

Supporting Information for

”New Insights into the Folding of a β -Sheet Miniprotein in a Reduced Space of Collective Hydrogen Bond Variables: Application to a Hydrodynamic Analysis of the Folding Flow”

by I. V. Kalgin, A. Caffisch, S. F. Chekmarev and M. Karplus

Collective variables

The algorithm to determine the collective variable that we used in the paper proceeds as follows:

1. Introduce a state vector $\mathbf{p}^m = (p_1^m, p_2^m, \dots, p_D^m)$ which is conjugate to the vector of the hydrogen bond distances $\mathbf{h}^m = (h_1^m, h_2^m, \dots, h_D^m)$, where D is the number of possible hydrogen bonds (dimension of conformation space), and m is the current number of the protein conformation among the total M conformations that are stored ($1 \leq m \leq M$). The bond distances are measured in angstroms.

2. Set the p_i^m to be equal to 1 if the hydrogen bond is formed and 0 otherwise. A bond was considered to be formed if $h_i^m \leq 3.6\text{\AA}$ and $\angle(\text{N} - \text{H} - \text{O})_i^m \geq 130^\circ$ [1]. To avoid the noise introduced by the spatial fluctuations characteristic of the C- and N-terminal residues, the hydrogen bonds that these residues formed with the other residues were not included (see Fig. 1 in the main text). This reduces the total number of possible hydrogen bonds to $D = 272$.

3. Apply the standard PCA algorithm [2] to the data set $\mathbf{p}^m = (p_1^m, p_2^m, \dots, p_D^m)$ ($1 \leq m \leq M$), i.e. calculate the means $\bar{p}_i = \text{E}(p_i)$, the covariance matrix $\mathbf{C}(p_i, p_j) = \text{E}[(p_i - \bar{p}_i)(p_j - \bar{p}_j)]$ and its eigenvalues λ and eigenvectors \mathbf{x} that satisfy the equation $\mathbf{C}\mathbf{x} = \lambda\mathbf{x}$ (the symbol E denotes the expected value).

4. Choose the desired dimension of the reduced space (K) and direct the collective variables g_1, g_2, \dots, g_K along the eigenvectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$ corresponding to the largest eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_K$. The protein conformation with bond distances h_1, h_2, \dots, h_D is determined by the values of the collective variables $g_j = \sum_{i=1}^D w_{ij} h_i$, where w_{ij} indicates the contribution of i bond into j variable.

5. Map the representative points of the protein in the original $\mathbf{h} = (h_1, h_2, \dots, h_D)$ space onto the reduced $\mathbf{g} = (g_1, g_2, \dots, g_K)$ space of collective variables.

As in the conventional PCA method, this algorithm transforms the original configuration space to a set of orthogonal collective variables (in our case, \mathbf{h} to \mathbf{g}) and is easy to implement. Since the collective variables are linear combinations of the original variables, they are measured in the same units as the latter, i.e. in angstroms. In what follows we refer to this algorithm to as the Hydrogen Bond PCA (HB PCA) method.

Principal Component Analysis (PCA), the atomic coordinate space

TABLE S1: Clusters of Protein Conformations

Cluster ^a	W_{clst} ^b	N_{str} ^c	Most populated structure ^d	W_{str} ^e	Cluster type ^f
1	24.29	467	-EEEEETEEEEETEEEE-	33.38	Native
			-EEEEETEEEEETEEEE--	31.77	
2	5.21	1613	-EEEEETEEEEEEEEEE--	6.37	Cs-or
			-EEEEETEEEEEEEEEE-	5.58	
3	11.73	1137	-EEE-SSS-EEEEETEEEE-	14.55	Native + Ns-or
			-EEEEETEEEEETEEEE-	9.77	
4	4.39	3881	--HHHHHHHHHHT-----	0.46	Helical
			--SS--HHHHTTT-----	0.32	
5	4.3	2991	-B-SSSSS-EEEEETESB-	1.98	Ch-curl
			-B-S-SSS-EEEEETESB-	1.84	
6	2.59	1017	-B-SSSSS--EEETTEE-B-	5.63	
			-B-S-SSS--EEETTESB-	3.86	
7	2.91	760	-B-SSSS--EESSSS-EE-	5.16	
			-B-SSSS-EEEEETEEEE-	4.95	
8	6.5	5502	-----SS--EEEEETEEEE-	0.57	Others
			----SSS--EEEEETEEEE-	0.31	
9	6.62	1992	--EEETEEEEETEEEE-	4.25	
			-EEEEETEEEEETEEEE-	4.07	
10	8.06	6247	-B--SSS---EEETTEE-B-	0.98	
			-B--SSS---EEETTESB-	0.6	
11	4.9	4132	-EEEEETEEEE---SS--B-	0.67	
			-EE-SSS-EE---SS---B-	0.55	
12	3.98	3721	--SSSS---EEEEETEEEE-	0.18	
			-EEEEETEEEEESSS-EE-	0.18	
13	3.96	3015	-B--SSS---EEETTEE-B-	1.14	
			----SSS--EEEEETEEEE-	0.76	
14	4.4	3911	-B-SSSSS--EEETTEE-B-	0.7	
			-B-SSSS---EETTEE-B-	0.43	
15	6.17	4086	-EEEEETEEEE-SSS-----	0.71	
			--EEETEEEEEEEEEE-	0.71	

^a Cluster number.

^b Cluster weight equal to the number of the representative points in the cluster relative to the total number of the points (in %).

^c The number of conformations that have different secondary structure strings.

^d The secondary structure strings of the most populated conformations.

^e Weight of the given conformation in the cluster (in %).

^f Corresponds to the clusters of Fig. S1.

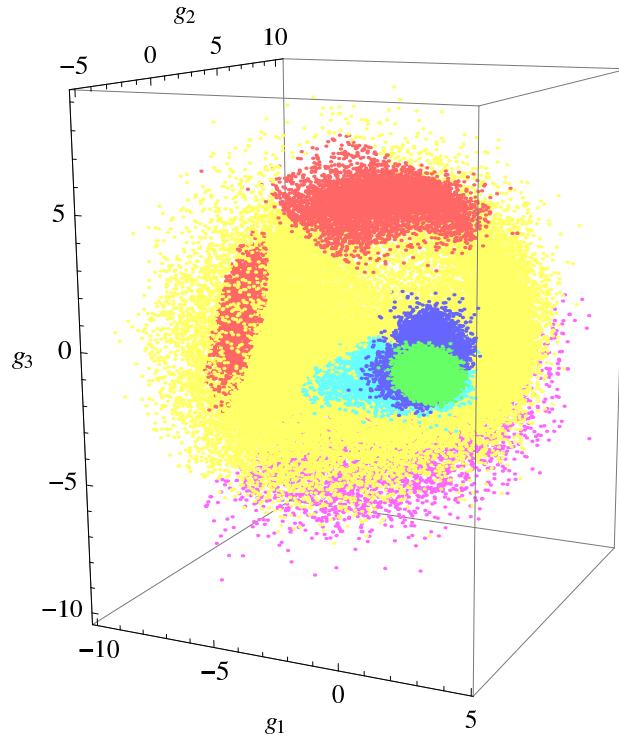


FIG. S1: Distribution of the representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space. The green, cyan, blue, magenta, red and yellow points correspond to Native, Cs-or, Native+Ns-or, Helical, Ch-curl and others conformations respectively. The representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space are clustered using MCLUST method [3].

Comments: The Ns-or cluster is not separated from the native cluster.

Principal Component Analysis (PCA), the hydrogen bond distance space

TABLE S2: Clusters of Protein Conformations (Notations are as in Table S1)

Cluster ^a	W_{clst} ^b	N_{str} ^c	Most populated structure ^d	W_{str} ^e	Cluster type ^f
1	14.72	1670	-EEEEETEEEEETEEEE--	36.34	Native
			-EEEEETEEEEETEEEE-	25.49	
			-EEEEETEEEEETEEEE-	19.11	
2	5.37	2572	--EEETEEEEETEEEE-	10.85	
3	5.81	2918	-EEEEETEEEEETEEEE--	36.1	
			-EEEEETEEEEETEE---	8.47	
4	2.75	3696	-EEEEETEEEEEEEEEE--	9.57	Cs-or
			-EEEEETEEEEEEEEEE-	9.42	
5	2.58	1272	-EEE-SSS-EEEETEEEE-	19.47	Ns-or
			-EEEESSSEEEETEEEE-	19.16	
6	11.73	2593	-EEEEETEEEEETEEEE-	35.03	Ns-or+Native
			-EEE-SSS-EEEETEEEE-	10.02	
7	6.61	35199	--HHHHHHHHHHT-----	0.52	Helical 1
			---HHHHHHHHHT-----	0.27	
8	4.36	33520	----SS---BTTTBSS---	0.21	Helical 2
			-----SHHHHHHTTTB---	0.16	
9	7.34	8920	-B-SSSS--EEETEE-B-	2.6	Ch-curl
			-B-SSSS-EEEETEEEE-	1.98	
10	9.91	34301	----SS--EEEETEEEE-	0.85	Others
			-B--SSS---EEETEE-B-	0.73	
11	4.92	8266	-EEEEETEEEEETEEEE-	4.7	
			---EETEEEEETEEEE-	3.05	
12	11.6	68766	--EEETEEEEETEEEE-	0.34	
			--EEETEEEEETTEES--	0.16	
13	4.7	13997	-EEEEETEEEEETEE---	2.42	
			-EEEEETEEEEETEEEE--	1.76	
14	2.56	18376	--EEETEEE--BTTB---	0.5	
			---EETEE---BTTB---	0.46	
15	5.04	35813	----SS-SS-EETEE---	0.18	
			-EEETTTEEETTTEEE--	0.16	

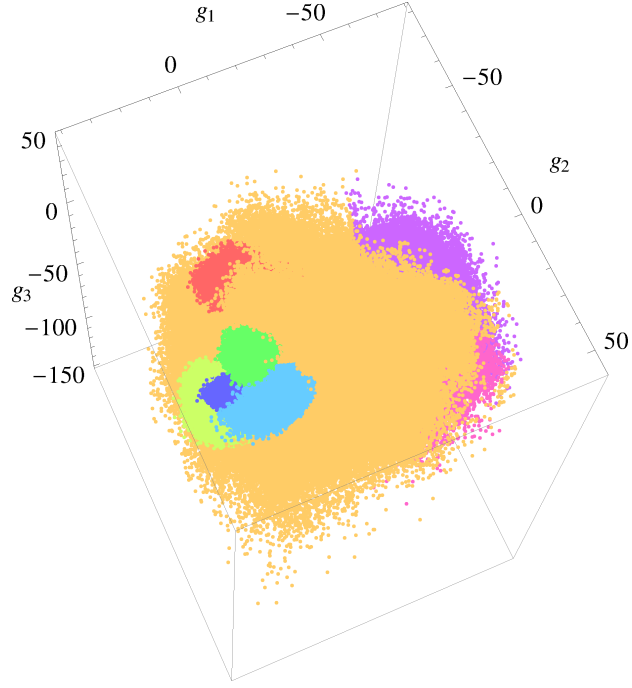


FIG. S2: Distribution of the representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space. The green-yellow, green, light blue, blue, purple, magenta, red and yellow points correspond to Native, Cs-or, Ns-or, Ns-or+Native, Helical 1, Helical 2, Ch-curl, and others conformations respectively. The representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space are clustered using MCLUST method [3].

Comments: The Cs-or and Ns-or clusters are not symmetric with respect to the Native cluster along g_2 . A cluster similar to the intermediate Ns-or+Native cluster, i.e. a Cs-or+Native cluster, is not observed.

Local Linear Embedding (LLE) [4], the atomic coordinate space

TABLE S3: Clusters of Protein Conformations (Notations are as in Table S1)

Cluster ^a	W_{clst} ^b	N_{str} ^c	Most populated structure ^d	W_{str} ^e	Cluster type ^f
1	9.44	854	-EEEEETEEEEETEEEE--	37.89	Native
			-EEEEETEEEEETEEEE-	15.18	
2	17.91	615	-EEEEETEEEEETEEEE-	39.27	
			-EEEEETEEEEETEEEE--	24.37	
3	3.08	1205	-EEEEETEEEEEEEEEE--	8.71	Cs-or
			-EEEEETEEEEEEEEEE-	8.54	
4	2.37	580	-EEEESSSEEEEEETEEEE-	19.25	Ns-or
			-EEE-SSS-EEEEETEEEE-	17.1	
5	7.06	3599	-EEE-SSS-EEEEETEEEE-	9.85	Ns-or + Others
			--EEETEEEEETEEEE-	1.56	
6	3.46	1013	-B-SSSS--EEEEETEEEE-	4.48	Ch-curl
			-B-SSSS--EESSSS-EE-	4.45	
7	2.96	1210	-B-SSSS--EEETTEE-B-	6.02	
			-B-S-SSS--EEETTEESB-	3.45	
8	3.02	1872	-EEE-SSS-EEEEETEEEE-	2.62	
			-B--SSS---EEETTEE-B-	2.35	
9	6.59	5503	-EEE-SSS-EEEEETEEEE-	0.52	Others
			-EEETEEEEETEEEE--	0.3	
10	7.49	6027	-EEE-SSS-EEEEETEEEE-	1.72	
			-EEEEETEEEEETTEE---	0.61	
11	6.93	2358	-EEEEETEEEEETEEEE--	6.48	
			-EEEEETEEEEETEEEE-	5.98	
12	4.72	1596	-EEEEETEEEEETEEEE-	12.94	
			--EEETEEEEETEEEE-	9.14	
13	3.59	3023	-EEEEETEEEEETTEE---	0.97	
			-EEEEETEEEEETEEEE--	0.95	
14	0.98	879	--EEEEETEEEEETEEEE-	0.92	
			--EEEE-SSEEEETEEEE-	0.72	
15	8.53	6357	-EEEEETEEEEETEEEE-	0.83	
			-EEE-SSS-EEEEETEEEE-	0.42	
16	4.72	3728	-EE-SSS-EE---SS---B-	1.17	
			--EEETEEEEETEEEE-	1.02	
17	7.16	6262	--SSSS---EEEEETEEEE-	0.36	
			--SSSS--EEEEETEEEE-	0.31	

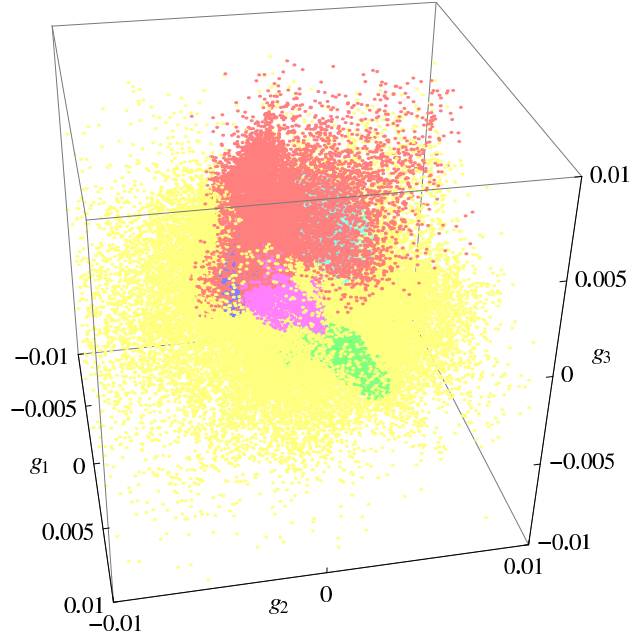


FIG. S3: Distribution of the representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space. The green, cyan, blue, magenta, red and yellow points correspond to Native, Cs-or, Ns-or, Ns-or+Native, Ch-curl and others conformations respectively. The representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space are clustered using MCLUST method [3].

Comments: All clusters are more or less separated, but the picture is vague.

Full Correlation Analysis (FCA) [5, 6], the atomic coordinate space

TABLE S4: Clusters of Protein Conformations (Notations are as in Table S1)

Cluster ^a	W_{clst}^b	N_{str}^c	Most populated structure ^d	W_{str}^e	Cluster type ^f
1	31.02	1584	-EEEEETEEEEETEEEE-	25.91	Native
			-EEEEETEEEEETEEEE--	23.42	
			-EEEEETEEEEETEEEE-	11.88	
			-EEEEETEEEEETTEE---	8.48	
2	3.76	1525	-EEEEETEEEEETEEEE-	16.45	
			-EEEEETEEEEETTEE---	5.14	
3	5.66	1420	-EEEEETEEEEETEEEE-	12.16	
			--EEETEEEEETEEEE-	5.24	
4	8.14	2713	-EEEEETEEEEETEEEE-	12.14	Cs-or
			-EEEEETEEEEETEEEE-	11.24	
5	2.45	605	-EEEEETEEEEETEEEE-	1.22	Ns-or
			-EEEEETEEEEETEEEE---	1.22	
6	2.21	1675	--HHHHHHHHHT-----	0.54	Helical
			--TTTHHHHHHS-----	0.31	
7	2.23	2009	-B--SSS---EEETEE-B-	2.67	Ch-curl
			-B--SSS---EEETTESB-	1.39	
8	3.82	3011	-B-SSSS---EEETEE-B-	5.47	
			-B-SSSS---EEETEE-B-	3.71	
9	2.85	1068	-B-SSSS-EEEEETTESB-	2.47	
			-B-S-SSS-EEEEETTESB-	2.24	
10	3.52	2357	-B-SSSS--EESSSS-EE-	5.3	
			-B-SSSS--EEEEETTEE-	5.02	
11	2.89	765	-EETTTEETTTEEEEE--	0.7	Others
			-EEEEETEEEEETEEEE--	0.56	
12	3.41	2824	--SSSS--EEEEETEEEE-	0.92	
			-EE-SSS-EE---SS---B-	0.8	
13	5.	3640	----SS--EEEEETEEEE-	0.34	
			-B---SSS-----SSS--B-	0.25	
14	3.53	3090	-EEEEETEEEEETTEE---	0.69	
			--EEETEEEEETEEEE-	0.43	
15	6.52	5303	-B-SSSS---EEETEE-B-	0.53	
			-EEEEETEEEE---SS--B-	0.51	
16	6.43	5444	-B--SSS---EEETEE-B-	0.24	
			----SSS-----SSS-----	0.24	
17	6.56	5952	-B--SSS---EEETEE-B-	0.24	
			----SSS-----SSS-----	0.24	

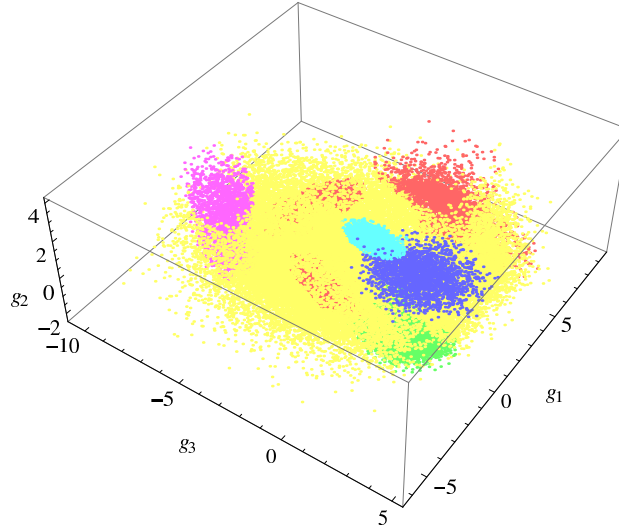


FIG. S4: Distribution of the representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space. The green, cyan, blue, magenta, red and yellow points correspond to Native, Cs-or, Ns-or, Helical, Ch-curl and others conformations respectively. The representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space are clustered using MCLUST method [3].

Comments: Separation into clusters is reasonable, but the Ns-or cluster (blue) is positioned between the Cs-or (cyan) and Native (green) clusters.

Manifold Sculpting (MS) [7], the hydrogen bond distance space

TABLE S5: Clusters of Protein Conformations (Notations are as in Table S1)

Cluster ^a	W_{clst} ^b	N_{str} ^c	Most populated structure ^d	W_{str} ^e	Cluster type ^f
1	6.07	425	-EEEEETEEEEETEEEE-	48.53	Native
			--EETEEEEETEEEE-	6.31	
			-EEEEETEEEEETEEEE-	35.25	
			-EEEEETEEEEETEEEE-	15.77	
2	12.8	820	-EEEEETEEEEETEEEE-	35.28	
			-EEEEETEEEEETEEEE-	21.59	
3	10.08	572	-EEEEETEEEEETEEEE-	53.52	
			-EEEEETEEEEETEEEE-	5.39	
4	3.4	239	-EEEEETEEEEETEEEE-	12.28	Cs-or
			-EEEEETEEEEETEEEE-	10.3	
5	2.58	577	-EEE-SSS-EEEEETEEEE-	43.28	Ns-or
			-EEEESSSEEEETEEEE-	5.48	
			-EEEESSSEEEETEEEE-	19.04	
6	3.25	320	-EEE-SSS-EEEEETEEEE-	18.59	
			-EEE-SSS-EEEEETEEEE-	18.59	
7	2.66	262	-EEEESSSEEEETEEEE-	0.53	Helical
			-EEE-SSS-EEEEETEEEE-	0.34	
8	6.47	5050	--HHHHHHHHHT-----	3.35	Ch-curl
			--HHHHHHHHHS-----	2.7	
9	0.93	748	-B-SSSS-EEETTTEEE-	8.49	
			-B-SSSS--EEETTTEEE-	7.9	
10	1.03	276	-B-SSSS-EEETTTEESB-	8.4	
			-B-S-SSS-EEETTTEESB-	8.4	
11	1.88	429	-B-SSSS--EEETTTEEE-	13.89	
			-B-SSSS--EESSSS-EE-	13.68	
12	0.94	170	-B-SSSS--EE-SSS--EE-	7.44	
			-B-SSSS-EE-SSS--EE-	4.48	
13	2.84	605	-B-SSSS--EETTTEE-B-	1.46	Others
			-B-SSSS--EETTTEE-B-	1.3	
14	1.92	1377	--EETEEEEETEEEE-	4.36	
			-EEEEETEEEE-SSS-----	3.11	
15	5.56	2160	-EEEEETEEEEETEEEE-	1.98	
			--EETEEEEETEEEE-	1.38	
16	2.32	1432	-EE--SSSEEEETEEEE-	11.33	
			----SSS--EETTTEEE--	7.73	
17	4.22	1339	-EEEEETEEEEETEEEE-	0.15	
			-EEEEETEEEEETTEE---	0.11	
18	9.27	8548	----SS--BTTTBSS---	1.01	
			----SSS----SSS-----	0.71	
19	6.04	4144	----SS--EEETTTEEE-	0.3	
			----SSS--EEETTTEEE-	0.29	
20	6.99	5782	----SS--S-EETTTEE---	0.57	
			----SSS--EETTTEEE-	0.51	
21	5.3	4346	-EEEEETEEEE--SS--B-	0.58	
			---EETTTEE--SSS-----	0.37	
22	3.48	3127	-B---SSS----SSS--B-		
			-B---SSS----BTTTB-B-		

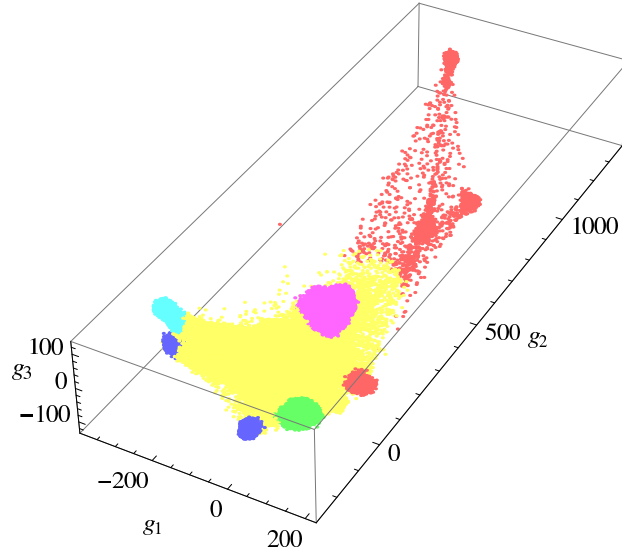


FIG. S5: Distribution of the representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space. The green, cyan, blue, magenta, red and yellow points correspond to Native, Cs-or, Ns-or, Helical, Ch-curl and others conformations respectively. The representative points of the protein in the $\mathbf{g} = (g_1, g_2, g_3)$ space are clustered using MCLUST method [3].

Comments: Clustering is good but the Ns-or and Ch-curl subclusters are very scattered.

TABLE S6: Comparison of the Weights (in %) of Consolidated Clusters

Cluster	KGA ^a	pfoldf ^b	REMD ^c	CTMD ^d	PCA ^e	PCA ^f	LLE	FCA	MS
Native	36.4	35.0	37.8	37.1	24.3	37.63	27.4	36.7-48.6 ^h	33.1
Cs-or	3.6	2.6	5.3	5.3	5.21	2.75	3.1	2.5	2.6
Ns-or	7.4	6.2	7.3	6.3	11.7	2.58	2.4 - 9.5 ^g	2.2	5.9
Helical	11.6	11.2	-	-	4.4	10.97	-	2.2	6.5
Ch-curl	6.0	4.9	3.9	1.8	9.8	7.34	9.4	13.1	7.6

^a Kinetic grouping analysis (KGA), Refs. [8, 9].

^b p_{fold} analysis based on an equilibrium kinetic network (pfoldf), Ref. [9].

^c Replica exchange molecular dynamics (REMD), Ref. [10].

^d Constant temperature molecular dynamics (CTMD), Ref. [10].

^e PCA in coordinate space.

^f PCA in hydrogen bonds space.

^g Ns-or and Ns-or+Others.

^h Clusters "1+3" and "1+2+3+4".

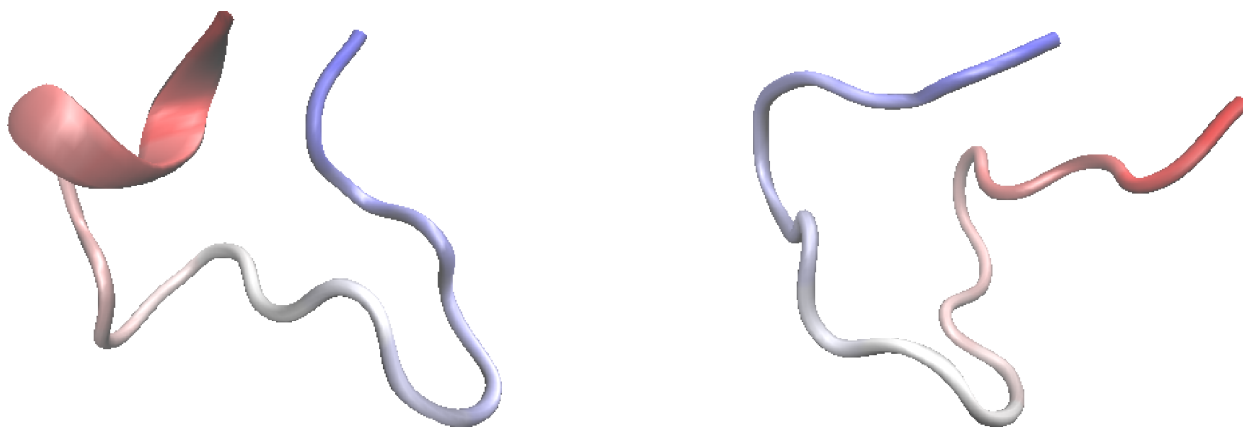


FIG. S6: Examples of unstructured conformers

All-atom RMSD vs the distance in the \mathbf{g} space

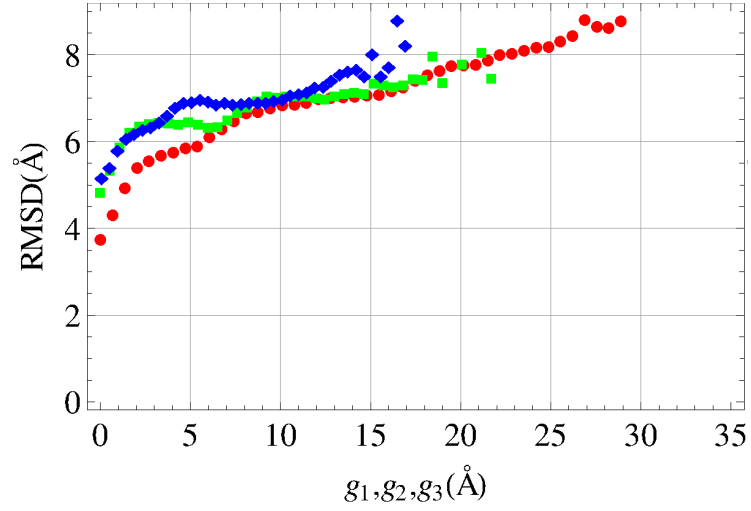


FIG. S7: Average all-atom RMSD between two conformations which correspond to different points (1 and 2) in the \mathbf{g} space as a function of the distance $g_i = |g_i^2 - g_i^1|$, where $i = 1, 2, 3$, and the upper indices denote the points. The averaging is taken over all points in the $g_i = g_i^1$ and $g_i = g_i^2$ planes. The red, green and blue curves are for g_1 , g_2 and g_3 , respectively.

Projections of the eigenvectors onto the hydrogen bond space

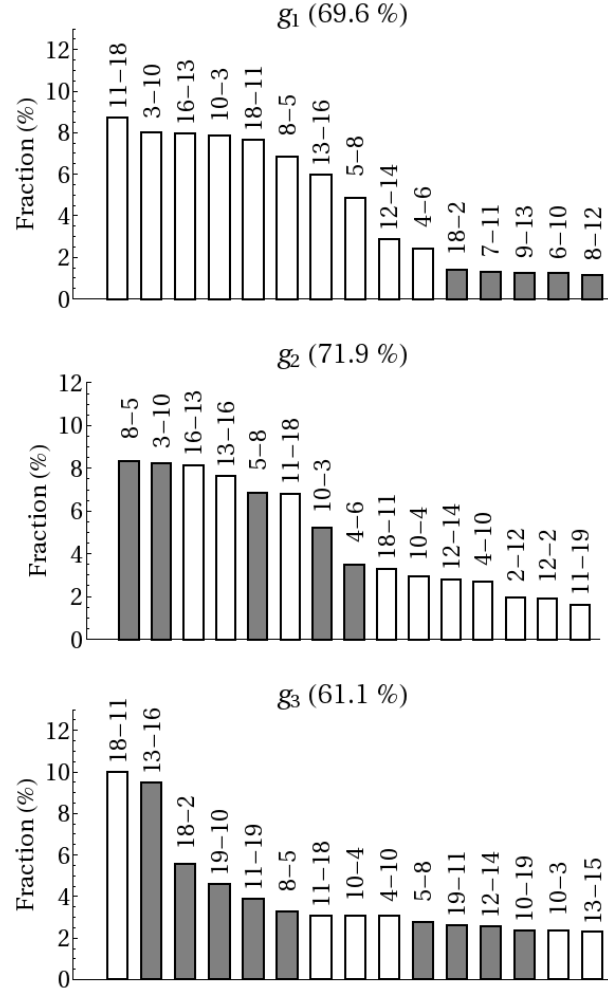


FIG. S8: Fractions of the hydrogen bonds which have a major contribution to the collective variables g_1 , g_2 and g_3 . The figures at the top of each bar denote the bond; the first figure is the number of the residue with the oxygen atom and the second figure is that with the nitrogen atom. The empty and solid bars are for the bond contributions to the negative and positive directions of the collective variable, respectively. The numbers in percentage at the top of each panel are the total contribution of the given bonds to the collective variable.

Contributions of hydrogen bonds to collective variables

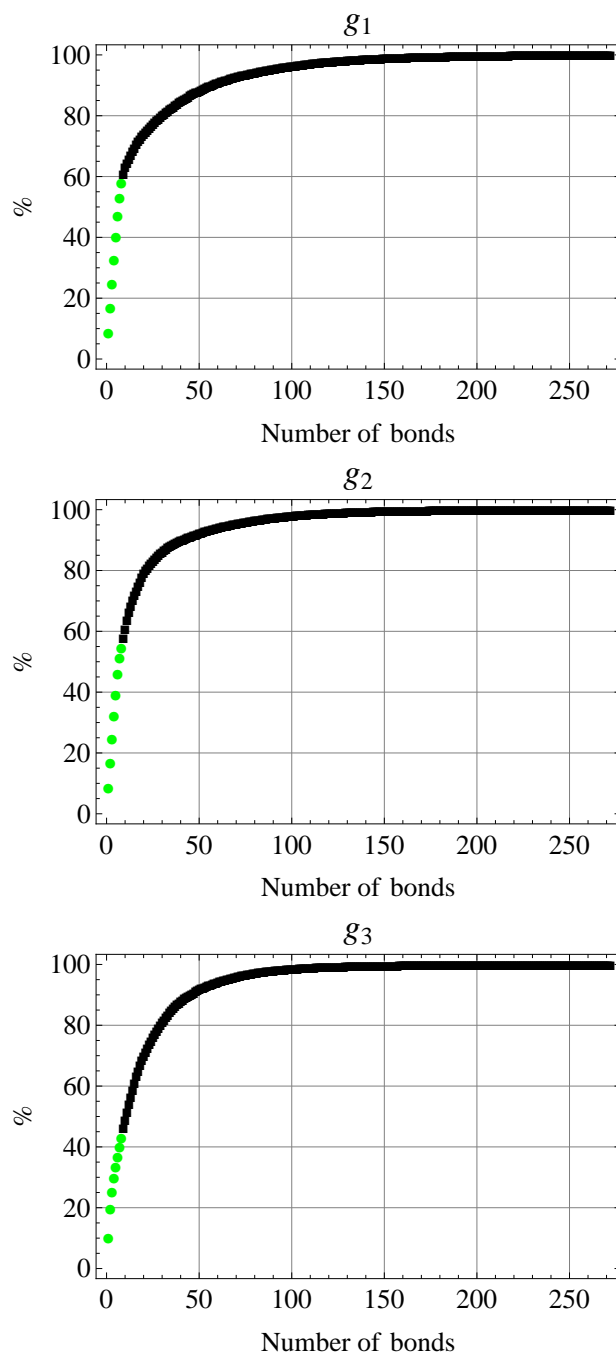


FIG. S9: Cumulative contributions of the hydrogen bonds to the collective variables g_1 , g_2 and g_3 . The green-colored points corresponds to the first eight bonds depicted in Fig. 6 of the main text.

Free energy profiles: Clustering over g_1 and the $\mathbf{g} = (g_1, g_2, g_3)$ space

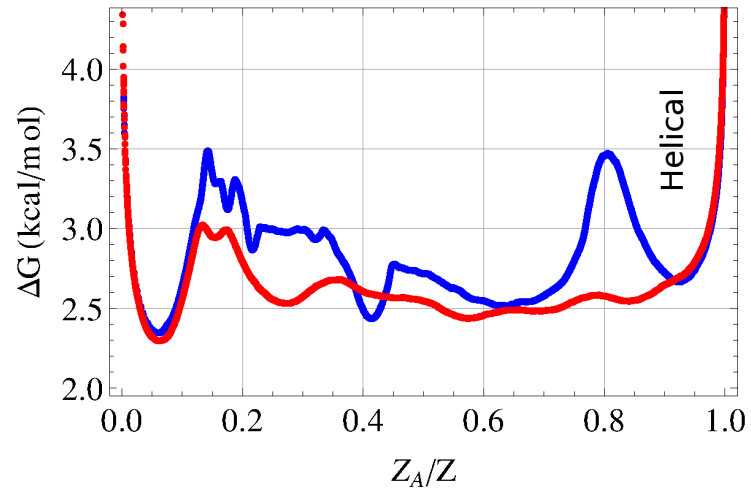


FIG. S10: Free energy profiles. The red curve corresponds to the clustering over the g_1 collective variable (Fig. 7 of the main text), and the blue curve to the clustering in the $\mathbf{g} = (g_1, g_2, g_3)$ space.

Comparison to the Zheng et al. study [12]

In their study of folding of beta3s miniprotein, Zheng et al. [12] used the LSDMap technique by Rohrdanz et al. [13] to reduce the conformation space of the protein. They have found that the first principal coordinate plays the role of the reaction coordinate for the folding process and the others, which correspond to smaller eigenvalues, discriminate between the clusters of representative conformations of the protein (basins on the FES). The LSDMap is a nonlinear dimensionality reduction technique, which treats the whole ensemble of configurations sampled from an MD simulation as a noisy data set that resides on a low-dimensional underlying manifold. Calculating the radius in configuration space around every point within which the manifold of the representative points can be approximated by a hyperplane tangent to the manifold, a Gaussian kernel determining the affinity of the points is introduced. This kernel is considered as a transition probability matrix for a Markov process to explore the configuration space of the system (equivalently, the process is described by the Fokker-Plank equation). The eigenfunctions of the kernel matrix present collective coordinates corresponding to the (diffusive) motion of the system on different time scales. Specifically, the zeroth eigenfunction corresponds to the Boltzmann distribution, the first eigenfunction to the collective motion with the slowest time scale, the second eigenfunction to the second slowest collective motion, etc., so that the variables of Ref. [12] allow discrimination of the protein dynamics on different time scales. However, a consequence of this useful feature of the variables of Ref. [12] is that the spatial distributions constructed with use of these variables (FESs, spatial kinetic networks, etc.) will be "biased", because different variables correspond to different time scales. In contrast, the variables we use allow the construction of "unbiased" spatial distributions. Different time scales can be distinguished by comparison of the cross-sections of the tubes that connect the clusters, i.e. the larger the cross-section, the faster the motion.

Rates of transitions between the clusters of conformations vs the distances in the atomic coordinate space

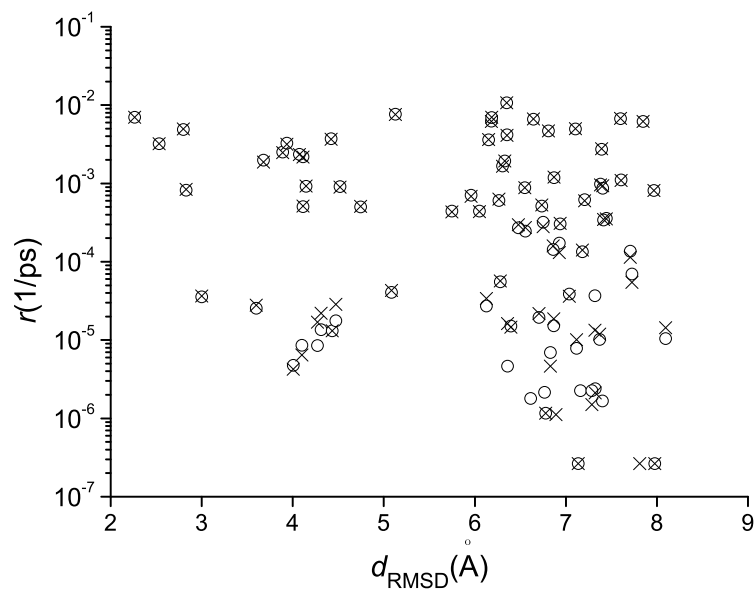


FIG. S11: Rates of transitions between the clusters of conformations vs the distances between the centers of the clusters in the atomic coordinate space. Crosses and circles are for the transitions from smaller and larger populated clusters, respectively.

TABLE S7: Numbers of transitions between the clusters and residence times.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
1	164687	9333	16	12161	138	0	26291	1	1	0	0	128	1923	55	8	1	0
2	9473	21614	13	1108	10	0	1869	0	0	0	0	371	4487	77	7	11	0
3	18	8	20627	239	0	0	3	0	0	0	0	4993	20	0	165	128	0
4	12080	1172	236	14019	12	0	2141	0	0	0	0	715	933	7	61	2	0
5	135	5	0	11	26124	1307	467	0	0	0	0	0	501	1076	89	0	0
6	0	0	0	0	1309	21040	6	0	0	0	0	0	13	2128	52	0	0
7	26293	1930	4	2111	458	6	16273	0	0	0	0	34	2195	177	8	0	0
8	0	0	0	0	0	0	0	57715	8238	0	0	0	0	2	1085	29	9303
9	1	0	0	0	0	0	0	8299	38533	0	0	0	0	9	857	1	3678
10	0	0	0	0	0	0	0	0	0	31748	603	0	0	575	201	283	0
11	0	0	0	0	0	0	0	0	0	600	41472	1	0	31	1029	585	454
12	103	356	5003	722	1	0	33	3	0	0	1	33735	735	43	863	3944	221
13	1886	4506	19	944	508	8	2221	0	0	0	0	684	17723	3231	292	168	0
14	57	102	0	2	1075	2125	162	3	7	574	33	55	3194	70172	5578	671	22
15	9	6	167	58	80	62	22	1080	875	198	1017	873	301	5550	63472	4131	9202
16	1	9	115	3	0	0	1	26	0	290	605	3922	165	683	4081	20111	3571
17	0	0	1	0	0	0	0	9245	3724	0	441	252	0	16	9255	3517	70631

The headers are the cluster numbers according to Table I of the main text. The off-diagonal elements are the numbers of transitions between the clusters, and the diagonal elements are the numbers of transitions within the clusters (the residence times) registered in the course of the $20\mu s$ trajectory.

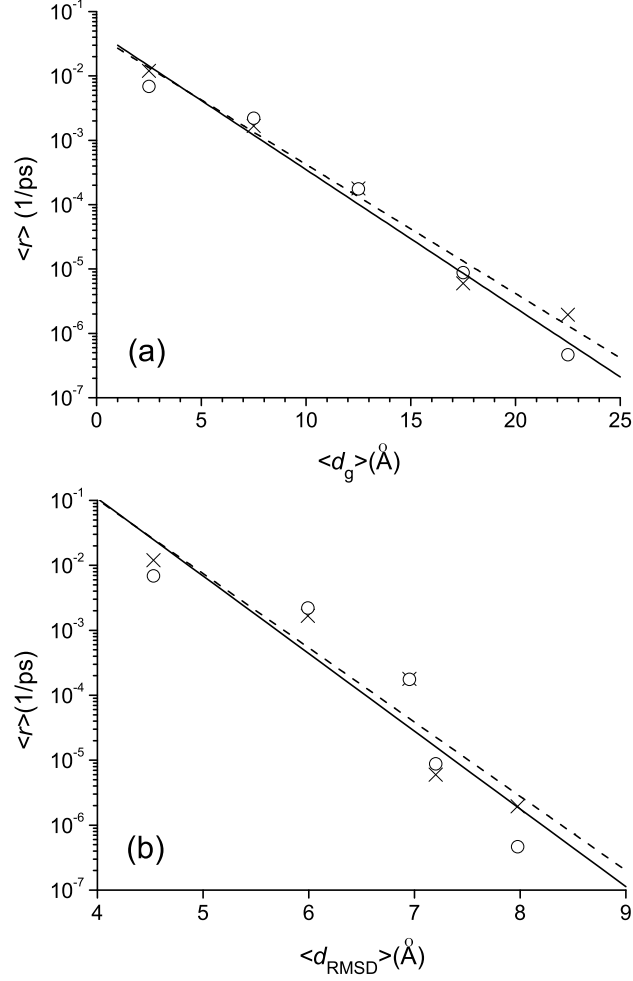


FIG. S12: Average rates of transitions between the clusters of conformations vs the average distances between the centers of the clusters in (a) \mathbf{g} space and (b) atomic coordinate space. Crosses and circles are for the transitions from smaller and larger populated clusters, respectively. In panel (a), the dashed line corresponds to the best fit for the crosses [$\langle r \rangle \sim \exp(-0.46\langle d_g \rangle)$], and the solid line to that for the circles [$\langle r \rangle \sim \exp(-0.49\langle d_g \rangle)$]. Similar in panel (b), where the dashed line is the best fit for the crosses [$\langle r \rangle \sim \exp(-2.63\langle d_{\text{RMSD}} \rangle)$], and the solid line is the best fit for the circles [$\langle r \rangle \sim \exp(-2.75\langle d_{\text{RMSD}} \rangle)$].

-
- [1] Seeber, M.; Cecchini, M.; Rao, F.; Settanni, G.; Caffisch, A. *Bioinformatics* **2007**, *23*, 2625-2627.
- [2] Jolliffe, I. T. *Principal Component Analysis*; 2nd ed.; Springer: New York, 2002.
- [3] Fraley, C.; Raftery, A. E. *J. Am. Stat. Assoc.*, **2002**, *97*, 611-631.
- [4] Roweis, S. T.; Saul, L.,K. *Science*, **2000**, *290*, 2323-2326.
- [5] Lange, O.; Grubmuller, H. *Proteins* **2006** *62*, 1053-1061.
- [6] Lange, O.; Grubmuller, H. *Proteins* **2008** *70*, 1294-1312.
- [7] Gashler, M.; Ventura, D.; Martinez, T. In *Advances in Neural Information Processing Systems 20*; Platt, J. C., Koller, D., Singer, Y., Roweis, S., Eds.; MIT Press: Cambridge, MA, 2008, pp. 513-520.
- [8] Muff, S.; Caffisch, A. *Proteins: Struct., Funct., Bioinform.* **2008**, *70*, 1185-1195.
- [9] Krivov, S. V.; Muff, S.; Caffisch, A.; Karplus, M. *J. Phys. Chem. B* **2008**, *112*, 8701-8714.
- [10] Muff, S.; Caffisch, A. *J. Phys. Chem. B* **2009**, *113*, 3218-3226.
- [11] Muff, S.; Caffisch, A. *J. Chem. Phys.* **2009**, *130*, 125104.
- [12] Zheng, W.; Qi, B.; Rohrdanz, M. A.; Caffisch, A.; Dinner, A. R.; Clementi, C. *J. Phys. Chem. B* **2011**, *115*, 13065-13074.
- [13] Rohrdanz, M. A.; Zheng, W.; Maggioni, M.; Clementi, C. *J. Chem. Phys.* **2011**, *134*, 124116.