

1 *Supplementary Materials for*

2

3 **A Suggestion of Converting Protein Intrinsic Disorder to Structural**

4 **Entropy using Shannon's Information Theory**

5

6

7 Hao-Bo Guo^{1,*}, Yue Ma², Gerald A. Taskan³, Hong Qin^{1,4}, Xiaohan Yang^{2,3}, Hong Guo²

8

9 ¹Department of Computer Science and Engineering, SimCenter, University of Tennessee, Chattanooga,
10 TN 37403

11 ²Department of Biochemistry and Cellular and Molecular Biology, University of Tennessee, Knoxville,
12 TN 37996

13 ³Biosciences Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831

14 ⁴Department of Biology, Geology, and Environmental Science, University of Tennessee Chattanooga, TN
15 37403

16

17

18

19 *The supplementary materials include*

20

21 • Fig. S1. Profile of the structural entropy of the residues in the giant protein Human ($C =$

22 34350).

23 • Appendix, on the derivation of the equations that convert the disorder contents to the
24 probabilities of states (with Figures S2 and S3)

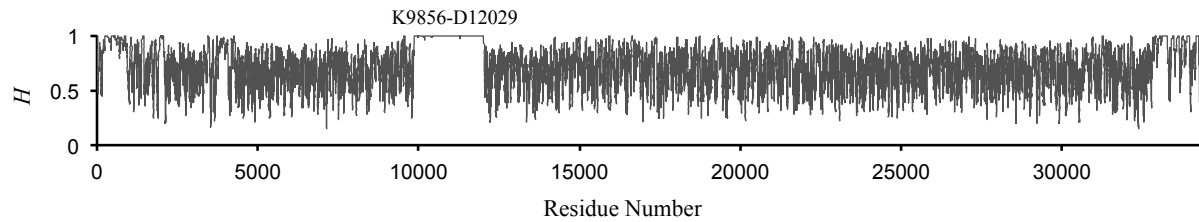
25 • Fig. S4. The exponential, gamma and power law fittings to the structural capacities of
26 the human and *JCVI-Syn3.0* proteomes

27 • Table S1, summarizations of the exponential, gamma and power law fittings of the
28 protein structural capacities of the proteomes studied in this paper

29

30

31



32
 33 **Figure S1.** Profile of the structural entropy of the residues in the giant protein Human ($C = 34350$). Residues
 34 K9856 to D12029 (2174 AA) are a long intrinsically disordered region (IDR) with $H > 0.95$ for all residues.
 35 The composition of residues in this IDR is C: 0, N: 0, A: 129, G: 16, L: 53, I 87, M: 11, V: 331, F: 24, W:
 36 3, S: 35, T: 55, Y: 32, Q: 28, K: 345, R: 64, H: 17, P: 456, D: 13, and E: 475. This region is abundant of
 37 disorder-promoting residues including 914 charged residues (K, R, H, D and E) and 456 P.
 38
 39

40 Appendix

41 In the present paper the protein intrinsic disorder contents at the residue level are used to quantify the
42 structural entropy and information. The quantities obtained therefore is also limited at the residue level,
43 despite that more sophisticated methods might be able to tackle the structural information at higher (such
44 as atomic and electronic) levels.

45 The Shannon equation¹ (eq. 1) might be a reasonable choice in studying the structural entropy of a
46 protein since its structure can be viewed as a linear sequence of amino acids linked by peptide bonds. The
47 function H of the Shannon entropy is statistical and derived from the state probabilities (p_i for the i -th state,
48 $i = 1, \dots, n$, and n is the number of total states) with three original criteria¹⁷ that

- 49 1) H is continuous in p_i ; i.e., p_i could be any value in range of $[0, 1]$ given that $\sum p_i = 1$;
- 50 2) H is a monotonic increasing function when all states are equally distributed with $p_i = 1/n$; it
51 should be noted that H achieves its maximal value of $H_{max} = C = \log n$ in this situation, where C
52 is the capacity;
- 53 3) H is additive, which is true for thermodynamic entropies, too. Shannon's definition came from
54 the statistical considerations; i.e., when the choice of a state was split into two states, the original
55 H should be weighted sum of the two individual values of H .

56 Here, for the structural entropy that concerns the intrinsic order or disorder of proteins, another criterion
57 need be added, i.e.

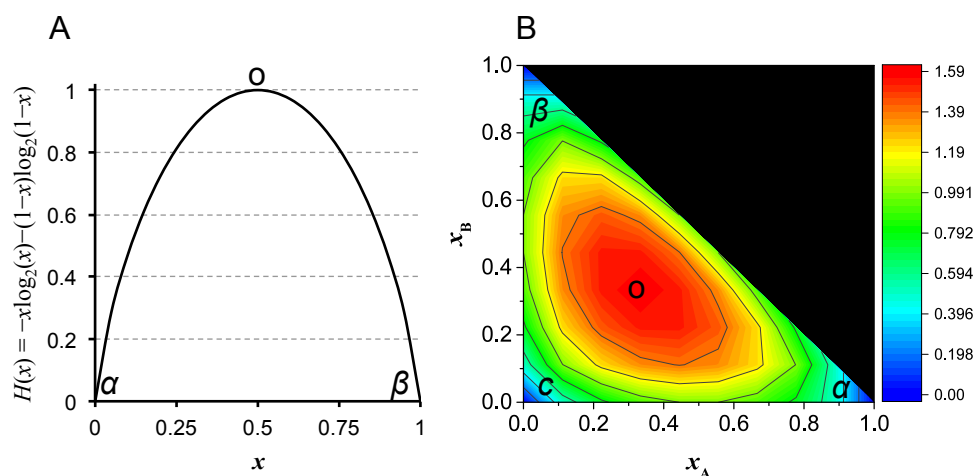
- 58 4) A totally disordered residue contributes the structural entropy of 1, whereas a totally ordered
59 residue contributes zero; the higher the disorder content, the higher the structural entropy a residue
60 has.

61 Intuitively, criterion 4 fits the definition of both thermodynamic and information entropies. In the former,
62 higher entropy corresponds to higher disorder, and in the latter entropy is synonymous to uncertainty. In
63 both definitions the residues with higher disorder contents should have higher structural entropies. It had
64 been proved¹ that the only H that satisfying criteria 1 to 3 is in the form of eq. 1, and therefore, to use this
65 equation to estimate the structural entropy of a protein the disorder contents of all residues must be con-
66 verted to probabilities of all states of the protein, in account of the criterion 4.

67 The disorder predictor gives a vector $\mathbf{d} = (d_1, d_2, \dots, d_L)$ that scores the disorder content of each se-
68 quence of a protein with L residues. The score d_i of the i -th residue distributes in range of $[0,1]$ with 0 for
69 fully ordered and 1 for fully disordered and that in between for a mixed state. However, considering the
70 structural entropy and information we cannot even treat a single residue as a two-state system (i.e., 0 for
71 the ordered and 1 for the disordered states) and apply eq. 1 such as

$$72 \quad H(X) = -x \log_2 x - (1-x) \log_2 (1-x), \quad (S1)$$

73 Where, x is the probability of the first (ordered) state and $(1-x)$ of the second (disordered) state, of that
 74 residue. Eq. S1 symmetrically assigns equal contributions to entropy for both states that fits the criterion 2;
 75 however, it fails to meet the criterion 4. Instead, the ordered and disordered states should respectively have
 76 zero (0) and full (1) contributions, respectively. To fit the criterion 4, we may suppose an imaginary two-
 77 state system as shown in Fig. S3A. The two states termed α ($x = 0$) and β ($x = 1$) contribute equally to the
 78 structural entropy and the entropy $H(x)$ is zero at both extrema. The fully mixed state at $x = 0.5$ has the
 79 maximal entropy of $H(x) = 1$, and this state should be regarded as the disordered state. Similarly, a three-
 80 state (or higher dimension) system may be supposed (Fig. S3B) with probabilities x_A for the α -, x_B for the
 81 β - and x_C for the c -states, respectively, with $\sum_{i=A,B,C} x_i = 1$. The fully mixed state ($x_i = 1/3$) has the maximal
 82 entropy of $H = \log_2 3$.



83
 84 **Figure S2.** Profiles of Shannon function for (A) a two-state system; both α - ($x = 0$) and β -states ($x = 1$) have zero
 85 entropies whereas the state with maximal entropy of 1.0 at $x = (1-x) = 0.5$; (B) a three-state system. The 2D contour
 86 map is a projection onto the probability space of x_A and x_B ; the black region is inaccessible with total probabilities
 87 larger than 1. All extreme states have zero entropies and the mixed state at $x_A = x_B = x_C = 1/3$ has the maximal entropy
 88 of $\log_2 3 = 1.585$.
 89

90 Therefore, the criterion 4 shown above gives two alternative approaches for converting the disorder
 91 contents d to probabilities of states. In the first approach, d is directly used in the estimations, i.e.,

$$92 \quad H(x_i) = d_i, \quad (S2)$$

93 d_i is the disorder content of the i^{th} residue. This approach (direct approach) is equivalent to a two-state
 94 approach and d_i automatically takes the value between 0 and 1, with 0 for the fully ordered and 1 for the
 95 fully disordered, fit well with criterion 4. However, a careful consideration of criterion 2 need be taken
 96 because the two extreme states (0 and 1) contribute unevenly to the entropy. Nevertheless, for a protein
 97 with L residues the maximal entropy or the capacity of the protein is $H_{max} = L$, when all residues are in the
 98 fully disordered state, which is consistent with the total state number of 2^L for the two-state system.

99 The second approach is based on Shannon's equation (Shannon-approach). Considering the two-state
 100 system in Fig. S3A, the α - and β -states (the 0 and 1 states) could be regarded as two representative second-
 101 ary structures. All mixed states between 0 and 1, therefore, have mixed secondary structural characteristics
 102 with the fully mixed state ($x = 0.5$) having the maximal entropy of $\log_2 2 = 1$. The symmetry of Shannon's
 103 function (eq. S1) provides that both states contribute equally to the entropy, and therefore criterion 2 holds.
 104 In this approach, the disorder contents are converted to the probabilities of states using

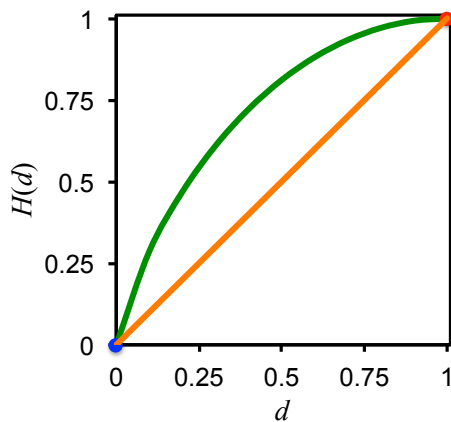
$$105 \quad H(X) = \sum_{i=1}^L -x_i \log_2 x_i - (1-x_i) \log_2 (1-x_i), \quad (S3)$$

$$x_i = d_i/2.$$

106 In both approaches the capacity C , or the maximal entropy H_{max} , of the protein equals to the residue number
 107 L ; i.e., the total number of the states of the protein is $n = 2^L$. The difference between the two approaches is
 108 that the direct-approach gives a linear function of the disorder content (orange in Fig. S4) and the Shannon-
 109 approach is a half function of the Shannon's equation in Eq. S1 (green in Fig. S4). It should be noted from
 110 that the disorder contents might underestimate the structural entropies.

111 The Shannon-approach is adopted in the main text. It should be noted from Fig. S3 that an alternative
 112 approach could be derived from the secondary structure predictions either use a two-state or three-state
 113 system or in higher dimensions. Moreover, this approach could be assisted by molecular dynamics (MD)
 114 simulations by providing an ensemble of configurations from which the probabilities of states could be
 115 extracted, which should be promising because the protein dynamics is involved.

116



117 **Figure S3.** The structural entropy $H(d)$ in function of the intrinsic disorder d . The orange line is from the direct-
 118 approach and the green line is from the Shannon-approach. Blue dot stands for the fully ordered state and red dot for
 119 the fully disordered state. Both profiles are based on two-state systems. In the direct-approach the two extreme states
 120 do not contribute equally to the entropy with the ordered state has entropy of 0 and disordered state has entropy of 1,
 121 respectively. In the Shannon-approach the fully ordered state could be served as either of two extreme states with
 122 entropies of 0, whereas the fully disordered state with entropy of 1 is the equally mixed state of both extreme states.
 123

124 The exponential model with $L = Ae^{bx}$, gamma model with $L = \Gamma^x(x/(n+1); \alpha, \beta)$, and power law model
 125 with $L = Ax^b$ have been used to fit the protein length L in the proteomes. Here x is the serial number of the
 126 protein in the hierarchical rank and n is the total number of proteins in the proteome. A and b are the fre-
 127 quency factors and exponential indexes in the exponential and power law models. The inverse gamma
 128 function was applied in the gamma model and the parameters α and β are calculated via

$$129 \quad \alpha = \left(\sum_{i=1}^n L_i \right)^2 / \left[n \sum_{i=1}^n L_i^2 - \left(\sum_{i=1}^n L_i \right)^2 \right], \quad (S4)$$

$$\beta = n \sum_{i=1}^n L_i^2 / (n-1) \sum_{i=1}^n L_i.$$

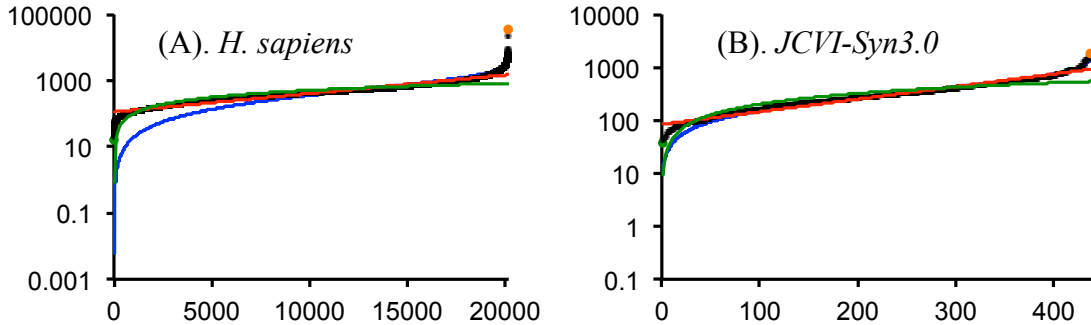
130 The coefficient of determination, R^2 , was calculated using the standard procedure of

$$131 \quad R^2 = 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (L_i - \bar{L})^2}, \quad (S5)$$

132 where, $e_i = f_i - L_i$ is the error for the i^{th} protein, and \bar{L} is the average protein length of the proteome.

133 Fig. S5 shows examples from the human (*H. sapiens*) and bacterial (*JCVI-Syn3.0*) proteomes. The
 134 fitting results of all proteomes assessed in the present work are summarized in Table S1. In all cases, the
 135 exponential model yield fittings with coefficient R^2 larger than 0.9; the gamma model gives good fittings
 136 except for the two animal models surveyed here. The power law model did fit well at the short- L side but
 137 had relatively large deviations at the long- L side. We may therefore use the exponential model for the fitting
 138 of all proteomes.

139



140
 141 **Figure S4.** Distribution of protein length L from (A) *H. sapiens* (human) and (B) *JCVI-Syn3.0* proteomes ranked in a
 142 hierarchical order (black dots) fitted with exponential (red), gamma (blue) and power law (green) models. The hori-
 143 zontal axis is the serious number of the proteins hierarchically ranked by the structural capacity, and the vertical length
 144 represents the structural capacity of the proteins. The proteins with largest and smallest structural capacities are shown
 145 in orange and green dot, respectively.
 146

147 **Table S1.** Fitting of the structural capacity L using different models

Species	Exponential ^a			Power law ^a			Gamma		
	A	b	R^2	A	b	R^2	α	β	R^2
<i>H. sapiens</i>	113.7	1.3E-4	0.939	0.844	0.695	0.814	0.858	654.2	0.792
<i>D. melanogaster</i>	94.8	2.0E-4	0.946	0.628	0.752	0.826	0.768	699.9	0.804
<i>S. cerevisiae</i>	102.9	4.4E-4	0.934	1.347	0.733	0.888	1.664	296.9	0.983
<i>A. thaliana</i>	88.0	9.0E-5	0.933	0.419	0.718	0.893	1.779	227.8	0.968
<i>O. sativa</i>	70.5	6.0E-5	0.969	0.206	0.735	0.837	1.418	265.3	0.986
<i>A. trichopoda</i>	59.8	1.0E-4	0.980	0.497	0.668	0.723	1.153	275.0	0.971
<i>P. patens</i>	46.1	1.0E-4	0.986	0.092	0.835	0.788	1.005	350.3	0.977
<i>Lokiarchaeum</i>	55.7	4.8E-4	0.959	0.929	0.710	0.854	1.517	177.0	0.939
<i>I. hospitalis</i>	80.0	1.5E-3	0.961	6.251	0.575	0.834	2.329	119.5	0.981
<i>N. equitans</i>	77.5	4.0E-3	0.961	10.231	0.586	0.811	1.895	147.8	0.940
<i>JCVI-Syn3.0</i>	84.2	5.5E-3	0.961	9.273	0.669	0.850	1.828	194.8	0.982
<i>Rickettsia</i>	72.9	1.3E-3	0.966	3.987	0.630	0.809	1.681	179.6	0.969
<i>S. elongatus</i>	79.8	9.0E-4	0.957	3.445	0.622	0.857	2.184	139.8	0.991
<i>Mimivirus</i>	81.4	2.5E-3	0.933	4.753	0.690	0.865	1.536	232.3	0.946
<i>Pandoravirus</i>	39.1	1.2E-3	0.990	0.792	0.793	0.793	1.271	203.9	0.980

148 ^a The functions used for the three models are shown above. For both the exponential and power law models A is the
 149 frequency factor (or pre-exponential factor) and b is the exponential index.
 150