

HCAtk and pyHCA: A Toolkit and Python API for the Hydrophobic Cluster Analysis of Protein Sequences.

Tristan Bitard-Feildel^{1, 2}, Isabelle Callebaut^{2, *},

1 Sorbonne Université, UPMC Université Paris 6, CNRS, IBPS, UMR 7238, Laboratoire de Biologie Computationnelle et Quantitative (LCQB), 75005 Paris, France

2 CNRS UMR7590, Sorbonne Université, Université Pierre et Marie Curie – Paris 6 – MNHN – IRD – IUC, Paris, France

* isabelle.callebaut@impmc.upmc.fr

Abstract

1 Motivation: Detecting protein domains sharing no similarity to known domains, as stored in domain
2 databases, is a challenging problem, particularly for unannotated proteomes, domains emerged recently,
3 fast diverging proteins or domains with intrinsically disordered regions.
4 Results: We developed pyHCA and HCAtk, a python API and standalone tool gathering together
5 improved versions of previously developed methodologies, with new functionalities. The developed tools
6 can be either used from command line or from a python API.
7 Availability: HCAtk and pyHCA are available at <https://github.com/T-B-F/pyHCA> under the CeCILL-
8 C license.

9 Introduction

10 The annotation of a protein sequence is very often the first step of many bioinformatics analyses, for
11 instance for studying the function of a gene or the evolution of organisms. Protein domain annotation
12 dominates analyses, describing a protein as a list of blocks corresponding to evolutionary and functional
13 conserved segments. Protein domain families have been extensively compiled through sequence or
14 structure similarity searches and stored in several public databases. These domain databases represent

15 state of the art of our current knowledge of the protein domain universe [11,18]. However, many protein
16 sequences escape, at least partially, domain annotation, particularly in non-model organisms and remain
17 in the so-called dark protein sequence universe [4]. Classical methodological approaches model protein
18 domain families as Hidden Markov Models (HMMs). However, to that aim, sequences need to be clustered
19 and aligned based on the identification of sequence similarities. Therefore, proteins from organisms distant
20 from the species considered in the model, fast diverging proteins, recently emerged domains and domains
21 containing disordered regions, are less likely to be annotated using this methodology [5]. Here, in order to
22 provide a comprehensive view of protein domain architectures, we present a standalone software named
23 HCA toolkit (HCAtk), and its associated python API pyHCA. The package is easily installable and
24 extend our previous developed tools making use of the Hydrophobic Cluster Analysis (HCA) of protein
25 sequences [6,8–10,12,24], with new functionalities. The HCA methodology, based on a two-dimensional
26 representation of protein sequences, highlights clusters of hydrophobic amino acids making up globular
27 domains. More on the HCA methodology can be found in the supplementary materials.

28 **Methods**

29 Seg-HCA [10] was developed to automatically delineate potential “foldable” domains within protein
30 sequences and is the core part of our package. Recently, Piovesan et al. [21] implemented an in-house
31 version of Seg-HCA in FIELDS, which allows to nicely visualize different properties of a protein sequence.
32 Our new version of Seg-HCA was rewritten for speed and a score is now computed, describing the general
33 composition in hydrophobic clusters of the delineated foldable domains. This score is compared to an
34 empirical distribution computed over 734 disordered protein sequences from DisProt v7 [20] to produce
35 a p-value. Figure 1A shows the distributions of scores computed using non redundant sequences of
36 the Protein Data Bank for the set of globular domains and the set of DisProt protein sequences. The
37 resulting p-value can thus be used to evaluate the likelihood of the delineated domains to fold into
38 globular structures. Interestingly, some Seg-HCA domains are reported with a high p-value. A closer
39 inspection revealed these sequences as partially disordered and undergoing possible folding upon binding.
40 A detailed description of scores with some examples is provided in the supplementary material.

41
42 The second methodology included in the package is our TREMOLO-HCA software (Traveling through
43 REMOte homoLOgy) [9]. Using as queries domains delineated using Seg-HCA, remote similarity search
44 is performed against protein sequences from the Uniprot database [29] using HHblits [25]. For each hit,

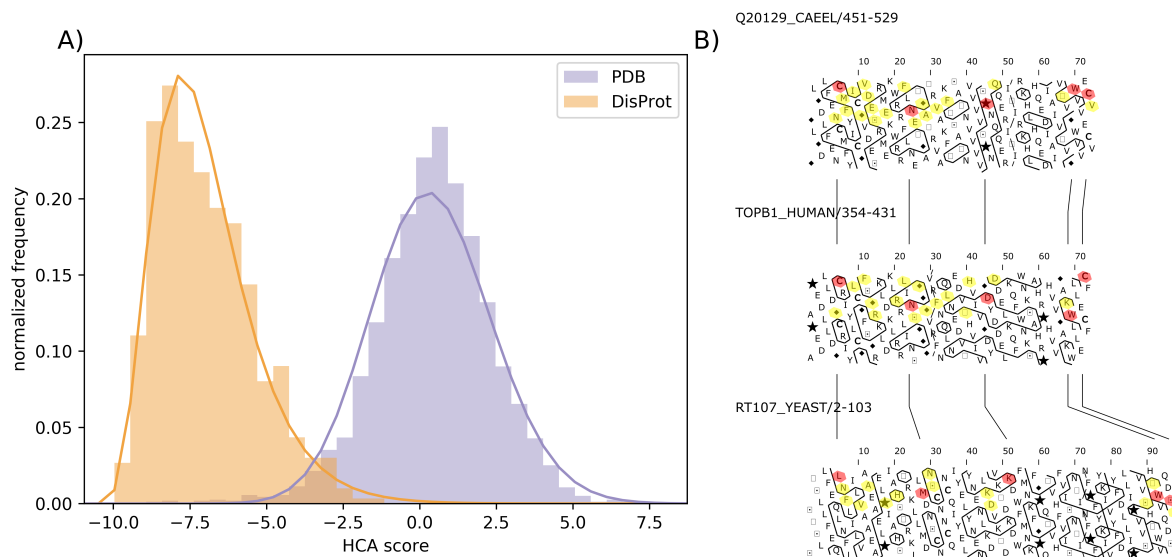


Figure 1: **HCA score and HCA plot example.** Panel A, left, shows the normalized HCA score distribution calculated for protein sequences from DisProt v7 (left, orange – disordered sequences) and from PDB (right, violet - globular domains). The HCA p-value assessing the globularity of delineated foldable segment, is computed using the empirical distribution from DisProt sequences. Panel B, right, shows the HCA plots of three BRCT domains from the Pfam family (PF00533). The aligned protein sequences were used as an input and conserved amino acids can be visualized in red (highly conserved) and yellow, in the context of hydrophobic clusters (HC), in order to evaluate the secondary structure conservation, relatively to the HC shapes.

45 domain arrangement of the Uniprot targets is retrieved from the Interpro database [17]. The final output
46 allows to directly link unknown protein domains, delineated by Seg-HCA, to existing annotations and to
47 analyze these unknown domains in the context of their domain arrangement. The original tool was based
48 on PSI-Blast and the CDD webserver. The new implementations based on HHBlits and Interpro allows a
49 more sensitive detection of protein sequence remote similarity combined with a larger coverage of the
50 protein domain universe thanks to the multiple sources of annotations integrated in Interpro. Several
51 scripts are also provided to easily parse and query the TREMOLO output and to quickly retrieve protein
52 domains of Interpro overlapping the unknown Seg-HCA domains or to retrieve the whole protein domain
53 architectures associated with the Seg-HCA domains.

54

55 Finally, two drawing functionalities were developed. The first functionality allows visualization of
56 hydrophobic clusters of protein sequences, whether these are aligned or not. For each protein provided in
57 a fasta file, an HCA plot is drawn, allowing the quick inspection of the hydrophobic cluster content of
58 a protein sequence, which gives information about its composition in regular secondary structures. A

59 detailed description of the drawing methodology is provided in the Supplementary Material. Moreover,
60 another new functionality was implemented to highlight conservation between aligned protein sequences on
61 their HCA plots. Conserved protein sequence positions can therefore be inspected in the context of their
62 hydrophobic cluster organization (Figure 1B). The second drawing functionality is a new methodology
63 built on the TREMOLO results to easily visualize the known protein domain annotation (from Interpro)
64 and the newly delineated domains in an evolutionary context by using the NCBI taxonomic database.
65 The tree is automatically built by fetching the taxonomic id of the Uniprot target sequences found
66 by TREMOLO thanks to the ete3 python package. The Seg-HCA domains of TREMOLO can then
67 be analyzed in the context of their protein domain arrangement and visualized in terms of taxonomic
68 specificity and domain association (Fig. S1). The HCA toolkit is written in python3 and is provided
69 under the CeCILL-C license agreement. The functions associated with the HCA analyses in the toolkit
70 can also be directly used through a python API and as such can easily be used in other software.

71 Funding

72 This work has been supported by the Agence Nationale de la Recherche (grant number ANR-14-CE10-
73 0021) and the Institut National du Cancer (grant number PLBIO14-299).

74 Conflict of Interest: none declared.

75 Supporting Information

76 Supplementary Figure 1 is accessible at [https://github.com/T-B-F/pyHCA/blob/master/img/Supplementary_](https://github.com/T-B-F/pyHCA/blob/master/img/Supplementary_Fig1.pdf)
77 [Fig1.pdf](https://github.com/T-B-F/pyHCA/blob/master/img/Supplementary_Fig1.pdf).

78 HCA methodology, HCA plot and Seg-HCA

79 HCA hydrophobic clusters, made of strong hydrophobic amino acids (V, I, L, F, M, Y, W), are different
80 from hydrophobic segments as they can incorporate other, non-hydrophobic residues. This property
81 originates from the use of a two-dimensional alpha-helical net, connecting hydrophobic amino acids
82 separated by up to three non-hydrophobic amino acids (or a proline) [12]. Hydrophobic clusters de-
83 fined in this way (with this hydrophobic alphabet and the connectivity distance associated with the
84 α -helix) have been shown to match at best regular secondary structures (α -helices and β -strands) and
85 to constitute hallmarks of folded domains [8, 30]. Sequence segments delineated by Seg-HCA, which

86 correspond to regions where a high density in hydrophobic clusters is detected, have been shown to
87 correspond to domains that have the ability to fold, either in an autonomous way or following contact
88 with partners [5, 10]; these segments are later referred to as HCA domains. The advantage of Seg-HCA
89 for the characterization of the dark proteome is to allow the prediction of these foldable domains from
90 the only information of a single amino acid sequence, without the prior knowledge of homologous sequences.

91

92 Figure S2 presents the Hydrophobic Cluster Analysis (HCA) methodology and indicates how are
93 generated the HCA plots shown in Figure 1B. From an original 1D amino acid sequence (panel A), a 2D
94 plot is created (panel D) by rolling the amino acid sequence around an α -helix (panel B) and cutting the
95 helix along the horizontal axis. The helix forms a plane (panel C) on which every line of amino acids
96 corresponds to an helix turn. The plane is duplicated and the hydrophobic clusters are defined by joining
97 contiguous strong hydrophobic amino acids (V, I, L, F, M, Y, W).

98

99 Regular secondary structure (RSS) elements can be easily visualized on the 2D plot, mainly corre-
100 sponding to hydrophobic clusters, which are separated from each other by loops. Vertical hydrophobic
101 clusters mainly correspond to β -strands and horizontal clusters to α -helices. A dictionary of the most
102 current hydrophobic clusters, established from a comprehensive analysis of experimental 3D structures,
103 can be considered for evaluating the main propensities of hydrophobic clusters towards RSS [8, 24].

104 **HCA score**

105 The HCA score, used to compute a p-value associated with each HCA domain, is defined as follow. Each
106 residue of an HCA segment is associated with a class regarding the residue type and hydrophobicity. Such
107 a residue is either in an hydrophobic cluster and hydrophobic, in an hydrophobic cluster and hydrophilic,
108 or outside an hydrophobic cluster. A value is attributed to each class and the HCA score is computed as
109 follow:

110 with $s(i)$ the function mapping the residue i to each class value.

111

112 Therefore, the HCA score scales with the density in hydrophobic clusters and in hydrophobic residues
113 inside the clusters. As the HCA score calculation motivation is to provide an estimation of the globular
114 character, i.e. the foldability of an HCA domain, the value of each of the three classes was optimized to
115 maximize the separation between the distributions of the HCA scores computed on disordered sequences

126 Examples of disordered regions with low HCA scores.

127 Supplementary Figures 3 and 4 show two examples of protein sequences taken from the left tail of the
128 HCA score distribution of DisProt sequences, displayed in Figure 1A of the main document. For each
129 figure, the HCA plot is drawn on top and the DisProt annotation taken from the DisProt webserver
130 is shown at the bottom. These two sequences have HCA patterns typical of disordered proteins, i.e.
131 proteins having very few hydrophobic residues, often gathered in HCA clusters of small length and spread
132 along the sequence, this one including many proline residues (star symbols). Both proteins regions shown
133 in Fig. S3 and S4 have low percentages of hydrophobic amino acids (13% and 6%).

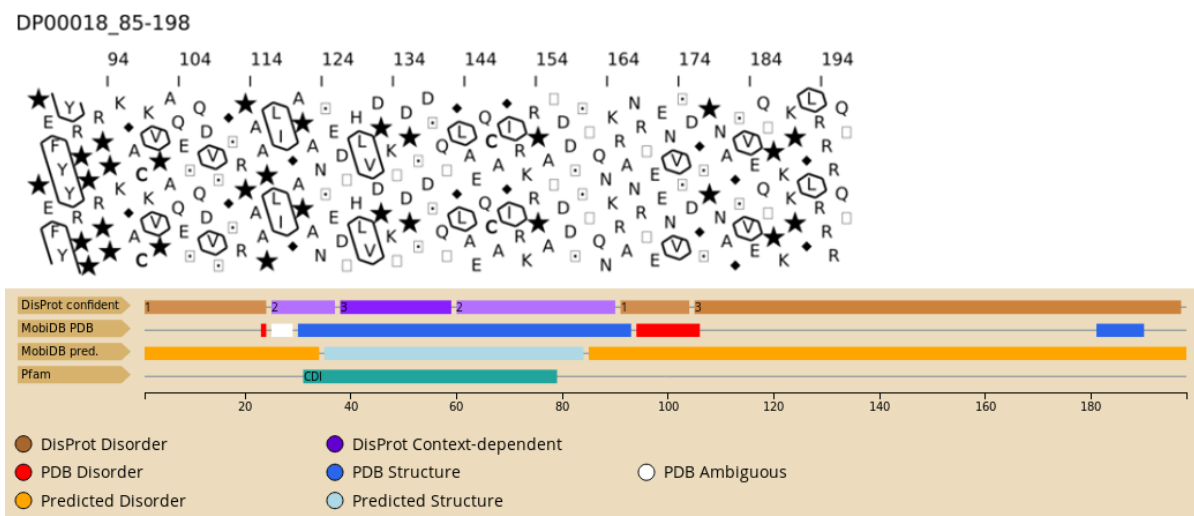


Figure S3: **Disordered regions with low HCA score - example n° 1.** The HCA pattern displayed by the disordered region is typical of non globular regions, with very few clusters and many proline (star symbols). The disordered region corresponds to the sequence segment from amino acids 85 to 198 of the human Cyclin-dependent kinase inhibitor p27(Kip1) (Uniprot P46527).

134 Figure S3 corresponds to the C-terminal sequence of the human Cyclin-dependent kinase inhibitor p27.
135 p27Kip1 controls eukaryotic cell division through interactions with cyclin-dependent kinases [22] and is
136 known as a flexible protein [13], whose stability is associated with phosphorylation. The C-terminal region
137 of p27 has a high flexibility, which provides the molecular basis for the sequential signal transduction
138 conduit that regulates its own degradation and cell division [3,13]).

139 Figure S4 corresponds to the C-terminal sequence of the chicken Histone H5 protein. Histone proteins
140 have well characterized intrinsic disordered regions which are necessary to their biological function [19] and
141 are targets for post-translational modifications recognized by specific readers. Two serine phosphorylation
142 sites have been identified at position 146 and 167. The abundance of lysine and arginine also suggests
143 possible acetylation/methylation sites. On the other hand, the C-terminal domain of chicken Histone H5

144 has a DNA binding motif [23], whose activity requires a high level of intrinsic flexibility.

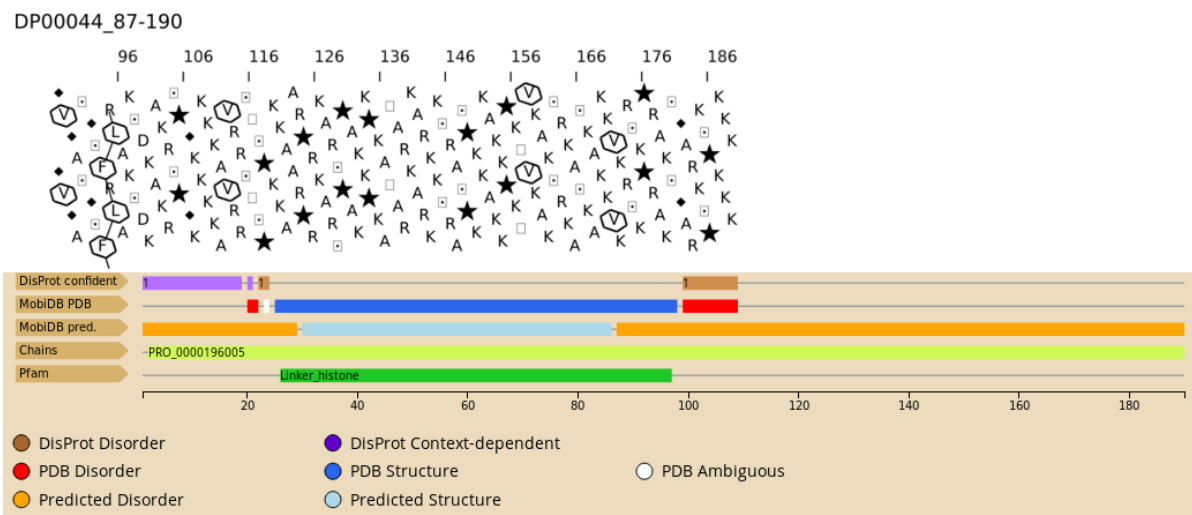


Figure S4: **Disordered regions with low HCA score - example n° 2.** The HCA pattern displayed by the disordered region is typical of non globular regions, with very few clusters and many proline (star symbols). The disordered region corresponds to the sequence segment from amino acids 87 to 190 of the chicken Histone H5 protein (Uniprot P02259).

145 **Examples of disordered regions with high HCA scores.**

146 Figures S5 and S6 show two examples extracted from the right tail of the DisProt HCA score distribution,
147 i.e. proteins with HCA scores similar to the lowest scores of the sequences extracted from PDB (named
148 PDB sequences below). These two examples display HCA patterns including larger hydrophobic clusters
149 (typical of regular secondary structures), as found in globular proteins, but with a slightly lower total
150 content in hydrophobic amino acids (24% for both against 30% expected). The first example concerns a
151 disordered region (amino acids 291 to 352) found in the chicken zing finger FYVE domain-containing
152 protein 9 (UniProt O95405, Fig. S5). The 3D structure of only the FYVE domain has been experimentally
153 characterized, corresponding to the FYVE zinc finger domain (PF01363) (amino acids 663 to 751), the
154 second domain corresponds, from amino acids 1048 to 1400, to a Pfam domain of unknown function
155 (PF11979). The protein regulates the subcellular localization of SMAD2/SMAD3 by recruiting them to
156 the TGF-beta receptor [7,28]. The HCA pattern displayed by the disordered region is similar to patterns
157 observed for foldable regions, suggesting that this small domain is able to fold, at least under particular
158 conditions. The absence of any clear annotation in the N-terminal part of the protein, including two
159 small regions predicted as disordered, but in which a potential order is detected, suggests the presence of
160 an un-detected domain of unknown function [4]).

DP00549_82-134

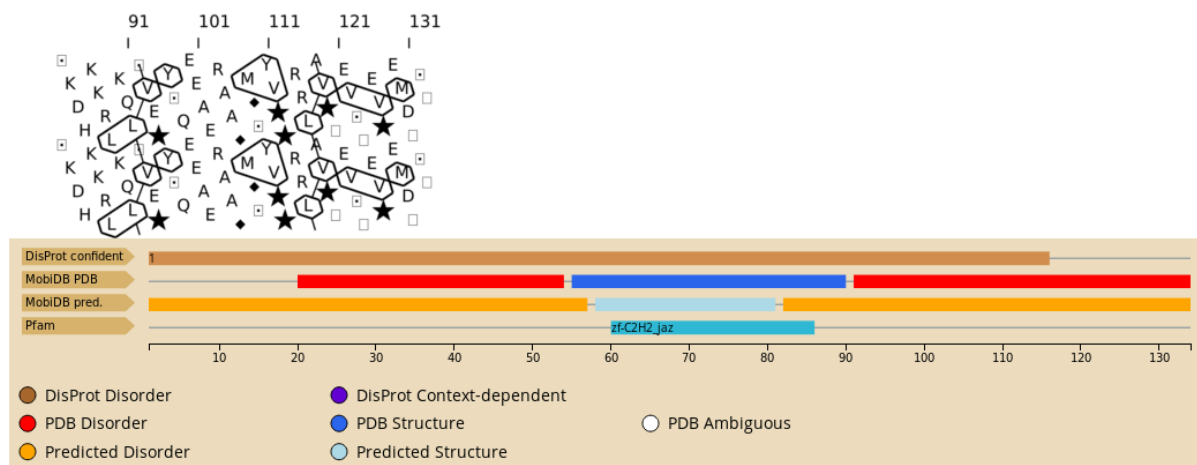


Figure S6: **Disordered regions with high HCA score - example n°2.** The HCA pattern displayed by this disordered region is similar to pattern of foldable regions. Clusters of several hydrophobic residues can be seen close together. The disordered region, (82-134) corresponds to the C-terminal of human Zinc finger protein 593 (UniProt O00488).

169 Examples of PDB sequences with high HCA scores.

170 Figures S7 and S8 show two examples of HCA plots for sequences extracted from the PDB. Figure S7
 171 corresponds to the HCA plot of the *Archeoglobus fulgidus* VapC ribonuclease (Uniprot O28590, amino
 172 acids 1 to 156) whose 3D structure has been solved X-ray crystallography (PDB entry 1W8I). This
 173 ribonuclease is involved in a toxin-antitoxin module with toxin activity [1] and includes one known domain
 174 (amino acids 3 to 127 corresponds to the PFAM domain PIN, PF01850). The corresponding structure
 175 includes 9 long α -helices with 40% of hydrophobic amino acids. The dense network of HCA clusters
 176 visible in Fig. S7 is typical of globular proteins and the long α -helices can be visualized as horizontal
 177 clusters on the 2D HCA plot.

O28590 1-156

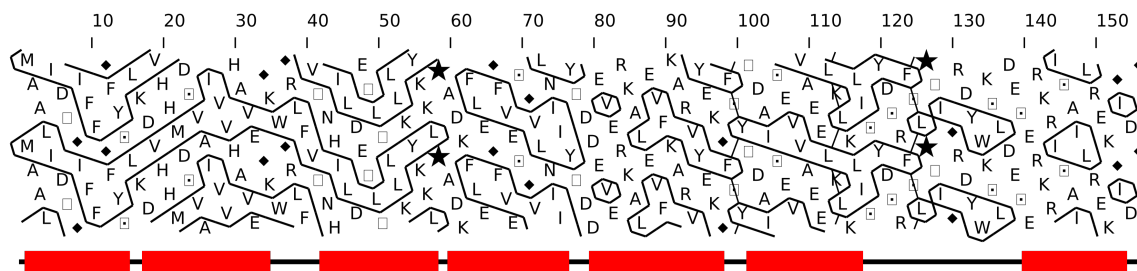


Figure S7: **PDB sequence with an high HCA score – example n°1.** This sequence (Uniprot O28590, PDB 1W8I), including amino acids 1 to 156, corresponds a typical globular protein. (α -helix: red rectangle, annotations extract from the experimental 3D structure 1W8I).

178 Figure S8 is another example of a globular protein HCA plot, i.e. with a high HCA score. The figure
179 shows the HCA plot of the mature mouse interferon beta (Uniprot P01575, amino acids 22 to 181, PDB
180 entry 1WU3 [26]). The protein is made of one domain (Pfam amino acids 27 to 179 (PF00143)), including
181 5 long α -helices with 42% of hydrophobic residues. As for Fig. S7, the protein contains large hydrophobic
182 clusters, typical of regular secondary structures, separated by loops.

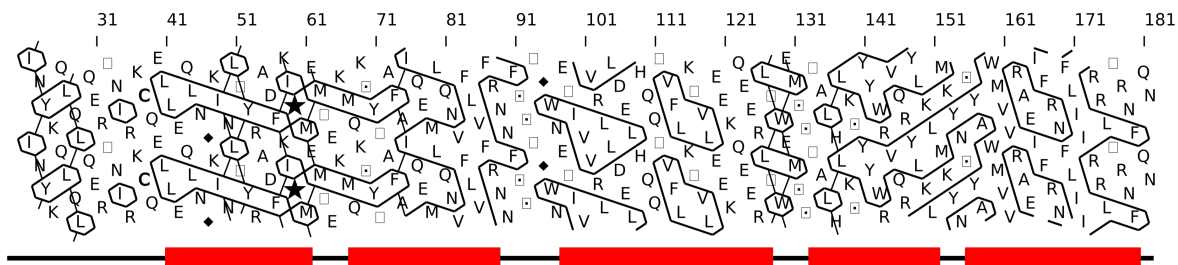


Figure S8: **PDB sequence with an high HCA score – example n°2.** The sequence region (Uniprot P01575, PDB entry 1WU3), including amino acids 22 to 181, corresponds to the PDB structure 1WU3 a typical globular protein. (α -helix: red rectangle, annotations extracted from the experimental 3D structure 1WU3).

183 Examples of PDB sequences with low HCA scores.

184 Figures S9 to S10 are two examples of PDB sequence HCA plots with low HCA scores. Fig. S9 corresponds
185 to the N-terminal domain (amino acids 50 to 174) of the nucleoprotein from human SARS coronavirus
186 (Uniprot protein P59595, PDB entry 2OFZ [16]). This nucleoprotein has RNA binding activity, packaging
187 the positive strand of the human SARS coronavirus RNA genome into a helical ribonucleocapsid [27].
188 The RNA binding activity is mediated by the region encompassing amino acids 45 to 181, such binding
189 activity is usually mediated by a high level of flexibility. The full-length protein is made of one or two
190 protein domains according to the Pfam database (PF00937, from 15 to 378) or the SUPERFAMILY
191 database (SSF110304, from 28 to 181 and SSF103068, from 252 to 365). The 3D structure of the first SSF
192 domain is made of four β -strands (from amino acids 61 to 59, 84 to 91, 102 to 113, and 130 to 135) and
193 one small α -helix (50 to 57), with a large number of flexible loops around the β -sheet core (Saikatendu et
194 al., 2007). According to the coverage of the sequence by large loops, this protein domain has a lower
195 percentage of hydrophobic residues (26%), than the regularly admitted of 30% limit, characteristic of
196 globular domains.

197 Fig. S10 is another example of PDB protein sequence with low HCA score. The HCA plot represents
198 a sequence segment (from amino acids 500 to 629) of the the *Staphylococcus aureus* surface protein G

P59595 50-174

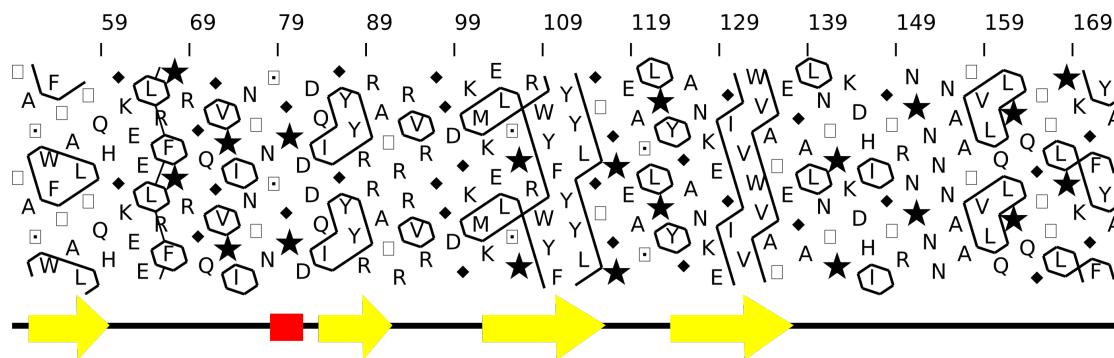


Figure S9: **PDB sequence with a low HCA score – example n°1**. The sequence region (Uniprot P5995), amino acids 49 to 174, PDB structure 2OFZ), corresponds to a RNA-binding domain, having large, flexible loops and a few regular secondary structures, constituting a β -sheet core. (α -helix: red rectangle β -strands: yellow arrows, annotation extracted from the experimental 3D structure 2OFZ).

199 (SasG) (Uniprot Q2G2B2 sequence, PDB entry 5DBL). The full-length protein is made of 19 domains.
200 The sequence starts with a signal peptide motif, followed by pairs of G5 domain/E domain (Pfam
201 PF04650, PF17041) and ends by a cell wall anchor domain (PF00746). Amino acids 501 to 548 and 547
202 to 629 corresponds to a pair of E domain/G5 domain. The G5 domain has only a few conserved amino
203 acids and is supposed to have an adhesive function [2]. As assessed by the presence of small clusters and
204 a relatively weak percentage in strong hydrophobic amino acids, approximately one half of the SasG
205 repetitive regions are intrinsically unfolded in isolation, but fold in the context of neighboring folded
206 G5 domains, highlighting the role of the intrinsically disordered region of the E/G5 domain pair as a
207 key factor for the cooperative folding multidomain protein [14]. Once folded, the two domains form an
208 elongated structure, made of small beta strands which correspond on the HCA plot to small clusters.
209 The small β -strands form triplets-stranded β -sheets connected by collagen-like triple helical regions. In
210 this particular case, several threonine are found included in β -beta strand.

Q2G2B2 500-629

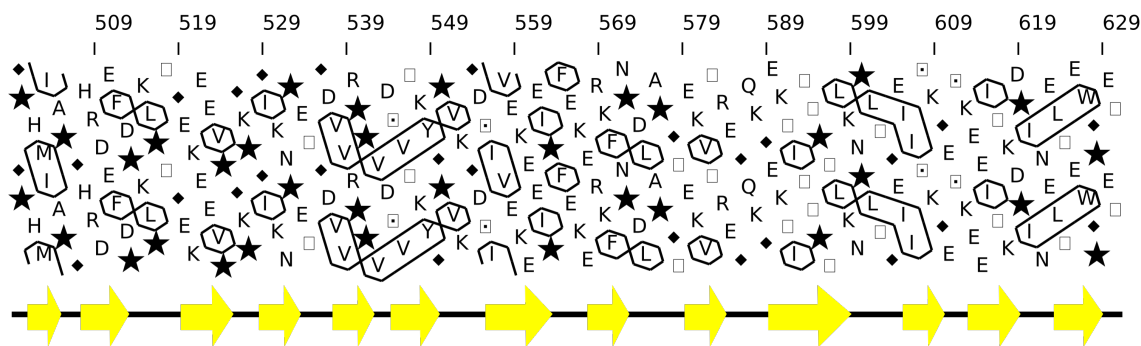


Figure S10: **PDB sequence with low HCA score – example n°2.** The sequence region, (Uniprot Q2G2B2, amino acids 500 to 629, PDB structure entry 5DBL), corresponds to a pair of E/G5 domains of the *S. aureus* surface protein G. (β -strand: yellow arrow, annotations extracted from the experimental 3D structure 5DBL).

References

1. V. L. Arcus, J. L. McKenzie, J. Robson, and G. M. Cook. The pin-domain ribonucleases and the prokaryotic vapbc toxin-antitoxin array. *Protein Engineering Design and Selection*, 24(1–2):33–40, Jan 2011.
2. A. Bateman, M. T. G. Holden, and C. Yeats. The g5 domain: a potential n-acetylglucosamine recognition domain involved in biofilm formation. *Bioinformatics*, 21(8):1301–1303, 2005.
3. E. A. Bienkiewicz, J. N. Adkins, and K. J. Lumb. Functional consequences of preorganized helical structure in the intrinsically disordered cell-cycle inhibitor p27kip1. *Biochemistry*, 41(3):752–759, 2002.
4. T. Bitard-Feildel and I. Callebaut. Exploring the dark foldable proteome by considering hydrophobic amino acids topology. *Scientific Reports*, 7:41425, 2017.
5. T. Bitard-Feildel, M. Heberlein, E. Bornberg-Bauer, and I. Callebaut. Detection of orphan domains in drosophila using “hydrophobic cluster analysis”. *Biochimie*, 119:244–253, 2015.
6. I. Callebaut, G. Labesse, P. Durand, A. Poupon, L. Canard, J. Chomilier, B. Henrissat, and J. P. Mornon. Deciphering protein sequence information through hydrophobic cluster analysis (hca): current status and perspectives. *Cellular and molecular life sciences*, 53(8):621–645, Sep 1997.

7. R. M. Durbin, D. L. Altshuler, R. M. Durbin, G. R. Abecasis, D. R. Bentley, A. Chakravarti, A. G. Clark, F. S. Collins, F. M. De La Vega, P. Donnelly, and et al. A map of human genome variation from population-scale sequencing. *Nature*, 467(7319):1061–1073, Oct 2010.
8. R. Eudes, K. Le Tuan, J. Delettré, J.-P. Mornon, and I. Callebaut. A generalized analysis of hydrophobic and loop clusters within globular protein sequences. *BMC Structural Biology*, 7(1):2, 2007.
9. G. Faure and I. Callebaut. Comprehensive repertoire of foldable regions within whole genomes. *PLoS Computational Biology*, 9(10):e1003280, Oct 2013.
10. G. Faure and I. Callebaut. Identification of hidden relationships from the coupling of hydrophobic cluster analysis and domain architecture information. *Bioinformatics*, 29(14):1726–33, Jul 2013.
11. R. D. Finn, P. Coghill, R. Y. Eberhardt, S. R. Eddy, J. Mistry, A. L. Mitchell, S. C. Potter, M. Punta, M. Qureshi, A. Sangrador-Vegas, and et al. The pfam protein families database: towards a more sustainable future. *Nucleic acids research*, 44(D1):D279–85, Jan 2016.
12. C. Gaboriaud, V. Bissery, T. Benchetrit, and J. Mornon. Hydrophobic cluster analysis: An efficient new way to compare and analyse amino acid sequences. *FEBS Letters*, 224(1):149–155, Nov 1987.
13. C. A. Galea, A. Nourse, Y. Wang, S. G. Sivakolundu, W. T. Heller, and R. W. Kriwacki. Role of intrinsic flexibility in signal transduction mediated by the cell cycle regulator, p27 kip1. *Journal of molecular biology*, 376(3):827–38, 2008.
14. D. T. Gruszka, C. A. T. F. Mendonça, E. Paci, F. Whelan, J. Hawkhead, J. R. Potts, and J. Clarke. Disorder drives cooperative folding in a multidomain protein. *Proceedings of the National Academy of Sciences of the United States of America*, 113(42):11841–11846, 2016.
15. P. L. Hayes, B. L. Lytle, B. F. Volkman, and F. C. Peterson. The solution structure of znf593 from homo sapiens reveals a zinc finger in a predominately unstructured protein. *Protein Science*, 17(3):571–576, 2008.
16. M. W. Martin, J. Newcomb, J. J. Nunes, J. E. Bemis, D. C. McGowan, R. D. White, J. L. Buchanan, E. F. Dimauro, C. Boucher, T. Faust, and et al. Discovery of novel 2,3-diarylfuro[2,3-b]pyridin-4-amines as potent and selective inhibitors of lck: Synthesis, sar, and pharmacokinetic properties. 17:2299–2304, 2007.

17. A. Mitchell, H.-Y. Chang, L. Daugherty, M. Fraser, S. Hunter, R. Lopez, C. McAnulla, C. McMennamin, G. Nuka, S. Pesseat, and et al. The interpro protein families database: the classification resource after 15 years. *Nucleic acids research*, 43(Database issue):D213–21, Jan 2015.
18. M. E. Oates, J. Stahlhacke, D. V. Vavoulis, B. Smithers, O. J. L. Rackham, A. J. Sardar, J. Zaucha, N. Thurlby, H. Fang, and J. Gough. The superfamily 1.75 database in 2014: a doubling of data. *Nucleic acids research*, 43(Database issue):D227–33, Jan 2015.
19. Z. Peng, M. J. Mizianty, B. Xue, L. Kurgan, and V. N. Uversky. More than just tails: intrinsic disorder in histone proteins. *Molecular BioSystems*, 8(7):1886, 2012.
20. D. Piovesan, F. Tabaro, I. Mičetić, M. Necci, F. Quaglia, C. J. Oldfield, M. C. Aspromonte, N. E. Davey, R. Davidović, Z. Dosztányi, and et al. Disprot 7.0: A major update of the database of disordered proteins. *Nucleic Acids Research*, 45(D1):D219–D227, 2017.
21. D. Piovesan, I. Walsh, G. Minervini, and S. C. Tosatto. Fells: fast estimator of latent local structure. *Bioinformatics*, 33(12):1889–1891, Jun 2017.
22. K. Polyak, M. Lee, H. Erdjument-Bromage, and A. Koff. Cloning of p27 kip1, a cyclin-dependent kinase inhibitor and a potential mediator of extracellular antimitogenic signals. *Cell*, 76:59–66, 1994.
23. V. Ramakrishnan, J. T. Finch, V. Graziano, P. L. Lee, and R. M. Sweet. Crystal structure of globular domain of histone h5 and its implications for nucleosome binding. *Nature*, 362(6417):219–223, 1993.
24. J. Rebehmed, F. Quintus, J.-P. Mornon, and I. Callebaut. The respective roles of polar/nonpolar binary patterns and amino acid composition in protein regular secondary structures explored exhaustively using hydrophobic cluster analysis. *Proteins: Structure, Function, and Bioinformatics*, 84(5):624–638, May 2016.
25. M. Remmert, A. Biegert, A. Hauser, and J. Söding. Hhblits: lightning-fast iterative protein sequence searching by hmm-hmm alignment. *Nature Methods*, 9(2):173–175, 2011.
26. T. Senda, S.-i. Saitoh, and Y. Mitsui. Refined crystal structure of recombinant murine interferon- β at 2.15 Å resolution. *J. Mol. Biol.*, 253:187–207, 1995.

27. S. Stertz, M. Reichelt, M. Spiegel, T. Kuri, L. Martínez-sobrido, A. García-sastre, F. Weber, and G. Kochs. The intracellular sites of early replication and budding of sars-coronavirus. *Virology*, 361:304–315, 2007.
28. T. Tsukazaki, T. A. Chiang, A. F. Davison, L. Attisano, and J. L. Wrana. Sara, a fyve domain protein that recruits smad2 to the $\text{tgf-}\beta$ receptor. *Cell*, 95(6):779–791, 1998.
29. UniProt Consortium. Uniprot: a hub for protein information. *Nucleic acids research*, 43(Database issue):D204–12, Jan 2015.
30. S. Woodcock, J. P. Mornon, and B. Henrissat. Detection of secondary structure elements in proteins by hydrophobic cluster analysis. *Protein engineering*, 5(7):629–35, Oct 1992.