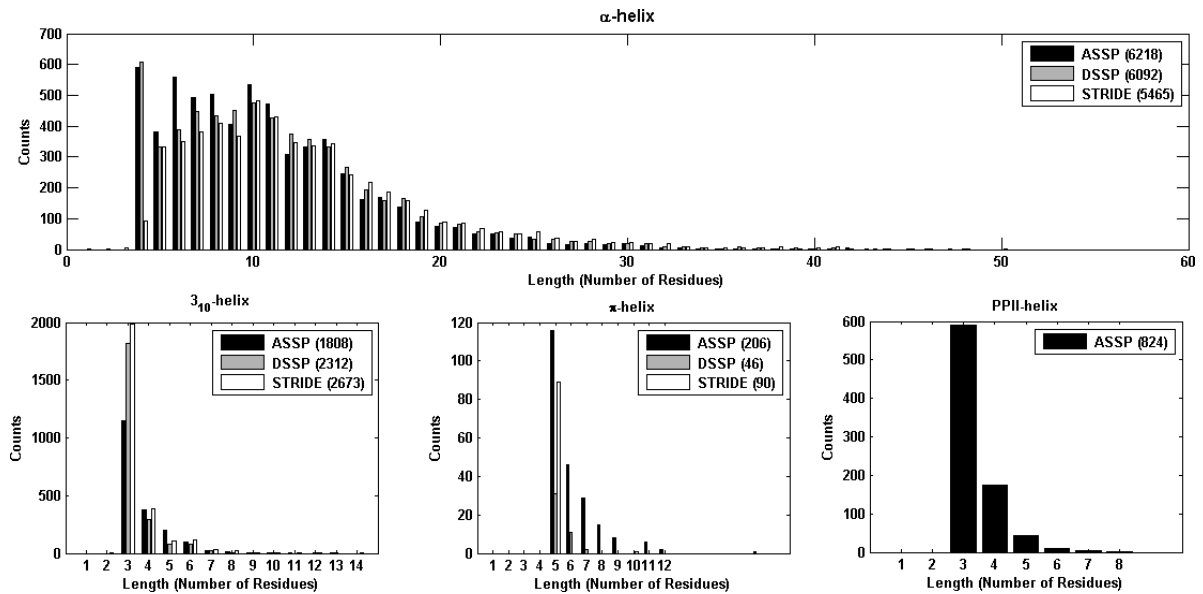
7	1OB1:C	Merozoite surface protein 1 (P. falciparum)	52-55 (3 <sub>10</sub> )	52-55 (3 <sub>10</sub> )
8	1RFN:_	Coagulation factor IX	91-94 (3 <sub>10</sub> )	91-94 (3 <sub>10</sub> )
9	1PB5:A	Lnr module from Notch	29-32 (3 <sub>10</sub> )	28-32 (3 <sub>10</sub> )
10	1KDG:A	Cellobiose dehydrogenase	532-535 (3 <sub>10</sub> )	532-535 (3 <sub>10</sub> )
11	1HZM:A	Protein phosphatase 6	60-63 ( $\alpha$ )	61-64 ( $\alpha$ )
12 <sup>s</sup>	1GTX:A	4-Aminobutyrate aminotransferase	70-73 ( $\alpha$ )	70-73 ( $\alpha$ )
13	1QJ5:A	7, 8-Diaminopelargonic acid synthase	50-53 ( $\alpha$ )	50-53 (3 <sub>10</sub> )
14	2GSA:A	Glutamate semialdehyde aminotransferase	65-67 (3 <sub>10</sub> )	65-68 (3 <sub>10</sub> )
15	2OAT:A	Ornithine aminotransferase	83-86 ( $\alpha$ )	83-86 (3 <sub>10</sub> )
16	1BNL:A	Endostatin (Homo sapiens)	135-138 ( $\alpha$ )	135-138 ( $\alpha$ )
17	1DY2:A	Endostatin (M. musculus)	207-210 ( $\alpha$ )	207-210 ( $\alpha$ )
18	1KOE:S	Endostatin (M. musculus)	266-269 ( $\alpha$ )	266-269 ( $\alpha$ )
19	1BQB:A	Aureolysin	223-226 ( $\alpha$ )	223-226 ( $\alpha$ )
20	1NPC:E	Neutral protease	227-230 ( $\alpha$ )	227-230 ( $\alpha$ )
21	8TLN:E	Thermolysin	266-269 ( $\alpha$ )	266-269 ( $\alpha$ )
22	1J9Q:A	Nitrate reductase (A. feacalis)	105-108 ( $\alpha$ )	105-108 (3 <sub>10</sub> )
23	1NIF:A	Nitrate reductase (A. cycloclates)	105-108 ( $\alpha$ )	105-108 (3 <sub>10</sub> )
24	1OE1:A	Nitrate reductase (A. xylosoxidans)	99-102 ( $\alpha$ )	99-102 (3 <sub>10</sub> )
25 <sup>†</sup>	1PTM:A	4-Hydroxythreonine-4-phosphate dehydrogenase	-	211-216 (3 <sub>10</sub> )
26	1AK0:E	P1 nuclease	131-134 ( $\alpha$ )	131-134 ( $\alpha$ )
27 <sup>†</sup>	1MZR:_	2, 5-Diketo-d-gluconate reductase	-	191-194 (3 <sub>10</sub> )
28	1H21:A	Split-soret cytochrome C	77-81 ( $\pi$ )	77-81 (3 <sub>10</sub> )
29	1HXX:A	Ompf porin	143-146 ( $\alpha$ )	143-146 (3 <sub>10</sub> )
30	1JV1:A	Glcnaclp uridyltransferase	182-185 (3 <sub>10</sub> )	182-185 (3 <sub>10</sub> )
31	1KWS:A	Beta-1,3-glucuronyltransferase 3	298-301 ( $\alpha$ )	298-301 ( $\alpha$ )

**Figure S1** Plots showing twist ( $^{\circ}$ ) vs. rise ( $\text{\AA}$ ) for the steps constituting the ASSP identified right handed  $\alpha$ ,  $3_{10}$ ,  $\pi$  and left handed PPII-helices. The corresponding total number of helices and steps is given in the parentheses above each plot. The mean values of parameters are listed in Table 1.

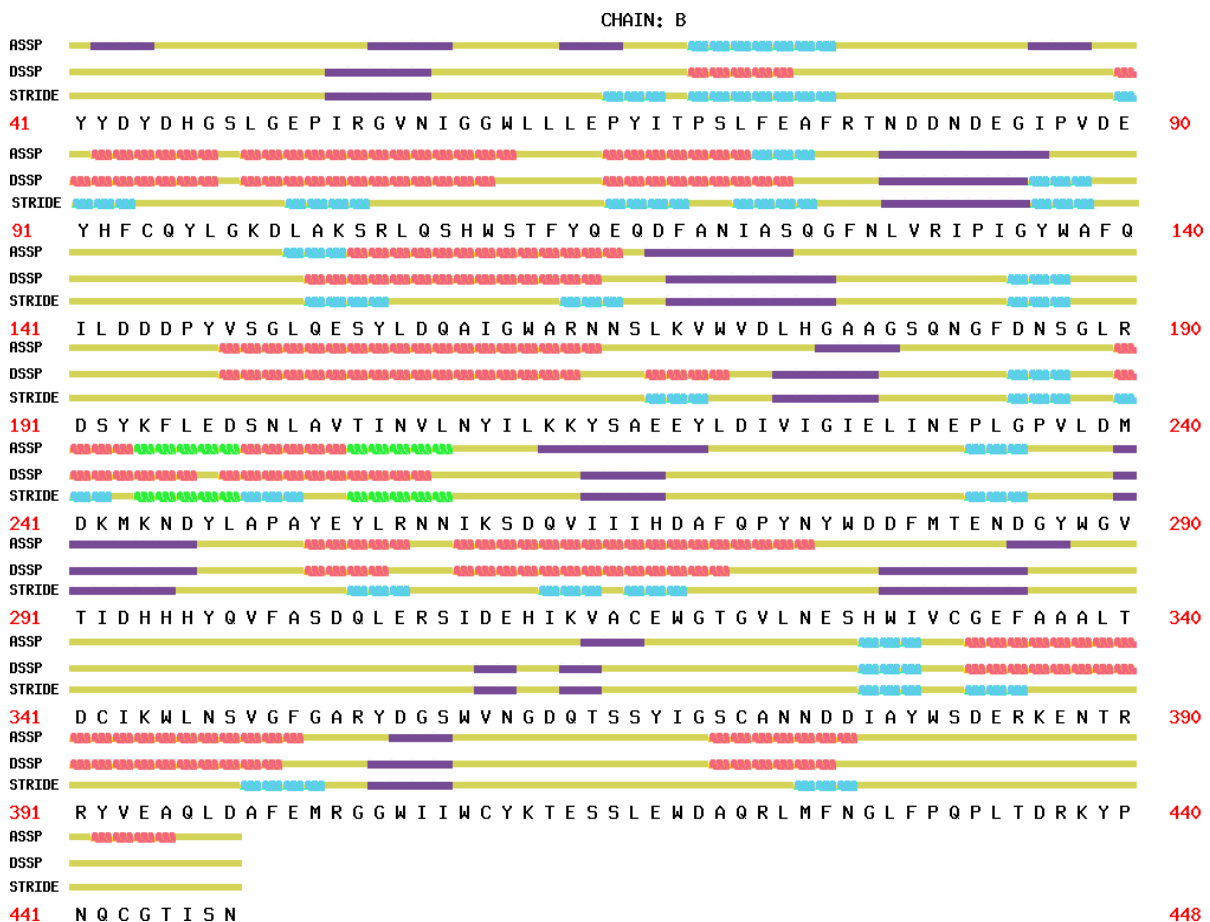




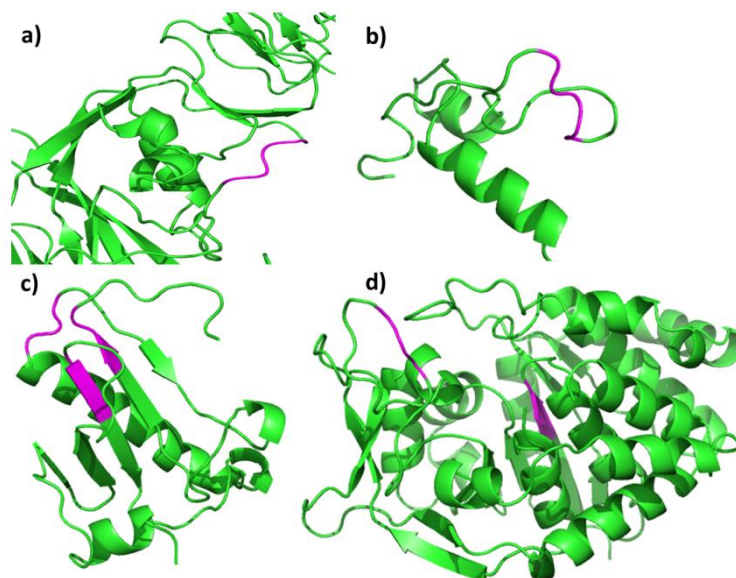
**Figure S2** Bar diagram showing length distribution for  $\alpha$ ,  $3_{10}$  and  $\pi$ -helices assigned by *ASSP*, *DSSP* and *STRIDE*. PPII-helices are not identified by *DSSP* and *STRIDE*.



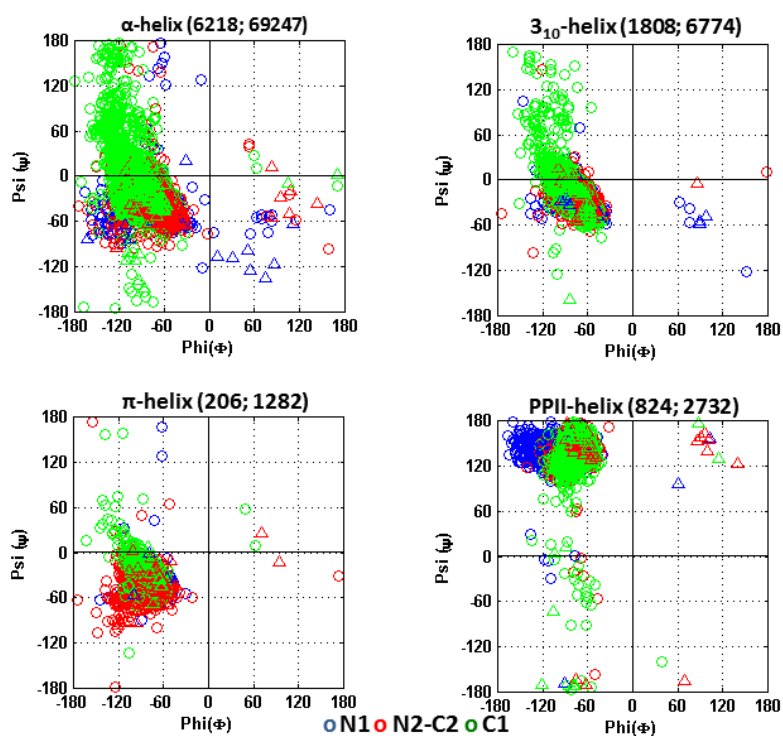
**Figure S3** Comparison of SSEs assigned by *ASSP*, *DSSP* and *STRIDE* for full length protein chain (PDB ID: 1H4P:B)



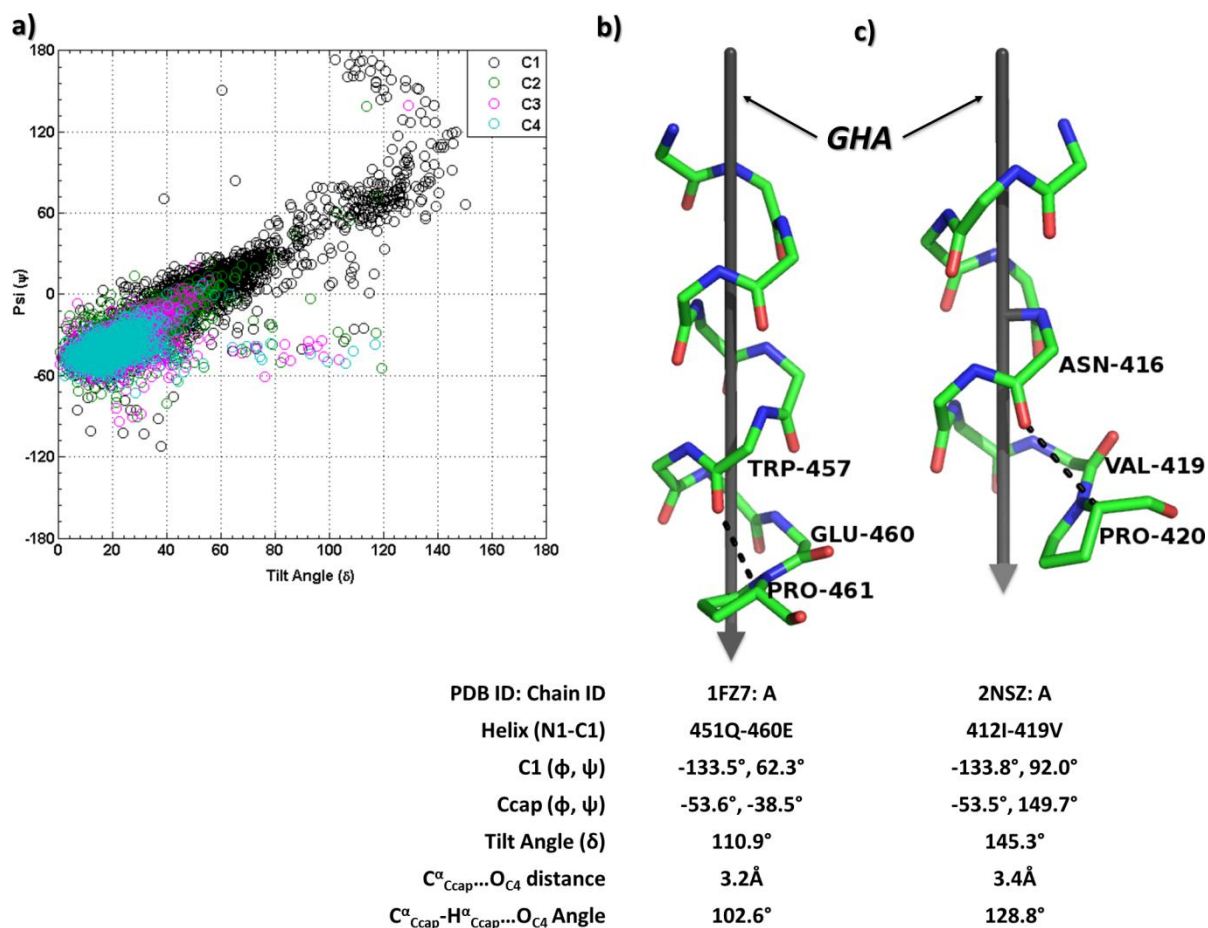
**Figure S4** Cartoon diagram showing the SSEs according to *PyMol* for a) 1JSD: A (84S-87N); b) 1EL6: A (25G-29V); c) 1KPF: A (50T-52F; 89G-91N; 116L-118G) and d) 1KIC: A (7L-9H; 174V-176L). Uniform segments identified as a part of cluster (iv) are shown in magenta.



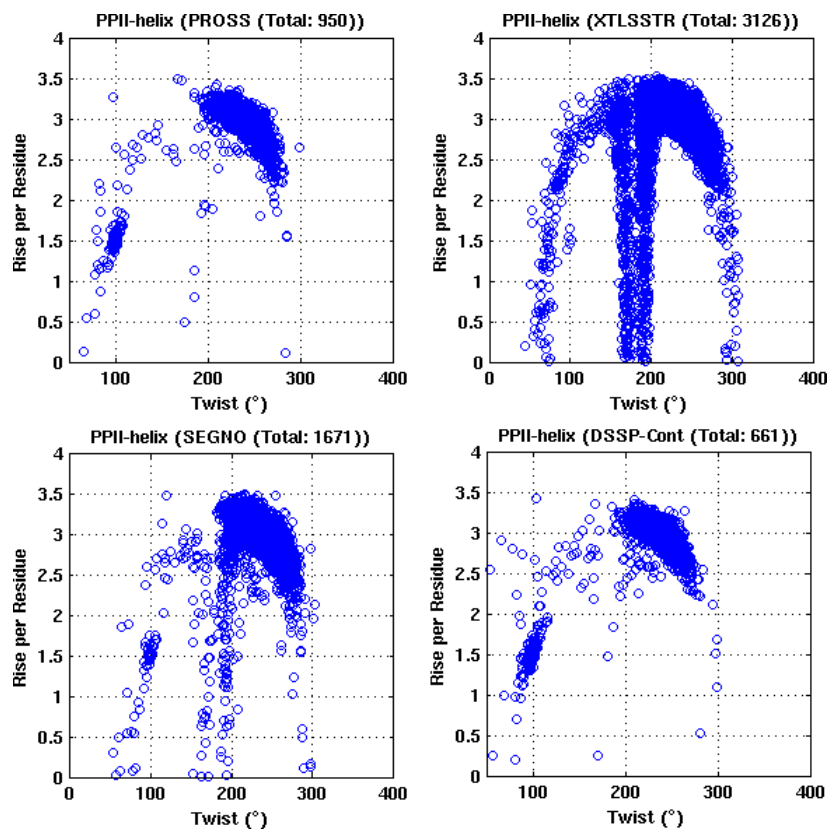
**Figure S5** Plots showing backbone torsion angles ( $\phi$ ,  $\psi$ ) for ASSP assigned right handed  $\alpha$ ,  $3_{10}$ ,  $\pi$  and left handed PPII helices. Only the helical residues (N1 to C1) are considered. GLY residues are shown as ' $\Delta$ '. The number of helices and residues constituting the helix are given within parentheses above each plot.



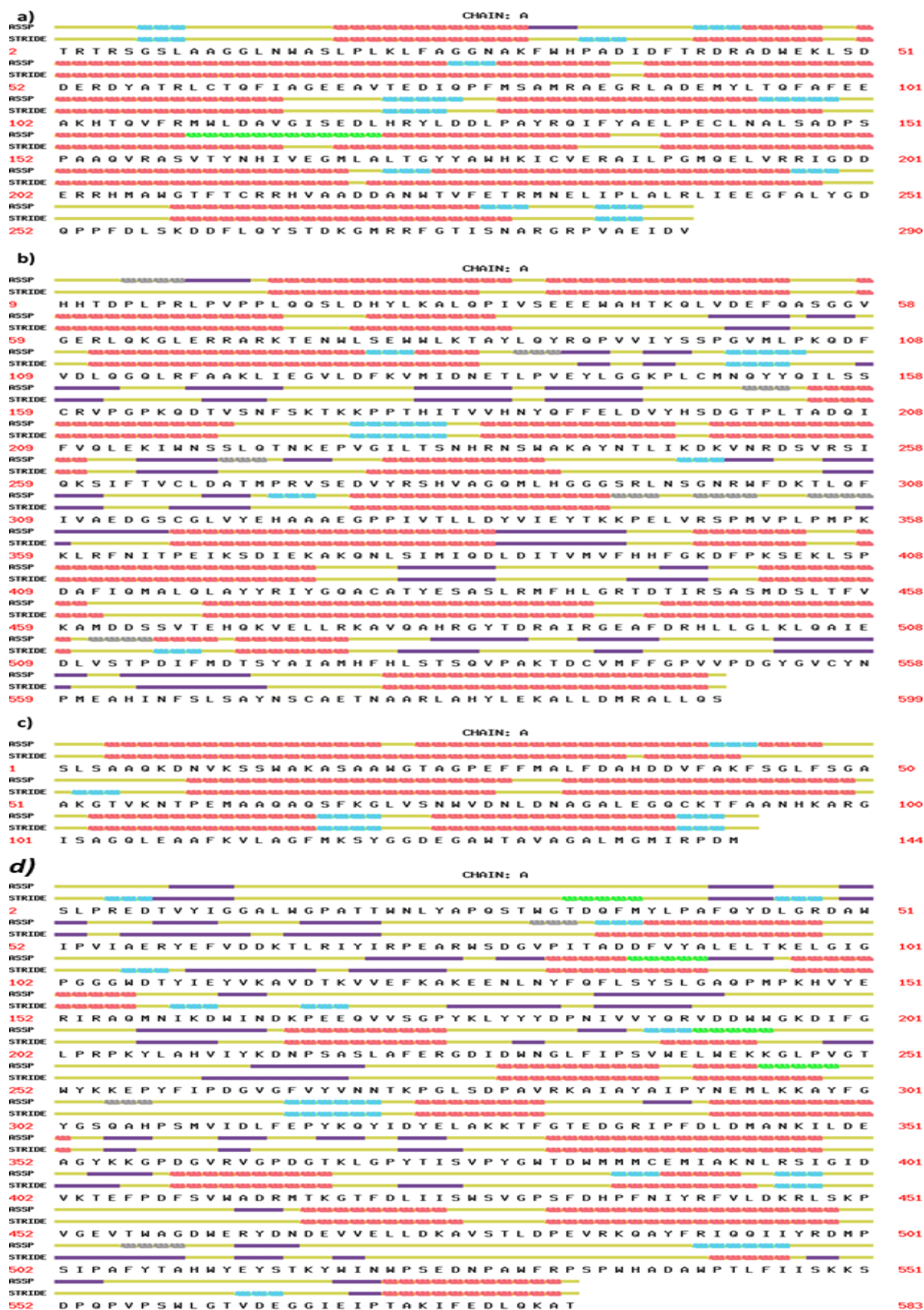
**Figure S6** a) Plot showing the correlation between tilt angle ( $\delta$ ) and backbone torsion angle Psi ( $\psi$ ) for residues at the C-terminal positions C4 to C1 of 4646  $\alpha$ -helices with length  $>6$  residues. b) and c) Representative examples of ASSP assigned  $\alpha$ -helices with  $C1_{\psi} > 40^{\circ}$  and  $C^{\alpha}$  of PRO at Ccap making a C-H...O hydrogen bond (shown in black dotted line) with backbone carbonyl oxygen of residue at C4 position.



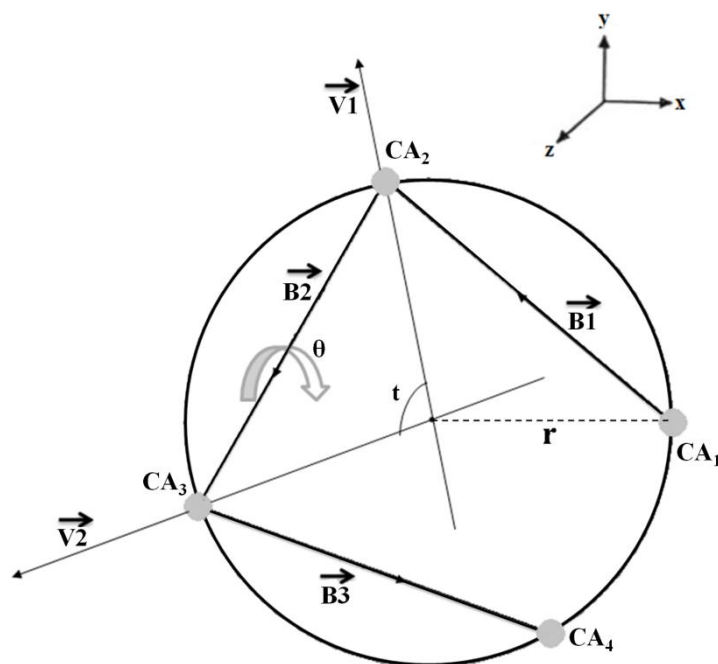
**Figure S7** Distribution of twist ( $^{\circ}$ ) vs rise ( $\text{\AA}$ ) plot for the steps constituting PPII-helices identified by different algorithms. Total number of helices is given in parentheses.



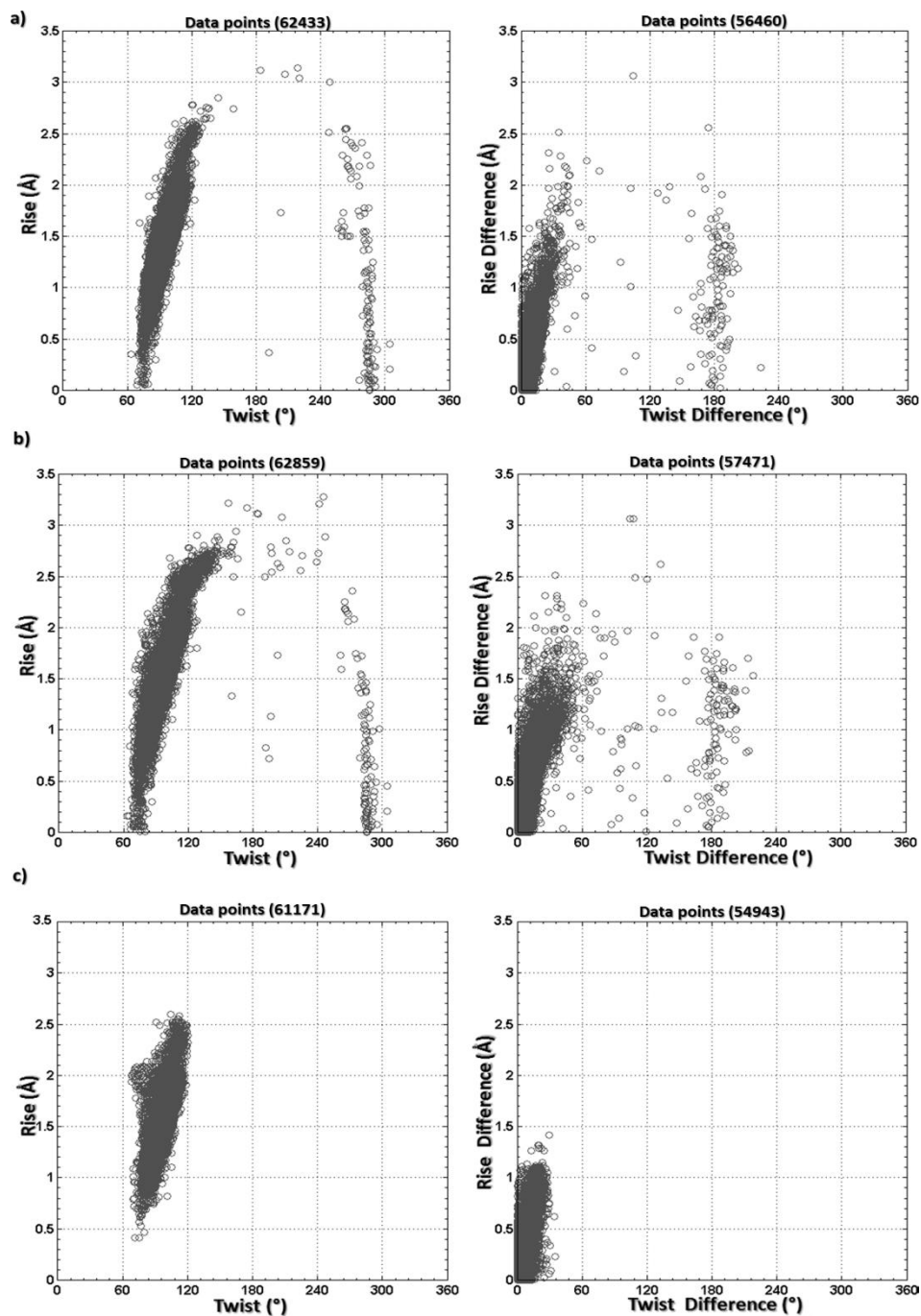
**Figure S8** Comparison of SSEs assigned by *ASSP* and *STRIDE* for full length protein chains: a) PDB ID: 3EE4:A b) 1NM8:A c) PDB ID:1B0B:A and d) 3I5O:A.



**Figure S9** Down the helix axis view of 4  $C^\alpha$  atoms in an ideal  $\alpha$ -helix. Various symbols used in the diagram are defined in the text.



**Figure S10** a) Distribution of twist ( $^{\circ}$ ) vs rise ( $\text{\AA}$ ) (column 1) and twist difference ( $^{\circ}$ ) vs rise difference ( $\text{\AA}$ ) (column 2) for the steps constituting  $\alpha$ -helices. a) *DSSP* assigned 5465  $\alpha$ -helices; b) *STRIDE* assigned 6092  $\alpha$ -helices and c) *ASSP* identified 6218  $\alpha$ -helices. The number of data points is given above each plot.



## References

- Cubellis, M. V., Cailliez, F. & Lovell, S. C. (2005). *BMC Bioinformatics* **6** Suppl 4, S8.  
Fodje, M. N. & Al-Karadaghi, S. (2002). *Protein Eng* **15**, 353-358.  
Frishman, D. & Argos, P. (1995). *Proteins* **23**, 566-579.

- Kabsch, W. & Sander, C. (1983). *Biopolymers* **22**, 2577-2637.
- King, S. M. & Johnson, W. C. (1999). *Proteins* **35**, 313-320.
- Konagurthu, A. S., Lesk, A. M. & Allison, L. (2012). *Bioinformatics* **28**, i97-105.
- Labesse, G., Colloc'h, N., Pothier, J. & Morion, J. P. (1997). *Comput Appl Biosci* **13**, 291-295.
- Levitt, M. & Greer, J. (1977). *J Mol Biol* **114**, 181-239.
- Majumdar, I., Krishna, S. S. & Grishin, N. V. (2005). *BMC Bioinformatics* **6**, 202.
- Mansiaux, Y., Joseph, A. P., Gelly, J. C. & de Brevern, A. G. (2011). *PLoS One* **6**, e18401.
- Martin, J., Letellier, G., Marin, A., Taly, J. F., de Brevern, A. G. & Gibrat, J. F. (2005). *BMC Struct Biol* **5**, 17.
- Novotny, M. & Kleywegt, G. J. (2005). *J Mol Biol* **347**, 231-241.
- Richards, F. M. & Kundrot, C. E. (1988). *Proteins* **3**, 71-84.
- Sklenar, H., Etchebest, C. & Lavery, R. (1989). *Proteins* **6**, 46-60.
- Srinivasan, R. & Rose, G. D. (1999). *Proceedings of the National Academy of Sciences of the United States of America* **96**, 14258-14263.
- Taylor, W. R. (2001). *J Mol Biol* **310**, 1135-1150.