

# Structural and functional analysis of “non-smelly” proteins

Jing Yan,<sup>1</sup> Jianlin Cheng,<sup>2</sup> Lukasz Kurgan,<sup>3\*</sup> and Vladimir N. Uversky<sup>4,5,\*</sup>

<sup>1</sup>Department of Electrical and Computer Engineering, University of Alberta, Edmonton, Canada

<sup>2</sup>Department of Electrical Engineering and Computer Science, University of Missouri – Columbia, USA

<sup>3</sup>Department of Computer Science, Virginia Commonwealth University, Richmond, USA

<sup>4</sup>Department of Molecular Medicine and USF Health Byrd Alzheimer's Research Institute, Morsani College of Medicine, University of South Florida, Tampa, FL 33647, USA.

<sup>5</sup>Protein Research Group, Institute for Biological Instrumentation of the Russian Academy of Sciences, 142290, Pushchino, Moscow region, Russia.

Emails      JY: [jyan5@ualberta.ca](mailto:jyan5@ualberta.ca)  
                JC: [chengji@missouri.edu](mailto:chengji@missouri.edu)  
                LK: [lkurgan@vcu.edu](mailto:lkurgan@vcu.edu)  
                VNU: [vuversky@health.usf.edu](mailto:vuversky@health.usf.edu)

\* **Corresponding Authors:** LK, Phone: +1-804-827-3986; Fax: +1-804-828-2771; Email: [lkurgan@vcu.edu](mailto:lkurgan@vcu.edu); 401 West Main Street, Room E4225, Richmond, Virginia 23284, USA; VNU, Phone: +1-813-974-5816; Fax: +1-813-974-7357; Email: [vuversky@health.usf.edu](mailto:vuversky@health.usf.edu); 12901 Bruce B. Downs Blvd., MDC07, Tampa, Florida 33612, USA

## Highlights

- Cysteine-depleted proteins are abundant in all domains of life
- Prokaryotes are significantly enriched in cysteine-depleted proteins compared to eukaryotes
- Only about 0.05% of proteins are depleted in aromatic residues and cysteine
- Proteins depleted in aromatic residues and cysteine have high levels of intrinsic disorder
- Organisms with higher levels of cysteine-depleted proteins have higher levels of the intrinsic disorder
- “Non-smelly” proteins are involved in translation, transcription, nucleosome assembly, protein folding, and transmembrane transport functions

## Abstract

Cysteine and aromatic residues are major structure-promoting residues. We assessed the abundance, structural coverage, and functional characteristics of the “non-smelly” proteins; i.e., proteins that do

not contain cysteine residues (C-depleted) or cysteine and aromatic residues (CFYWH-depleted), across 817 proteomes from all domains of life. The analysis revealed that although these proteomes contained significant levels of the C-depleted proteins, with prokaryotes being significantly more enriched in such proteins than eukaryotes, the CFYWH-depleted proteins were relatively rare, accounting for about 0.05% of proteomes. Furthermore, CFYWH-depleted proteins were virtually never found in PDB. Depletion in cysteine and in aromatic residues was associated with the substantially increased intrinsic disorder levels across all domains of life. Archaeal and Eukaryotic organisms with higher levels of the C-depleted proteins were shown to have higher levels of the intrinsic disorder and lower levels of structural coverage. We also showed that the “non-smelly” proteins typically did not independently fold into monomeric structures, and instead they fold by interacting with nucleic acids as constituents of the ribosome and nucleosome complexes. They were shown to be involved in translation, transcription, nucleosome assembly, transmembrane transport, and protein folding functions, all of which are known to be associated with the intrinsic disorder. Our data suggested that, in general, structure of monomeric proteins is crucially dependent on the presence of cysteine and aromatic residues.

**Keywords:** Intrinsically disordered proteins; cysteine-depleted proteins; nucleic acid binding proteins; proteins depleted in cysteine and aromatic residues; protein structure; protein function

## Introduction

It is accepted now that intrinsically disordered proteins (IDPs) and hybrid proteins containing ordered domains and functionally important intrinsically disordered proteins regions (IDPRs) occupy a significant part of any proteome across all kingdoms of life and viruses [1-6], being especially abundant in eukaryotes [7,2]. Under physiological conditions, IDPs/IDPRs lack rigid 3D structure and therefore are typically not amenable to experimental structure determination by X-ray crystallography [8-10], which is by far the most commonly used technology to solve protein structures. As a result, they are considered as major constituents of the dark proteome [11,12,8]. While being disordered as a whole or in localized regions, these proteins have a number of important biological roles, especially in transcriptional and translational regulation, splicing, and signaling via cellular protein networks [13-15]. Furthermore, enhanced structural plasticity and exceptional spatiotemporal heterogeneity of IDPs/IDPRs define their mosaic structures, where different regions are disordered to different degrees. IDPs contain a multitude of potentially foldable, partially foldable, differently foldable or not foldable at all segments playing different roles in protein functionality [16,17], and even containing ordered regions that need to undergo order-to-disorder transition in order to make protein active [16,18,19]. In cellular protein-protein interaction networks, IDPs/IDPRs often play a role of hubs [20-24] that are engaged in promiscuous interactions and regulate the structural and functional integrity of these networks [25,26,15]. Furthermore, because of this binding promiscuity [27] and the ability to gain very different structures at binding to different partners [28], IDPs/IDPRs can “rewire” protein-protein interaction networks in response to environmental changes [29].

Systematic comparative analyses of amino acid sequences of ordered proteins and IDPs revealed the presence of numerous important differences [30-34,13]. For examples, extended IDPs/IDPRs from different kingdoms of life were shown to be rich in polar and charged amino acids and deficient in hydrophobic residues [30,35,34,36]. This also resulted in the elaboration of the concept of “order-promoting” (C, W, Y, I, F, V, L, H, T, and N) and “disorder-promoting” residues (A, G, D, M, K, R, S, Q, P, and E); i.e., residues more commonly found in ordered and disordered proteins/regions, respectively [37]. Because of the high relative enrichment of the amino acid sequences of ordered proteins and domains in cysteine, tryptophan, tyrosine, phenylalanine, and histidine, these residues are typically considered as strong order-promoting residues. Based on these observations, we hypothesized that structure and functionality of proteins can be noticeably dependent on the presence cysteine and aromatic residues in their amino acid sequences. One can argue that cysteine is important for protein structural stability only when another cysteine is present in the same chain, to enable disulfide bond formation. Observations below provide important evidence that this is not always correct. Intramolecular disulfide bonds are surely important stabilizing factors. For example, proteins and peptides containing cystine knot, which is a rotaxane-like structural motif containing three disulfide bridges, where a polypeptide region between two of those disulfides forms a loop, through which a third disulfide bond is threaded, are known to show a particularly high degree of structural stability [38,39]. There are also numerous examples in the literature, where the importance of intramolecular disulfide bonds for protein thermal stability was demonstrated (as systemized in [40]). As a result, introduction of additional disulfide bonds is considered as an attractive protein engineering strategy for generating proteins (e.g., antibodies) with enhanced conformational stability [41]. Furthermore, dysregulated cellular redox conditions leading to the alterations in the formation of native disulfide bonds are directly linked to various human diseases [42]. However, even a single cysteine contributes to protein structure stability, as it can be engaged in the intermolecular disulfide bond, or can exist as a free thiol and serve as a part of a protein catalytic site, or as a site of various posttranslational modifications (e.g., S-hydroxylation (S-OH), disulfide bond formation, phosphorylation, S-acylation, S-prenylation, protein splicing, N-acetylation, N-ADP-ribosylation, amidation, S-archaeol cysteine, cysteine sulfinic acid (-SO<sub>2</sub>H) formation, methylation, N-myristoylation, nitrosylation, N-palmitoylation, S-palmitoylation, and S-glutathionylation) [43]. Furthermore, a single cysteine can be used for specific coordination of various ligands, e.g., metal ions. In fact, cysteine is known to show high affinity toward zinc ions (Zn<sup>2+</sup>), and the resulting cysteine-Zn<sup>2+</sup> complexes are important for protein structure, catalysis, and regulation [44], as seen in the CH<sub>3</sub>-type zinc finger proteins [45,46] and in redox switches [44]. On the other hand, CD<sub>3</sub> motifs serves as a Mn<sup>2+</sup> coordination group [47].

To check the validity of the hypothesis that the presence cysteine and aromatic residues is crucial for protein structure, we conducted here a comprehensive bioinformatics analysis of the “non-smelly” proteins; i.e., proteins depleted in cysteine and aromatic residues. Since cysteines are known to smell like rotten eggs [48] and since the side chains of W, Y, F, and H are aromatic (i.e., they contain aromatic ring systems, which are stable, cyclic, planar compounds with a ring of resonance bonds and which, unlike pure saturated hydrocarbons, might have specific odors/aroma), the proteins depleted in these residues are defined here as “non-smelly”. In this study, we assembled a dataset of “non-smelly” proteins found in 817 complete proteomes, and also looked for such proteins in the Protein Data Bank (PDB) [49,50].

# Materials and methods

## Datasets

We analyzed a dataset of 817 complete proteomes, which are defined as collections of proteins encoded by the fully sequenced genome of a specific organism. We obtained these proteomes from the UniProt resource [51,52]. They cover 276,733 proteins from 64 Archaeal organisms, 5,077,609 proteins from 552 Bacterial organisms and 4,208,817 proteins from 201 Eukaryotic organisms, for the total of 9,563,159 proteins. A complete list of the considered species is given in the Supplementary Materials. Table 1 provides a breakdown of these organisms into specific kingdoms/phyla.

We also examined proteins with solved structures collected from PDB. We limited our analysis to the wild-type protein chains that have the expression tags removed and that exclude peptides (chain length > 30 residues), and which cover majority of the corresponding full protein chain from UniProt (>60% coverage). We clustered the sequences of the considered PDB structures at 100% identity to remove duplicates. These steps ensure compatibility with the proteome-level analysis. We collected 99,461 chains that satisfy the aforementioned criteria (*PDB* dataset), as well as two of its subsets that include 50,301 chains that are in complex with nucleic acids (*PDB NA* dataset) and 7,413 monomers; i.e., single-chain structures that do not interact with other proteins and nucleic acids (*PDB monomer* dataset).

## Computation of structural characteristics

We used computational methods to quantify content of putative intrinsic disorder (the fraction of residues that are predicted to be intrinsically disordered) and the current structural coverage (the fraction of proteins for which structure is available) on the whole-proteome scale. We evaluated the quality of the disorder content predictions using a large benchmark dataset that was recently used in [53,54] and which was originally published in [55]. We quantify the predictive quality by computing the mean absolute error (MAE) and Pearson correlation coefficient (PCC) between the disorder content predicted with the consensus and the native disorder content. The resulting MAE = 5.5% and PCC = 0.43, which suggests that the content predictions are relatively accurate and correlated with the native disorder content. Our consensus secures similar results for the subset of the benchmark proteins that have cysteine (MAE = 4.9% and PCC = 0.46) and that have above average cysteine content (MAE = 5.0% and PCC = 0.42).

Recent studies demonstrate that intrinsic disorder can be accurately predicted from protein sequences [55-59]. Furthermore, consensus-based approaches that combine outputs of several disorder predictors were shown to provide more accurate predictions when compared to single predictors [60-62]. For instance, consensus predictors more precisely quantify the disorder content (fraction of the disordered residues in a given protein sequence), reducing error by about 4% when compared to single predictors [61]. We applied a consensus of five complementary predictions produced by two popular tools, IUPred [63] and ESpritz [64]. They include results produced with two versions of the IUPred method, which

were designed to predict long ( $\geq 30$  consecutive residues) and short disordered regions, and three versions of ESpritz that focus on the three types of annotations of disorder using: DisProt database [65,66], crystal structures from PDB, and NMR structures from PDB. These tools are characterized by competitive levels of predictive quality [57,55] and short runtime, which is critical to facilitate processing of over 9.5 million protein sequences. The consensus prediction of disorder requires that at least 3 out of 5 predictions indicate intrinsic disorder. The same consensus was applied in several related studies [2,67-72,8]. Our methodology is also similar to the consensus-derived putative disorder in MobiDB [73,74] and D<sup>2</sup>P<sup>2</sup> [16] databases. We calculated the disorder content of a given dataset of proteins (e.g., proteome) which is defined as a fraction of residues predicted as disordered among all residues in that dataset.

We estimated the current structural coverage based on a computationally tractable approach proposed in [75] and recently used in [8,2,76]. For each protein we ran three rounds of PSI-BLAST [77] searches against the sequences of protein structures from PDB. A given proteins sequence that has  $>50$  residues in length is annotated as having structure if it registers a hit in PDB with the E-value  $<0.001$ . In other words, structurally solved proteins are assumed to have at least one long segment of residues (representing at least one domain) that is sufficiently similar to a sequence of an already solved structure. The structural coverage of a proteome is defined as the fraction of the structurally solved sequences among all sequences in this proteome. Research shows that such PSI-BLAST-based estimates provide relatively accurate results. For instance, a similar PSI-BLAST-based approach failed to find templates (similar sequences that are structured) for only 3 out of 120 target proteins in CASP9 [78]. We recognize that there are more precise approaches to estimate structural coverage that are capable of finding remote homologs, such as I-TASSER [79,80], HHpred [81,82], and MODELLER [83,84]. However, these tools could not be scaled to the size of our dataset. To the best of our knowledge, the largest such attempt is the MODBASE resource that covers only 76 organisms [85]. We note that our estimates of the structural coverage are slightly underestimated by inadequately considering remote homologs. Nevertheless, this bias should be equally distributed across different proteomes, allowing us to perform comparative analyses between the corresponding organisms and domains of life.

## **Functional annotations using GO terms**

We annotated protein functions and cellular locations that are associated with the proteins depleted in major order-promoting residues, such as cysteine, phenylalanine, tyrosine, tryptophan, and histidine. Such proteins were split into two groups, depleted in cysteine (C-depleted) and depleted in cysteine and aromatic residues: phenylalanine, tyrosine, tryptophan, and histidine (CFYWH-depleted). The corresponding functions/locations are significantly enriched in these proteins sets when compared with the proteins from the same domain of life. This analysis relies on the GO terms [86] collected from the UniProt resource. We excluded annotations with “potential”, “probable” and “by similarity” qualifiers that are generated using computer-predictions or indirect experimental evidence. We evaluated magnitude and statistical significance of the differences in the rates of occurrence of GO terms between the C-depleted (or CFYWH-depleted) proteins and a generic set of proteins in the same domain of life by following protocols defined in earlier related analyses [2,3,71]. This analysis was performed for each of the three types of GO terms: cellular components, biological processes and molecular functions. We

randomly selected half of the GO-annotated chains for a given C-/CFYWH-depleted protein set and compared them with the same number of chains/residues drawn at random from the same taxonomic domain. We ensured that proteins drawn from the same domain of life have the same chain length (with  $\pm 10\%$  tolerance), since the amount of intrinsic disorder, which indirectly affects protein function and location, is dependent on the chain length [87]. We repeated this 10 times and evaluated significance of the differences in the 10 sets of counts for the corresponding GO terms. If these measurements are normal, based on the Anderson-Darling test at 0.05 significance, then we applied the paired *t*-test (proteins sets are paired to match chain lengths) to evaluate the statistical significance of differences; otherwise we utilized the non-parametric paired Wilcoxon rank sum test. We considered only the differences with *p*-value<0.001 which also have large magnitude, i.e., the average enrichment in the C-/CFYWH-depleted protein set must be larger than 30%. We analyzed the enrichment of GO terms for the entire set of C-/CFYWH-depleted proteins as well as for the subsets of fully disordered C-/CFYWH-depleted proteins.

## Results and Discussion

### Abundance of C-depleted and CFYWH-depleted proteins

We measured fraction of the C-depleted and the CFYWH-depleted proteins across the 817 proteomes and among the proteins in the three PDB-derived datasets. Table 1 summarizes these values across each domain of life and several larger kingdoms and phyla. About 28% proteins in Archaea, 19% in Bacteria and 8% in Eukaryota do not have cysteines. While we observe substantial differences in the abundance of the C-depleted proteins across the three domains of life while, these values are consistent across the kingdoms and phyla within each domain of life (Table 1). This observation suggests that this trend is broadly associated with domains of life. Only about 0.06% proteins in Archaea and 0.04% in Bacteria and Eukaryota do not have cysteine and aromatic residues. Table 1 shows that the abundance of the CFYWH-depleted proteins is similar across Bacteria and Eukaryota, with Archaea having somehow elevated levels of such proteins.

**Table 1.** Amount of cysteine-depleted (C-depleted) and cysteine and aromatic residues-depleted (CFYWH-depleted) proteins, disorder content and structural coverage for the considered 817 proteomes. We report median of these per-proteome measurements for the entire domains of life (in bold font) and several larger kingdoms/phyla. The domains of life and the phyla/kingdoms within each domain are arranged according to their overall fraction of C-depleted proteins.

Domain of life	Kingdom/phylum of life	Number of Species	Median (per proteome) fraction of C-depleted proteins [%]	Median (per proteome) fraction of CFYWH-depleted proteins [%]	Median (per proteome) disorder content [%]	Median (per proteome) structural coverage [%]
Archaea	<b>All</b>	<b>64</b>	<b>28.48</b>	<b>0.06</b>	<b>5.88</b>	<b>53.10</b>
	Crenarchaeota	17	34.05	0.06	3.00	53.56
	Other	4	24.71	0.11	4.45	55.85
	Euryarchaeata	43	19.81	0.07	5.10	53.96
Bacteria	<b>All</b>	<b>552</b>	<b>18.65</b>	<b>0.04</b>	<b>5.45</b>	<b>55.85</b>
	Firmicutes	76	22.62	0.04	4.60	60.07
	Actinobacteria	70	21.97	0.07	10.30	59.45
	Other	108	19.20	0.03	4.70	57.70
	Bacteroidetes	44	18.67	0.02	3.45	54.02
	Proteobacteria	254	16.99	0.04	5.80	61.37
Eukaryota	<b>All</b>	<b>201</b>	<b>7.87</b>	<b>0.04</b>	<b>19.70</b>	<b>47.70</b>
	Fungi	84	9.64	0.03	21.55	46.33
	Viridiplantae	15	7.37	0.04	16.40	45.84
	Other	31	5.40	0.04	16.80	41.71
	Metazoa	71	4.64	0.08	19.50	62.78

**Figure 1** summarizes the distribution of the per-proteome abundance of the C-/CFYWH-depleted proteins for the three domains of life. The numbers of the C-depleted proteins vary significantly between Archaea, Bacteria and Eukaryota ( $p$ -values  $< 0.0001$ ; **Figure 1A**) while the numbers of the CFYWH-depleted proteins are not significantly different ( $p$ -values  $\geq 0.01$ ; **Figure 1B**). Our analysis has revealed that prokaryotes harbor significantly larger numbers of the C-depleted proteins compared to eukaryotes. Furthermore, we compared these rates with the corresponding rates for the proteins with known structures collected from PDB. About 35% of proteins in the PDB dataset are depleted in cysteine (**Figure 1A**), while only two proteins (0.002%) are depleted in cysteine and aromatic residues (**Figure 1B**). The relatively high rate of the C-depleted proteins in PDB can be explained by two observations: about 2/3 of the PDB dataset is composed of the prokaryotic proteins; and because 51% of the proteins in this dataset were solved in complex with nucleic acids (PDB NA dataset). The effect of the latter factor is supported by our empirical finding that about 50% of the proteins in the PDB NA dataset are depleted in cysteine, which is a substantial enrichment particularly when compared to the PDB monomer dataset that has only about 19% of the C-depleted proteins (**Figure 1A**). We also emphasize the lack of the CFYWH-depleted proteins in PDB (**Figure 1B**). We found only two of them overall, with none in the PDB NA dataset and only one among the monomers. Importantly, the levels of the presence of the CFYWH-depleted proteins in PDB are substantially lower when compared to the rates of the CFYWH-depleted proteins in whole proteomes; i.e., 0.002% in PDB *vs.* 0.04% in Eukaryotic and Bacterial proteomes (20-fold decrease) and 0.06% in Archaeal proteomes (30-fold decrease).

Altogether, these results demonstrate that the C-depleted proteins are significantly enriched in prokaryotes compared to eukaryotes and that they are often involved in protein-nucleic acids interactions and relatively rarely fold into monomer structures. On the other hand, the CFYWH-depleted proteins are equally abundant across the three domains of life and virtually never found in PDB. The latter suggest that they are hard to solve structurally.

## **CFYWH- and FYWH-depleted proteins in the PDB dataset and their intrinsic disorder status**

Our search for the CFYWH-depleted proteins in the PDB dataset produced only two hits, a deletion mutant of the transcarboxylase biotin carrier subunit (also known as biotin carboxyl carrier protein, BCCP) from *Propionibacterium freudenreichii* *subsp. shermanii* (PDB ID: 1O78) and a molybdenum-pterin-binding protein 2 (molbindin-2 or MopII) from *Clostridium pasteurianum* (PDB ID: 1GUT). Functionally, BCCP serves as a carrier subunit of the transcarboxylase, which is a biotin-dependent 1200 kDa multi-subunit enzyme composed of 30 separable polypeptides [88]. Here, BCCP functions as a carboxyl group carrier to which biotin is covalently attached at Lys89. BCCP also binds the other two subunits of transcarboxylase to assist in the overall assembly of the enzyme [89]. NMR solution structure analysis revealed that the BCCP C-terminal domain (residues 51–123) is characterized by a compact  $\beta$ -sandwich structure, whereas the N-terminal region of the protein (residues 1–50) is disordered and does not have detectable structure [90]. **Figure 2A** represents the NMR solution structure of a CFYWH-depleted 10-48 deletion mutant (residues 1–9/49–123) of BCCP and shows that this protein contains six anti-parallel  $\beta$ -strands forming  $\beta$ -sandwich and a rather disordered N-terminal region. Furthermore, high flexibility was also detected at the C-terminal ‘ $\beta$ -finger’ segment of this deletion



mutant that contains the Lys89 biotinylation site [91]. The second CFYWH-depleted protein in PDB, molbindin-2, is a bacterial protein that serves as an intracellular storage facility for molybdate. **Figure 2B** shows that this protein exists as a hexamer assembled as a trimer of dimers and binds up to eight molybdate ions with high affinity [92]. A protomer of this protein has a twisted antiparallel  $\beta$ -sheet structure formed by five  $\beta$ -strands [92] (see **Figure 2C**).

Our analysis showed that the number of proteins in the PDB dataset that are depleted in the aromatic residues (FYWH-depleted) is also very low. In fact, we found only 7 such proteins, which, in addition to the aforementioned BCCP and molbindin-2 were a ribosomal protein S33 from *Trypanosoma brucei brucei* (strain 927/4 *GUTat10.1*) (which is a part of the bacterial ribosome, high resolution structure of which was solved by cryo-electron microscopy, PDB ID: 4V8M-AZ, see **Figure 2D**), a bacterial ribosomal protein S28E from *Methanobacterium thermoautotrophicum* (PDB ID: 1NE3, see **Figure 2E**), an eukaryotic 40S ribosomal protein rpS28e from *Tetrahymena thermophila* (PDB ID: 4V5O-A1/B1 and PDB ID: 4BTS-A1/B1/C1/D1, see **Figure 2F**), the pulmonary surfactant-associated polypeptide C (SP-C, PDB ID: 1SPF, see **Figure 2G**), and rat metallothionein-2 (PDB ID: 4MT2, see **Figure 2H**). Three of these FYWH-depleted proteins are ribosomal proteins, with the solution NMR structure of one of which (S28E) being solved at pH 4.5, and with two others (S33 and rpS28e) being a part of the ribosomal subunit). One of them (metallothionein-2) is a metal-binding protein that does not have any regular secondary structure elements and whose 3D structure is stabilized by homodimerization and coordination of five cadmium ions, two zinc ions, and one sodium ion [93]. The last one is a membrane-embedded protein (SP-C), whose structure in apolar solvent (a mixed solvent of C<sub>2</sub>H<sub>3</sub>Cl/C<sub>2</sub>H<sub>3</sub>OH/ 1 M HCl 32:64:5 (v/v)) was solved by NMR [94].

Since all CFYWH- and FYWH-depleted proteins in the PDB dataset are rather small and are characterized by strong amino acid biases (for example, metallothionein-2 possesses extremely high content of cysteine residues (32.8%)), next, we analyzed their intrinsic disorder predispositions using a set of commonly used per-residue disorder predictors, such as PONDR<sup>®</sup> VLXT [32], PONDR<sup>®</sup> VL3 [95], PONDR<sup>®</sup> VSL2 [96], IUPred\_short [97] (yellow curve), IUPred\_long [97], and PONDR<sup>®</sup> FIT [98]. **Figure 3** indicates that many of these proteins are predicted to have high levels of intrinsic disorder. In fact, according to their mean disorder predisposition, they can be ranged as follows: S33 (0.60±0.16) > rpS28e (0.52±0.17) > BCCP (0.47±0.14) = metallothionein-2 (0.47±0.44) > S28E (0.46±0.15) > molbindin (0.33±0.17) > SP-C (0.16±0.15). Low level of intrinsic disorder in SP-C was expected, since this is a transmembrane protein characterized by the high content of hydrophobic, order-promoting residues. **Figure 3** also shows that although, generally, the outputs of the predictors used in this study agree with each other, the disorder profile generated for metallothionein-2 reflects noticeable “confusion”, where PONDR<sup>®</sup> VL3, PONDR<sup>®</sup> VSL2, and PONDR<sup>®</sup> FIT predicted this protein to be completely disordered, whereas IUPred\_short and IUPred\_long suggested that the metallothionein-2 is absolutely ordered. This discrepancy is defined by the highly biased amino acid sequence of this protein, which does not have aromatic residues, being instead heavily enriched in cysteine residues (20 of its 61 residues (32.8%) are cysteines).

One can argue that CWYFH-depleted proteins could contain a higher number of other (non-WYFH) hydrophobic amino acids and still be folded. Unfortunately, the amount of currently available data related to such proteins is not sufficient for conducting reliable statistical analysis to check this hypothesis. In fact, almost complete lack of the non-smelly proteins in PDB, which has only two

CWYFH-depleted and seven WYFH-depleted proteins, serves as an important indication that unique (foldable) protein structure requires cysteines and aromatic residues. Composition profiler-based [34] comparison of the amino acid compositions of two CWYFH-depleted proteins, BCCP and molbindin-2, with the amino acid compositions of globular proteins in PDB revealed that these non-smelly proteins are significantly enriched in valines ( $p\text{-value} < 0.05$ ). Extending this analysis to all seven WYFH-depleted proteins showed that they are significantly enriched in valines and methionines. However, the levels of other hydrophobic residues (leucines and isoleucines) were not significantly increased. These data are insufficient for making unambiguous conclusion on the presence of the compensatory increase in the number of non-CWYFH hydrophobic amino acids in the non-smelly proteins. Generally, since hydrophobic residues are order-promoting [37], one would expect that if such compensation would take place, then the resulting WYFH-depleted proteins with the increased content of non-WYFH hydrophobic residues would still be mostly ordered. However, we are showing here that proteins without CWYFH are more disordered than proteins with CWYFH (see below). This indicates that the proposed compensation is not globally observed. Of course there could be some exceptions from the rule, but there is no such compensation, in general. Furthermore, there is a logical limit on how many hydrophobic residues one can put into a sequence that can fold into a soluble structure (there is the surface to volume ratio limiting the number of hydrophobic groups that can be protected from water by a surface layer of hydrophilic residues upon formation of a globular structure).

In summary, the number of the CFYWH- and FYWH-depleted proteins in PDB is vanishingly small. Despite being structurally characterized, these proteins typically (with the noticeable exception of the pulmonary surfactant-associated polypeptide C, which is a highly hydrophobic, membrane binding protein) have rather content of intrinsic disorder. None of these proteins are enzymes. They are either oligomeric metal-binding proteins, or ribosomal proteins engaged in interaction with ribosomal RNA and other ribosomal proteins, or parts of protein complexes, or transmembrane proteins. In other words, none of these seven protein exist as a non-interacting monomer, suggesting that their structure is stabilized by interaction with binding partners. Therefore, it is safe to conclude that stable monomeric protein structure requires inclusion of cysteine and aromatic residues.

## C- and CFYWH-depleted proteins are enriched in intrinsic disorder

The empirical observation that C-/CFYWH-depleted proteins are relatively rare in PDB suggests that they could be intrinsically disordered [99,8,10]. We tested this hypothesis utilizing accurate putative annotations of disorder. **Figure 4A** compares the putative disorder content (% of disordered residues) in all complete proteomes with the putative disorder content in the C-depleted and the CFYWH-depleted proteins for each domains of life. We found that proteins depleted in cysteine have relatively high disorder content at 7.8% in Archaea, 11.0% in Bacteria and 44.4% in Eukaryota. These are substantially higher amounts when compared to the corresponding complete proteomes. The increases relative to the complete proteomes range between  $(7.8-5.7)/5.7 = 37\%$  in Achaea and  $(44.4-20.2)/20.2 = 120\%$  in Eukaryota. The amounts of the putative intrinsic disorder are event higher among the CFYWH-depleted proteins, with 40.1% disorder content in Archaea, 60.9% in Bacteria, and over 80% in Eukaryota. When compared to the proteome-level disorder content, this corresponds to the relative increases by 604%, 867%, and 312%, respectively. **Figure 4B** compares fractions of the fully disordered proteins between the complete proteomes and the C-depleted and CFYWH-depleted protein datasets. The enrichment in

the number of fully disordered protein is even more substantial than for the disorder content. About 0.2% of all proteins vs. 0.7% of the C-depleted proteins in Archaea are fully disordered (250% increase), 0.4% vs. 1.4% in Bacteria (250% increase), and 1.0% vs. 8.2% in Eukaryota (820% increase). The corresponding increases when comparing the whole proteome-level amounts with the subset of the CFYWH-depleted proteins are approximately 8900% in Archaea and Bacteria and 6700% in Eukaryota. These results clearly demonstrate that the depletion in cysteine and in aromatic residues is associated with substantially elevated levels of intrinsic disorder across all domains of life.

We further analyzed proteome-level relation between the disorder content and the abundance of the C-depleted proteins, see **Figure 5A**. We did not pursue this analysis for the CFYWH-depleted proteins since their numbers are small relative to the proteome sizes (**Figure 1A**), and, therefore, they do not make sufficient impact on the proteome-level measurements. **Figure 5A** reveals a slight increase in the disorder content for organisms with high levels of C-depleted proteins; i.e., the linear fit is sloped upwards and the Pearson correlation coefficients (PCCs) are positive consistently across the three domains of life. This is in agreement with the domain-level increase in the intrinsic disorder for the C-depleted proteins shown in **Figure 4A**. We investigated whether this trends correlates with the current levels of structural coverage (% of proteins with at least partially known structures). **Figure 5B** shows relation between the current structural coverage and the amount of the C-depleted proteins for the three domains of life. We observe modest correlations for Archaea (PCC = -0.34) and Eukaryota (PCC = -0.43), and no correlation for Bacteria (PCC = 0.03). This suggests that archaeal and eukaryotic organisms with higher levels of the C-depleted proteins are characterized by lower levels of structural coverage. **Table 1** reveals that in case of the Eukaryotes this is driven by the high structural coverage and low fraction of the C-depleted proteins in Metazoa. This, in turn, is related to a strong taxonomic bias in the PDB, where 44% of protein structures (61,323 out of the total of 138,194) are from metazoan organisms, in spite of the fact that only 10% of currently sequenced proteins (13,484,303 out of 134,315,728 in UniProt) are from this kingdom of life. One possible explanation for the lack of the correlation in Bacteria is that these proteins have high propensity for crystallization, particularly in contrast to the eukaryotic proteins [76]. This has substantial influence since X-ray crystallography is the single biggest contributor to the protein structure determination efforts [100]; i.e., 90.3% (124,770 out of 138,194) protein structures in PDB were solved using X-ray crystallography. The visible decline in the structural coverage for Archaeal proteins, which also have high propensity for crystallization [76], is likely a result of the significantly higher amount of the C-depleted proteins (**Figure 1A**), when compared to the Bacterial proteins. Overall, our empirical analysis reveals that archaeal and eukaryotic organisms with higher levels of the C-depleted proteins are characterized by higher levels of the intrinsic disorder and lower levels of the structural coverage.

## Functional analysis of C-depleted and CFYWH-depleted proteins

**Figure 6** lists cellular location and functions that are enriched among the C-depleted proteins. The analysis is broken into three types of annotations: cellular components (at the top of the figure), molecular functions (in the middle of the figure) and biological processes (at the bottom of the figure) and performed separately for each domain of life. The C-depleted proteins in Archaea and Bacteria are primarily localized in membranes and ribosome while in Eukaryota they are also found in the nucleosome and nucleolus. These subcellular locations point to a high likelihood that C-depleted

proteins are involved in the protein-RNA and protein-DNA interactions. Molecular functions listed in **Figure 6** reveal that indeed they interact with the rRNAs and, in Eukaryotes, with DNA, while also being involved in the transporter and motor functions. Since the C-depleted proteins are enriched in the intrinsic disorder, these observations are further supported by literature that suggests that IDPs and IDPRs play key roles in the protein-nucleic acids interactions [71,101,69,102,68,103,104,30,105,32,106,107]. The biological processes associated with the C-depleted proteins are consistent with the aforementioned observations, and they cover translation, protein folding, nucleosome assembly, and protein transport. The subset of fully disordered C-depleted proteins (FD lines in **Figure 6**) can be found in ribosome, nucleosome, and, in Eukaryotes, in the chromatin. These proteins implement translation in Archaea and Bacteria, and they also carry out several other functions in Eukaryota, such as response to stress and spermatogenesis. Overall, our analysis reveals that protein-nucleic acids interactions underlie cellular functions and locations of the C-depleted proteins.

**Figure 7** summarizes the major subcellular locations and functions that are enriched in the CFYWH-depleted proteins. These proteins are primarily found in ribosome in Archaea and Bacteria, while in Eukaryota, they are also located in the nucleus, particularly in the chromatin, by being part of the nucleosome complex. This is in agreement with the observations that these proteins are significantly enriched in disorder (**Figure 4**), and that the nucleosome and ribosome complexes contain proteins enriched in disordered regions [71,69,102,68]. The molecular functions and processes associated with the CFYWH-depleted proteins involve RNA and DNA binding in the context of translation, nucleosome assembly, and transcription. This is again consistent with earlier studies that revealed that the high levels of intrinsic disorder represent one of the important characteristics of the nucleic acid binding proteins [71,101,69,102,68,103,104,30,105,32,106,107]. Furthermore, the spatially and temporally coordinated action of many macromolecular complexes and proteins containing functionally significant IDPRs represents an important means for the control of transcription [108]. The major stages of transcription include chromatin remodeling that regulates the global accessibility of promoter DNA, action of regulatory transcription factors, co-activators/co-repressors, and the basal transcription machinery, and at each of these stages, intrinsically disordered proteins or proteins with IDPRs play very important regulatory roles [108]. Next, we discuss the role of disordered nucleic acid-binding proteins in each of these stages.

Formation of the nucleosomes, which are the basic structural units of chromatin, represents the primary step in the DNA condensation that is strongly protein intrinsic disorder-dependent. Nucleosomes are formed via association of small, highly basic nuclear proteins, core histones, with DNA in a specific stoichiometry. The formed nucleosomes are condensed together via action of the linker histones. Comprehensive bioinformatics analyses of 2007 histones from 746 species revealed that all the members of the histone family are highly disordered and utilize disorder for various functions, such as heterodimerization, formation of higher order oligomers, interaction with DNA and other proteins, and posttranslational modifications [102]. Among nuclear proteins that bind to nucleosomes, alter the structure of chromatin, and affect transcription are the members of a high mobility group N (HMGN) protein family of highly disordered chromatin modifying proteins [109]. In addition to HMGNs, many other IDPs and proteins containing functionally important IDPRs, including various chromatin modifying enzymes, are involved in the regulation of DNA accessibility [108].

Among the most illustrative examples of IDPs related to the regulation of transcription (after the chromatin environment becomes accessible due to the action of chromatin remodeling proteins) are transcription factors (TFs, which are also known as sequence-specific DNA-binding factors). TFs are multifunctional proteins which are crucial for the control of expression of specific genes and for the regulation of the gene activity in response to specific stimuli. They deliver their effects via binding to specific DNA sequences, recruiting the RNA polymerase to specific genes, controlling the transfer of genetic information from DNA to mRNA, and positively or negatively influencing the gene transcription either alone or in a complex with other proteins [110]. Generally, the modular structure of TFs includes one or more DNA-binding domains (DBDs) for recognition and binding of the specific DNA sequences adjacent to the genes that they regulate, and one or more transactivation domains for recognition of the co-activators and/or other transcription factors. Computational analysis of several TF datasets revealed that between 82.6% and 94.1% of TFs possess long IDPRs, with the degree of disorder being significantly higher in eukaryotic TFs in comparison with their prokaryotic counterparts [111]. TFs also contain high levels of disorder-based protein interaction sites, molecular recognition features (MoRFs) [111]. Intrinsic disorder is not distributed evenly within the sequences of TFs. In fact, although in general the DNA-binding domains are noticeably less disordered than the TF activation regions (or transactivator domains), the AT-hooks and basic regions of DNA-binding domains of TFs are highly disordered [111]. In human TFs, almost 50% of the entire sequences are occupied by IDPRs [112]. Intrinsic disorder of transactivator domains is used in communication of TFs with other regulatory transcriptional proteins and has an important role in orchestrating the transcriptional assemblies [108]. Based on the high prevalence and versatility of intrinsic disorder in eukaryotic TFs, it has been concluded that these proteins can be used as important illustrations of various aspects of intrinsic disorder-based functionality [113].

At the next stage of transcription, co-activators and co-repressors define a cross-talk between chromatin, transcription factors, and the basal transcription machinery. Some of the co-activators can be considered as scaffolds containing multiple transcription factor binding sites and thereby processing multiple transcriptional regulatory inputs. One of such co-activators, p300, is known to interact with over 50 proteins and possesses histone acetyltransferase activity. Another illustrative example of the importance of intrinsic disorder in the transcription regulation is the Mediator complex. This complex serves as an interface between gene-specific regulatory proteins and the general transcription machinery and it contains high levels of functional intrinsic disorder [114].

Similarly, many proteins related to translation (i.e., the process of ribosome-mediated biosynthesis of proteins from mRNA) are either intrinsically disordered or contain long IDPRs. For example, ribosomal proteins are considered as an important example of the exceptional functional versatility of the RNA-binding IDPs. Based on the comprehensive bioinformatics analyses of the 3,411 ribosomal proteins from 32 species it has been concluded that many ribosomal proteins are either intrinsically disordered as a whole or represent hybrids containing ordered and disordered domains, and that intrinsic disorder is absolutely crucial for their various functions [69]. In agreement with these observations, our analysis showed that three of seven FYWH-depleted proteins whose structure is present in PDB are ribosomal proteins.

Taken together, our analysis revealed that although proteins that do not contain cysteines constitute rather large fraction of the analyzed proteomes (content of such C-depleted proteins is ranging from 8%

proteins in Eukaryota to 19% in Bacteria and to 28% in Archaea), proteins that do not have cysteine and aromatic residues (CFYWH-depleted proteins) constitute only very minor fractions of 817 complete proteomes (about 0.06% proteins in Archaea and 0.04% proteins in Bacteria and Eukaryota). Archaeal and Eukaryotic organisms with higher levels of the C-depleted proteins are predicted to have higher intrinsic disorder levels and lower structural coverage levels, whereas CFYWH-depleted proteins across all domains of life are characterized by the substantially increased levels of intrinsic disorder. Functional analysis revealed that the “non-smelly” proteins are often involved in protein-nucleic acids interactions. They are rarely present as independently folded monomeric structures and often serve as parts of the ribosome and nucleosome complexes. They are also found in cellular membranes. These C- and CFYWH-depleted “non-smelly” proteins are involved in translation, transcription, nucleosome assembly, transmembrane transport, and protein folding functions, all of which are known to be associated with the intrinsic disorder.

Finally, described in this article general inability of the “non-smelly” proteins to fold into self-organizing monomeric structures provides support to the hypothesis on the highly disordered nature of the primordial proteins, which is based on an intriguing correlation between the evolutions of genetic code and protein structure [115-117]. In fact, it was pointed out that the “prebiotic set” of amino acids (i.e., a set of amino acids that were generated by various abiotic processes) likely included 10 of 20 modern amino acids, such as A, D, E, G, I, L, P, S, T, and V [118,119], many of which were disorder-promoting. Based on a combination of 40 different factors, Eduard Trifonov proposed the following temporal order of addition of the amino acids to the genetic code: G/A, V/D, P, S, E/L, T, R, N, K, Q, I, C, H, F, M, Y, W [120]. This sorting underscores the correlation between the appearance of early amino acids (such as G, D, E, P, and S) in the primordial soup and their disorder-promoting tendencies in IDPs. In contrast, it seems that the major order-promoting residues, such as C, W, Y, F, and H, have been added to the genetic code at later evolutionary stages [116,115]. In other words, primordial proteins were “non-smelly”. Similar inferences were also made by Brooks *et al.* in their study on the amino acid composition of last universal ancestral genomes [121]. Also, it was pointed out that the emergence of the biosynthesis of aromatic amino acid enabled an early halophile-to-mesophile transition, emphasizing the potential role of aromatic residues in the adaptive spread of early life and suggesting a selective advantage for the incorporation of aromatic amino acids into the codon table [122]. Furthermore, the high prevalence of nucleic acid binding-related functions among the modern “non-smelly” proteins can be considered as a kind of functional fossil, since nucleic acid binding and RNA chaperoning were proposed to be the first functions of primordial polypeptides [123,124]. Such RNA chaperone activities of early proteins provided their carriers a significant selective advantage in the RNA world, where RNA, which is especially prone to misfolding [125,126], was used for both information storage and catalysis [127].

## Conclusions

We report the results of a comprehensive bioinformatics analysis of the prevalence and functionality of the “non-smelly” proteins (i.e., proteins that do not contain cysteine and aromatic residues, C- and CFYWH-depleted proteins) among 9,563,159 proteins from the 817 complete proteomes, and among the 99,461 PDB proteins with known 3D structures. This analysis revealed that prokaryotes are significantly enriched in the C-depleted proteins compared to eukaryotes. In fact, 28% proteins in

Archaea and 19% in Bacteria vs. only 8% in Eukaryota do not have cysteines. In general, C-depleted proteins are often involved in protein-nucleic acids interactions and they relatively rarely fold into monomer structures. On the other hand, CFYWH-depleted proteins are rather rare, are equally distributed across the three domains of life, and are virtually never found in PDB. Only about 0.05% of proteins do not have cysteine and aromatic residues. Depletion in cysteine and in aromatic residues is associated with the substantially elevated levels of intrinsic disorder in proteins across all domains of life. Archaeal and Eukaryotic organisms with higher levels of the C-depleted proteins have higher levels of the intrinsic disorder and lower levels of structural coverage. The C- and CFYWH-depleted proteins are part of the ribosome and nucleosome complexes and are also found in cellular membranes. They are involved in translation, transcription, nucleosome assembly, transmembrane transport and protein folding functions, all of which are known to be associated with the intrinsic disorder.

In line with highly disordered nature of the “non-smelly” proteins, is an important observation that such proteins are highly underrepresented in PDB. As a matter of fact, there are only two CFYWH-depleted proteins and five FYWH-depleted proteins among the ten thousand proteins in the PDB datasets which were solved by X-ray crystallography, or NMR, or cryo-EM. Furthermore, only two of these proteins, deletion mutant of BCCP and ribosomal protein S28E, have structures that can be considered as a result of a spontaneous folding of a single polypeptide chain, whereas structures of the other “non-smelly” proteins are stabilized by binding of metal ions and self-oligomerization (metallothionein-2 and molbindin-2) or by inclusion into large ribonucleoprotein complexes (ribosomal proteins S33 and rpS28e), or by placing into the non-polar solvent (pulmonary surfactant-associated polypeptide C). These observations indicate that a self-foldable unique 3D-structure in a globular protein is crucially dependent on the presence of cysteine and aromatic residues in its amino acid sequence.

## Acknowledgments

This research was supported in part by the Robert J. Mattauch Endowment funds and the National Science Foundation grant 1617369 to Lukasz Kurgan.

## Conflict of interest

The authors have declared no conflict of interest.

## References

1. Xue B, Williams RW, Oldfield CJ, Dunker AK, Uversky VN (2010) Archaic chaos: intrinsically disordered proteins in Archaea. *BMC Syst Biol* 4 Suppl 1:S1. doi:10.1186/1752-0509-4-S1-S1
2. Peng Z, Yan J, Fan X, Mizianty MJ, Xue B, Wang K, Hu G, Uversky VN, Kurgan L (2015) Exceptionally abundant exceptions: comprehensive characterization of intrinsic disorder in all domains of life. *Cell Mol Life Sci* 72 (1):137-151. doi:10.1007/s00018-014-1661-9

3. Ward JJ, Sodhi JS, McGuffin LJ, Buxton BF, Jones DT (2004) Prediction and functional analysis of native disorder in proteins from the three kingdoms of life. *J Mol Biol* 337 (3):635-645. doi:10.1016/j.jmb.2004.02.002
4. Dunker AK, Obradovic Z, Romero P, Garner EC, Brown CJ (2000) Intrinsic protein disorder in complete genomes. *Genome Inform Ser Workshop Genome Inform* 11:161-171
5. Peng Z, Mizianty MJ, Kurgan L (2014) Genome-scale prediction of proteins with long intrinsically disordered regions. *Proteins* 82 (1):145-158. doi:10.1002/prot.24348
6. Yan J, Mizianty MJ, Filipow PL, Uversky VN, Kurgan L (2013) RAPID: fast and accurate sequence-based prediction of intrinsic disorder content on proteomic scale. *Biochim Biophys Acta* 1834 (8):1671-1680. doi:10.1016/j.bbapap.2013.05.022
7. Dunker AK, Lawson JD, Brown CJ, Williams RM, Romero P, Oh JS, Oldfield CJ, Campen AM, Ratliff CM, Hips KW, Ausio J, Nissen MS, Reeves R, Kang C-H, Kissinger CR, Bailey RW, Griswold MD, Chiu W, Garner EC, Obradovic Z (2001) Intrinsically disordered protein. *J Mol Graph Model* 19 (1):26-59. doi:10.1016/S1093-3263(00)00138-8
8. Hu G, Wang K, Song J, Uversky VN, Kurgan L (2018) Taxonomic Landscape of the Dark Proteomes: Whole-Proteome Scale Interplay Between Structural Darkness, Intrinsic Disorder, and Crystallization Propensity. *Proteomics*:e1800243. doi:10.1002/pmic.201800243
9. Linding R, Jensen LJ, Diella F, Bork P, Gibson TJ, Russell RB (2003) Protein disorder prediction: implications for structural proteomics. *Structure* 11 (11):1453-1459
10. Oldfield CJ, Xue B, Van YY, Ulrich EL, Markley JL, Dunker AK, Uversky VN (2013) Utilization of protein intrinsic disorder knowledge in structural proteomics. *Biochim Biophys Acta* 1834 (2):487-498. doi:10.1016/j.bbapap.2012.12.003
11. Bhowmick A, Brookes DH, Yost SR, Dyson HJ, Forman-Kay JD, Gunter D, Head-Gordon M, Hura GL, Pande VS, Wemmer DE, Wright PE, Head-Gordon T (2016) Finding Our Way in the Dark Proteome. *J Am Chem Soc* 138 (31):9730-9742. doi:10.1021/jacs.6b06543
12. Kruger R (2016) Illuminating the Dark Proteome. *Cell* 166 (5):1074-1077. doi:10.1016/j.cell.2016.08.012
13. Uversky VN, Dunker AK (2010) Understanding protein non-folding. *Biochim Biophys Acta* 1804 (6):1231-1264. doi:10.1016/j.bbapap.2010.01.017
14. Uversky VN, Oldfield CJ, Dunker AK (2008) Intrinsically disordered proteins in human diseases: introducing the D2 concept. *Annu Rev Biophys* 37:215-246. doi:10.1146/annurev.biophys.37.032807.125924
15. Fuxreiter M, Toth-Petroczy A, Kraut DA, Matouschek A, Lim RY, Xue B, Kurgan L, Uversky VN (2014) Disordered proteinaceous machines. *Chem Rev* 114 (13):6806-6843. doi:10.1021/cr4007329
16. Oates ME, Romero P, Ishida T, Ghalwash M, Mizianty MJ, Xue B, Dosztanyi Z, Uversky VN, Obradovic Z, Kurgan L, Dunker AK, Gough J (2013) D(2)P(2): database of disordered protein predictions. *Nucleic Acids Res* 41 (Database issue):D508-516. doi:10.1093/nar/gks1226
17. Romero PR, Zaidi S, Fang YY, Uversky VN, Radivojac P, Oldfield CJ, Cortese MS, Sickmeier M, LeGall T, Obradovic Z, Dunker AK (2006) Alternative splicing in concert with protein intrinsic disorder enables increased functional diversity in multicellular organisms. *Proc Natl Acad Sci U S A* 103 (22):8390-8395. doi:10.1073/pnas.0507916103
18. Jakob U, Kriwacki R, Uversky VN (2014) Conditionally and transiently disordered proteins: awakening cryptic disorder to regulate protein function. *Chem Rev* 114 (13):6779-6805. doi:10.1021/cr400459c
19. Yan J, Dunker AK, Uversky VN, Kurgan L (2016) Molecular recognition features (MoRFs) in three domains of life. *Mol Biosyst* 12 (3):697-710. doi:10.1039/c5mb00640f
20. Patil A, Kinoshita K, Nakamura H (2010) Hub promiscuity in protein-protein interaction networks. *Int J Mol Sci* 11 (4):1930-1943. doi:10.3390/ijms11041930



21. Gsponer J, Babu MM (2009) The rules of disorder or why disorder rules. *Prog Biophys Mol Biol* 99 (2-3):94-103. doi:10.1016/j.pbiomolbio.2009.03.001
22. Haynes C, Oldfield CJ, Ji F, Klitgord N, Cusick ME, Radivojac P, Uversky VN, Vidal M, Iakoucheva LM (2006) Intrinsic disorder is a common feature of hub proteins from four eukaryotic interactomes. *PLoS Comput Biol* 2 (8):e100. doi:10.1371/journal.pcbi.0020100
23. Hu G, Wu Z, Uversky VN, Kurgan L (2017) Functional Analysis of Human Hub Proteins and Their Interactors Involved in the Intrinsic Disorder-Enriched Interactions. *Int J Mol Sci* 18 (12). doi:10.3390/ijms18122761
24. Dunker AK, Cortese MS, Romero P, Iakoucheva LM, Uversky VN (2005) Flexible nets. The roles of intrinsic disorder in protein interaction networks. *FEBS J* 272 (20):5129-5148. doi:10.1111/j.1742-4658.2005.04948.x
25. Barabasi AL, Oltvai ZN (2004) Network biology: understanding the cell's functional organization. *Nat Rev Genet* 5 (2):101-113. doi:10.1038/nrg1272
26. Barabasi AL, Albert R (1999) Emergence of scaling in random networks. *Science* 286 (5439):509-512
27. Vavouri T, Semple JI, Garcia-Verdugo R, Lehner B (2009) Intrinsic protein disorder and interaction promiscuity are widely associated with dosage sensitivity. *Cell* 138 (1):198-208. doi:10.1016/j.cell.2009.04.029
28. Oldfield CJ, Meng J, Yang JY, Yang MQ, Uversky VN, Dunker AK (2008) Flexible nets: disorder and induced fit in the associations of p53 and 14-3-3 with their partners. *BMC Genomics* 9 Suppl 1:S1. doi:10.1186/1471-2164-9-S1-S1
29. Buljan M, Chalancon G, Dunker AK, Bateman A, Balaji S, Fuxreiter M, Babu MM (2013) Alternative splicing of intrinsically disordered regions and rewiring of protein interactions. *Curr Opin Struct Biol* 23 (3):443-450. doi:10.1016/j.sbi.2013.03.006
30. Uversky VN, Gillespie JR, Fink AL (2000) Why are "natively unfolded" proteins unstructured under physiologic conditions? *Proteins* 41 (3):415-427. doi:10.1002/1097-0134(20001115)41:3<415::AID-PROT130>3.0.CO;2-7 [pii]
31. Dunker AK, Garner E, Guillot S, Romero P, Albrecht K, Hart J, Obradovic Z, Kissinger C, Villafranca JE (1998) Protein disorder and the evolution of molecular recognition: theory, predictions and observations. *Pac Symp Biocomput*:473-484
32. Dunker AK, Lawson JD, Brown CJ, Williams RM, Romero P, Oh JS, Oldfield CJ, Campen AM, Ratliff CM, Hips KW, Ausio J, Nissen MS, Reeves R, Kang C, Kissinger CR, Bailey RW, Griswold MD, Chiu W, Garner EC, Obradovic Z (2001) Intrinsically disordered protein. *J Mol Graph Model* 19 (1):26-59. doi:10.1016/S1093-3263(00)00138-8
33. Radivojac P, Iakoucheva LM, Oldfield CJ, Obradovic Z, Uversky VN, Dunker AK (2007) Intrinsic disorder and functional proteomics. *Biophys J* 92 (5):1439-1456. doi:S0006-3495(07)70955-4 [pii] 10.1529/biophysj.106.094045
34. Vacic V, Uversky VN, Dunker AK, Lonardi S (2007) Composition Profiler: a tool for discovery and visualization of amino acid composition differences. *BMC Bioinformatics* 8:211. doi:1471-2105-8-211 [pii] 10.1186/1471-2105-8-211
35. Romero P, Obradovic Z, Li X, Garner EC, Brown CJ, Dunker AK (2001) Sequence complexity of disordered protein. *Proteins* 42 (1):38-48
36. Campen A, Williams RM, Brown CJ, Meng J, Uversky VN, Dunker AK (2008) TOP-IDP-scale: a new amino acid scale measuring propensity for intrinsic disorder. *Protein Pept Lett* 15 (9):956-963
37. Williams RM, Obradovi Z, Mathura V, Braun W, Garner EC, Young J, Takayama S, Brown CJ, Dunker AK (2001) The protein non-folding problem: amino acid determinants of intrinsic order and disorder. *Pac Symp Biocomput*:89-100

38. Daly NL, Craik DJ (2011) Bioactive cystine knot proteins. *Curr Opin Chem Biol* 15 (3):362-368. doi:10.1016/j.cbpa.2011.02.008
39. Craik DJ, Daly NL, Waine C (2001) The cystine knot motif in toxins and implications for drug design. *Toxicol* 39 (1):43-60
40. Trivedi MV, Laurence JS, Siahaan TJ (2009) The role of thiols and disulfides on protein stability. *Curr Protein Pept Sci* 10 (6):614-625
41. Hagihara Y, Saerens D (2014) Engineering disulfide bonds within an antibody. *Biochim Biophys Acta* 1844 (11):2016-2023. doi:10.1016/j.bbapap.2014.07.005
42. Bechtel TJ, Weerapana E (2017) From structure to redox: The diverse functional roles of disulfides and implications in disease. *Proteomics* 17 (6). doi:10.1002/pmic.201600391
43. Darling AL, Uversky VN (2018) Intrinsic Disorder and Posttranslational Modifications: The Darker Side of the Biological Dark Matter. *Front Genet* 9:158. doi:10.3389/fgene.2018.00158
44. Pace NJ, Weerapana E (2014) Zinc-binding cysteines: diverse functions and structural motifs. *Biomolecules* 4 (2):419-434. doi:10.3390/biom4020419
45. Krishna SS, Majumdar I, Grishin NV (2003) Structural classification of zinc fingers: survey and summary. *Nucleic Acids Res* 31 (2):532-550. doi:10.1093/nar/gkg161
46. Negi S, Itazu M, Imanishi M, Nomura A, Sugiura Y (2004) Creation and characteristics of unnatural CysHis(3)-type zinc finger protein. *Biochem Biophys Res Commun* 325 (2):421-425. doi:10.1016/j.bbrc.2004.10.045
47. Harding MM (2004) The architecture of metal coordination groups in proteins. *Acta Crystallogr D Biol Crystallogr* 60 (Pt 5):849-859. doi:10.1107/S0907444904004081
48. Laska M (2010) Olfactory perception of 6 amino acids by human subjects. *Chem Senses* 35 (4):279-287. doi:10.1093/chemse/bjq017
49. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE (2000) The Protein Data Bank. *Nucleic Acids Res* 28 (1):235-242
50. Burley SK, Berman HM, Kleywegt GJ, Markley JL, Nakamura H, Velankar S (2017) Protein Data Bank (PDB): The Single Global Macromolecular Structure Archive. *Methods Mol Biol* 1607:627-641. doi:10.1007/978-1-4939-7000-1\_26
51. The UniProt C (2017) UniProt: the universal protein knowledgebase. *Nucleic Acids Res* 45 (D1):D158-D169. doi:10.1093/nar/gkw1099
52. UniProt C (2015) UniProt: a hub for protein information. *Nucleic Acids Res* 43 (Database issue):D204-212. doi:10.1093/nar/gku989
53. Hu G, Wu Z, Oldfield CJ, Wang C, Kurgan L (2019) Quality assessment for the putative intrinsic disorder in proteins. *Bioinformatics* 35 (10):1692-1700. doi:10.1093/bioinformatics/bty881
54. Katuwawala A, Oldfield CJ, Kurgan L (2019) Accuracy of protein-level disorder predictions. *Brief Bioinformatics*
55. Walsh I, Giollo M, Di Domenico T, Ferrari C, Zimmermann O, Tosatto SC (2015) Comprehensive large-scale assessment of intrinsic protein disorder. *Bioinformatics* 31 (2):201-208. doi:10.1093/bioinformatics/btu625
56. Monastyrskyy B, Kryshchafovich A, Moulton J, Tramontano A, Fidelis K (2014) Assessment of protein disorder region predictions in CASP10. *Proteins* 82 Suppl 2:127-137. doi:10.1002/prot.24391
57. Peng ZL, Kurgan L (2012) Comprehensive comparative assessment of in-silico predictors of disordered regions. *Curr Protein Pept Sci* 13 (1):6-18
58. Meng F, Uversky VN, Kurgan L (2017) Comprehensive review of methods for prediction of intrinsic disorder and its molecular functions. *Cell Mol Life Sci* 74 (17):3069-3090. doi:10.1007/s00018-017-2555-4
59. Meng F, Uversky V, Kurgan L (2017) Computational Prediction of Intrinsic Disorder in Proteins. *Curr Protoc*

60. Peng Z, Kurgan L (2012) On the complementarity of the consensus-based disorder prediction. *Pac Symp Biocomput*:176-187
61. Fan X, Kurgan L (2014) Accurate prediction of disorder in protein chains with a comprehensive and empirically designed consensus. *J Biomol Struct Dyn* 32 (3):448-464. doi:10.1080/07391102.2013.775969
62. Necci M, Piovesan D, Dosztanyi Z, Tosatto SCE (2017) MobiDB-lite: fast and highly specific consensus prediction of intrinsic disorder in proteins. *Bioinformatics* 33 (9):1402-1404. doi:10.1093/bioinformatics/btx015
63. Dosztanyi Z, Csizmok V, Tompa P, Simon I (2005) The pairwise energy content estimated from amino acid composition discriminates between folded and intrinsically unstructured proteins. *J Mol Biol* 347 (4):827-839. doi:10.1016/j.jmb.2005.01.071
64. Walsh I, Martin AJ, Di Domenico T, Tosatto SC (2012) ESpritz: accurate and fast prediction of protein disorder. *Bioinformatics* 28 (4):503-509. doi:10.1093/bioinformatics/btr682
65. Piovesan D, Tabaro F, Micetic I, Necci M, Quaglia F, Oldfield CJ, Aspromonte MC, Davey NE, Davidovic R, Dosztanyi Z, Elofsson A, Gasparini A, Hatos A, Kajava AV, Kalmar L, Leonardi E, Lazar T, Macedo-Ribeiro S, Macossay-Castillo M, Meszaros A, Minervini G, Murvai N, Pujols J, Roche DB, Salladini E, Schad E, Schramm A, Szabo B, Tantos A, Tonello F, Tsigos KD, Veljkovic N, Ventura S, Vranken W, Warholm P, Uversky VN, Dunker AK, Longhi S, Tompa P, Tosatto SC (2016) DisProt 7.0: a major update of the database of disordered proteins. *Nucleic Acids Res* D1:D219-D227. doi:10.1093/nar/gkw1056
66. Vucetic S, Obradovic Z, Vacic V, Radivojac P, Peng K, Iakoucheva LM, Cortese MS, Lawson JD, Brown CJ, Sikes JG, Newton CD, Dunker AK (2005) DisProt: a database of protein disorder. *Bioinformatics* 21 (1):137-140. doi:10.1093/bioinformatics/bth476
67. Na I, Meng F, Kurgan L, Uversky VN (2016) Autophagy-related intrinsically disordered proteins in intranuclear compartments. *Mol Biosyst* 12 (9):2798-2817. doi:10.1039/c6mb00069j
68. Meng F, Na I, Kurgan L, Uversky VN (2016) Compartmentalization and Functionality of Nuclear Disorder: Intrinsic Disorder and Protein-Protein Interactions in Intra-Nuclear Compartments. *Int J Mol Sci* 17 (1). doi:10.3390/ijms17010024
69. Peng Z, Oldfield CJ, Xue B, Mizianty MJ, Dunker AK, Kurgan L, Uversky VN (2014) A creature with a hundred waggly tails: intrinsically disordered proteins in the ribosome. *Cell Mol Life Sci* 71 (8):1477-1504. doi:10.1007/s00018-013-1446-6
70. Hu G, Wu Z, Wang K, Uversky VN, Kurgan L (2016) Untapped Potential of Disordered Proteins in Current Druggable Human Proteome. *Curr Drug Targets* 17 (10):1198-1205
71. Wang C, Uversky VN, Kurgan L (2016) Disordered nucleome: Abundance of intrinsic disorder in the DNA- and RNA-binding proteins in 1121 species from Eukaryota, Bacteria and Archaea. *Proteomics* 16 (10):1486-1498. doi:10.1002/pmic.201500177
72. Peng Z, Uversky VN, Kurgan L (2016) Genes encoding intrinsic disorder in Eukaryota have high GC content. *Intrinsically Disordered Proteins* 4 (1):e1262225. doi:10.1080/21690707.2016.1262225
73. Di Domenico T, Walsh I, Martin AJM, Tosatto SCE (2012) MobiDB: a comprehensive database of intrinsic protein disorder annotations. *Bioinformatics* 28 (15):2080-2081. doi:10.1093/bioinformatics/bts327
74. Potenza E, Di Domenico T, Walsh I, Tosatto SC (2015) MobiDB 2.0: an improved database of intrinsically disordered and mobile proteins. *Nucleic Acids Res* 43 (Database issue):D315-320. doi:10.1093/nar/gku982
75. Vitkup D, Melamud E, Moulton J, Sander C (2001) Completeness in structural genomics. *Nat Struct Biol* 8 (6):559-566. doi:10.1038/88640
76. Mizianty MJ, Fan X, Yan J, Chalmers E, Woloschuk C, Joachimiak A, Kurgan L (2014) Covering complete proteomes with X-ray structures: a current snapshot. *Acta Crystallogr D Biol Crystallogr* 70 (Pt 11):2781-

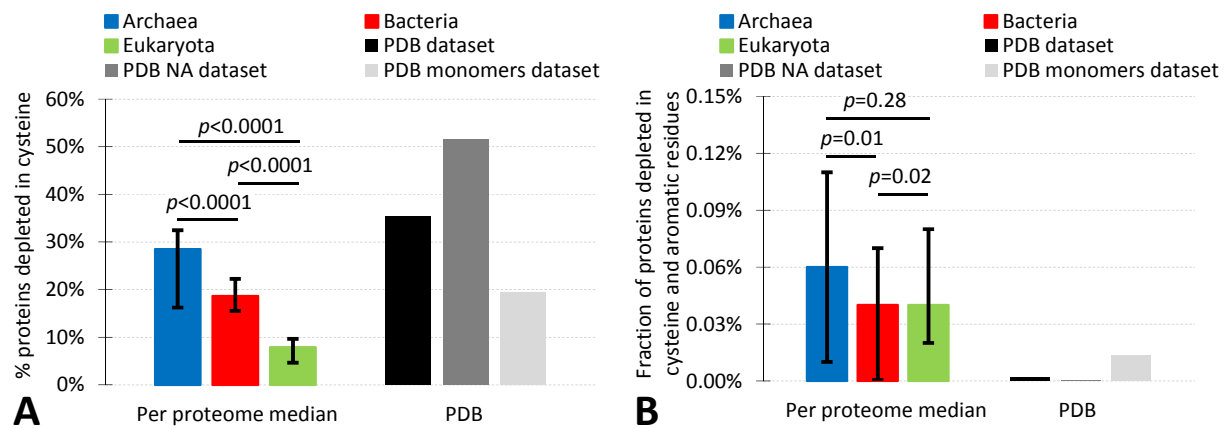
2793. doi:10.1107/S1399004714019427

77. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 25 (17):3389-3402
78. Mariani V, Kiefer F, Schmidt T, Haas J, Schwede T (2011) Assessment of template based protein structure predictions in CASP9. *Proteins* 79 Suppl 10:37-58. doi:10.1002/prot.23177
79. Yang J, Yan R, Roy A, Xu D, Poisson J, Zhang Y (2015) The I-TASSER Suite: protein structure and function prediction. *Nat Meth* 12 (1):7-8. doi:10.1038/nmeth.3213  
<http://www.nature.com/nmeth/journal/v12/n1/abs/nmeth.3213.html#supplementary-information>
80. Yang J, Zhang Y (2015) I-TASSER server: new development for protein structure and function predictions. *Nucleic Acids Res* 43 (W1):W174-181. doi:10.1093/nar/gkv342
81. Soding J, Biegert A, Lupas AN (2005) The HHpred interactive server for protein homology detection and structure prediction. *Nucleic Acids Res* 33 (Web Server issue):W244-248. doi:10.1093/nar/gki408
82. Hildebrand A, Remmert M, Biegert A, Soding J (2009) Fast and accurate automatic structure prediction with HHpred. *Proteins* 77 Suppl 9:128-132. doi:10.1002/prot.22499
83. Webb B, Sali A (2017) Protein Structure Modeling with MODELLER. *Methods Mol Biol* 1654:39-54. doi:10.1007/978-1-4939-7231-9\_4
84. Sali A, Potterton L, Yuan F, van Vlijmen H, Karplus M (1995) Evaluation of comparative protein modeling by MODELLER. *Proteins* 23 (3):318-326. doi:10.1002/prot.340230306
85. Pieper U, Webb BM, Dong GQ, Schneidman-Duhovny D, Fan H, Kim SJ, Khuri N, Spill YG, Weinkam P, Hammel M, Tainer JA, Nilges M, Sali A (2014) ModBase, a database of annotated comparative protein structure models and associated resources. *Nucleic Acids Res* 42 (Database issue):D336-346. doi:10.1093/nar/gkt1144
86. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT, Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis S, Matese JC, Richardson JE, Ringwald M, Rubin GM, Sherlock G (2000) Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet* 25 (1):25-29. doi:10.1038/75556
87. Howell M, Green R, Killeen A, Wedderburn L, Picascio V, Rabionet A, Peng ZL, Larina M, Xue B, Kurgan L, Uversky VN (2012) Not That Rigid Midgets and Not So Flexible Giants: On the Abundance and Roles of Intrinsic Disorder in Short and Long Proteins. *J Biol Syst* 20 (4):471-511. doi:10.1142/S0218339012400086
88. Hennessey JP, Jr., Johnson WC, Jr., Bahler C, Wood HG (1982) Subunit interactions of transcarboxylase as studied by circular dichroism. *Biochemistry* 21 (4):642-646
89. Shenoy BC, Wood HG (1988) Purification and properties of the synthetase catalyzing the biotinylation of the apoprotein of transcarboxylase from *Propionibacterium shermanii*. *FASEB J* 2 (8):2396-2401
90. Reddy DV, Shenoy BC, Carey PR, Sonnichsen FD (2000) High resolution solution structure of the 1.3S subunit of transcarboxylase from *Propionibacterium shermanii*. *Biochemistry* 39 (10):2509-2516
91. Jank MM, Sadowsky JD, Peikert C, Berger S (2002) NMR studies on the solution structure of a deletion mutant of the transcarboxylase biotin carrier subunit. *Int J Biol Macromol* 30 (5):233-242
92. Schuttelkopf AW, Harrison JA, Boxer DH, Hunter WN (2002) Passive acquisition of ligand by the MopII molbindin from *Clostridium pasteurianum*: structures of apo and oxyanion-bound forms. *J Biol Chem* 277 (17):15013-15020. doi:10.1074/jbc.M201005200
93. Braun W, Vasak M, Robbins AH, Stout CD, Wagner G, Kagi JH, Wuthrich K (1992) Comparison of the NMR solution structure and the x-ray crystal structure of rat metallothionein-2. *Proc Natl Acad Sci U S A* 89 (21):10124-10128
94. Johansson J, Szyperski T, Curstedt T, Wuthrich K (1994) The NMR structure of the pulmonary surfactant-associated polypeptide SP-C in an apolar solvent contains a valyl-rich alpha-helix. *Biochemistry* 33

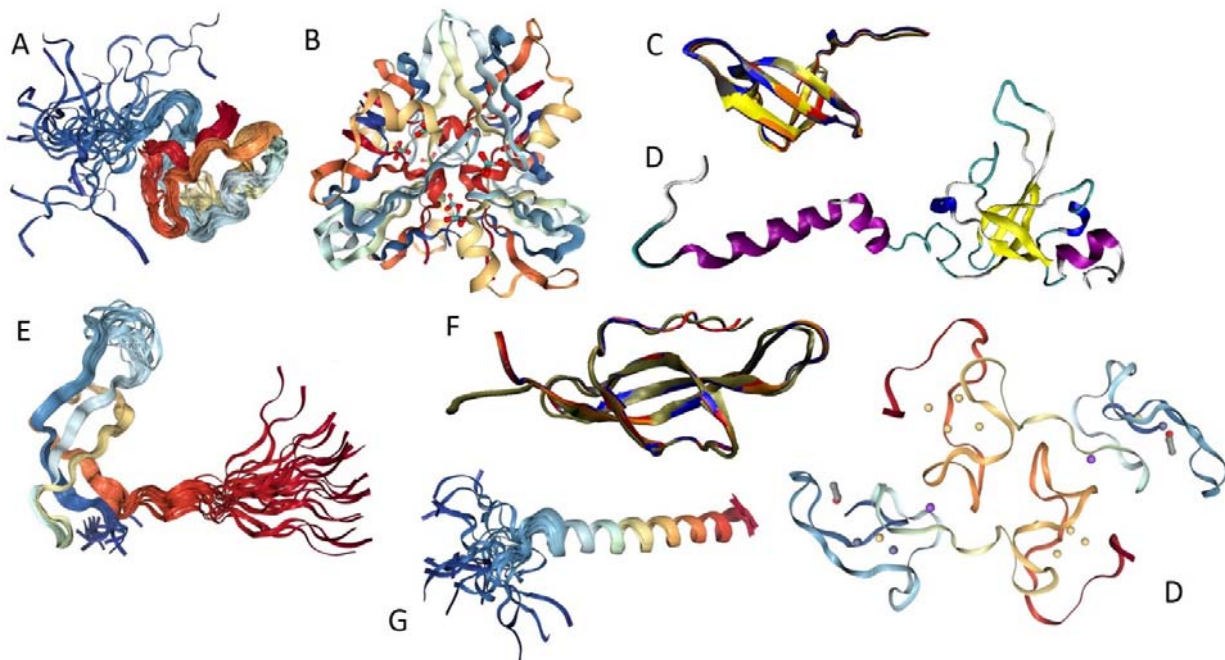
(19):6015-6023

95. Peng K, Radivojac P, Vucetic S, Dunker AK, Obradovic Z (2006) Length-dependent prediction of protein intrinsic disorder. *BMC Bioinformatics* 7:208. doi:10.1186/1471-2105-7-208
96. Peng K, Vucetic S, Radivojac P, Brown CJ, Dunker AK, Obradovic Z (2005) Optimizing long intrinsic disorder predictors with protein evolutionary information. *J Bioinform Comput Biol* 3 (1):35-60. doi:10.1142/s0219720005000886
97. Dosztanyi Z, Csizmek V, Tompa P, Simon I (2005) IUPred: web server for the prediction of intrinsically unstructured regions of proteins based on estimated energy content. *Bioinformatics* 21 (16):3433-3434. doi:10.1093/bioinformatics/bti541
98. Xue B, Dunbrack RL, Williams RW, Dunker AK, Uversky VN (2010) PONDR-FIT: a meta-predictor of intrinsically disordered amino acids. *Biochim Biophys Acta* 1804 (4):996-1010. doi:10.1016/j.bbapap.2010.01.011
99. Oldfield CJ, Ulrich EL, Cheng Y, Dunker AK, Markley JL (2005) Addressing the intrinsic disorder bottleneck in structural proteomics. *Proteins* 59 (3):444-453. doi:10.1002/prot.20446
100. Grabowski M, Niedzialkowska E, Zimmerman MD, Minor W (2016) The impact of structural genomics: the first quinquennial. *J Struct Funct Genomics* 17 (1):1-16. doi:10.1007/s10969-016-9201-5
101. Basu S, Bahadur RP (2016) A structural perspective of RNA recognition by intrinsically disordered proteins. *Cell Mol Life Sci* 73 (21):4075-4084. doi:10.1007/s00018-016-2283-1
102. Peng Z, Mizianty MJ, Xue B, Kurgan L, Uversky VN (2012) More than just tails: intrinsic disorder in histone proteins. *Mol Biosyst* 8 (7):1886-1901. doi:10.1039/c2mb25102g
103. Varadi M, Zsolyomi F, Guharoy M, Tompa P (2015) Functional Advantages of Conserved Intrinsic Disorder in RNA-Binding Proteins. *PLoS One* 10 (10):e0139731. doi:10.1371/journal.pone.0139731
104. Tompa P (2002) Intrinsically unstructured proteins. *Trends Biochem Sci* 27 (10):527-533
105. Dunker AK, Brown CJ, Lawson JD, Iakoucheva LM, Obradovic Z (2002) Intrinsic disorder and protein function. *Biochemistry* 41 (21):6573-6582
106. Chowdhury S, Zhang J, Kurgan L (2018) In Silico Prediction and Validation of Novel RNA Binding Proteins and Residues in the Human Proteome. *Proteomics*:e1800064. doi:10.1002/pmic.201800064
107. Wu Z, Hu G, Yang J, Peng Z, Uversky VN, Kurgan L (2015) In various protein complexes, disordered protomers have large per-residue surface areas and area of protein-, DNA- and RNA-binding interfaces. *FEBS Lett* 589 (19 Pt A):2561-2569. doi:10.1016/j.febslet.2015.08.014
108. Fuxreiter M, Tompa P, Simon I, Uversky VN, Hansen JC, Asturias FJ (2008) Malleable machines take shape in eukaryotic transcriptional regulation. *Nat Chem Biol* 4 (12):728-737. doi:10.1038/nchembio.127
109. Rochman M, Taher L, Kurahashi T, Cherukuri S, Uversky VN, Landsman D, Ovcharenko I, Bustin M (2011) Effects of HMGN variants on the cellular transcription profile. *Nucleic Acids Res* 39 (10):4076-4087. doi:10.1093/nar/gkq1343
110. Latchman DS (1997) Transcription factors: an overview. *Int J Biochem Cell Biol* 29 (12):1305-1312. doi:S1357-2725(97)00085-X [pii]
111. Liu J, Perumal NB, Oldfield CJ, Su EW, Uversky VN, Dunker AK (2006) Intrinsic disorder in transcription factors. *Biochemistry* 45 (22):6873-6888. doi:10.1021/bi0602718
112. Minezaki Y, Homma K, Kinjo AR, Nishikawa K (2006) Human transcription factors contain a high fraction of intrinsically disordered regions essential for transcriptional regulation. *J Mol Biol* 359 (4):1137-1149. doi:S0022-2836(06)00466-9 [pii]
- 10.1016/j.jmb.2006.04.016
113. Staby L, O'Shea C, Willemoes M, Theisen F, Kragelund BB, Skriver K (2017) Eukaryotic transcription factors: paradigms of protein intrinsic disorder. *Biochem J* 474 (15):2509-2532. doi:10.1042/BCJ20160631

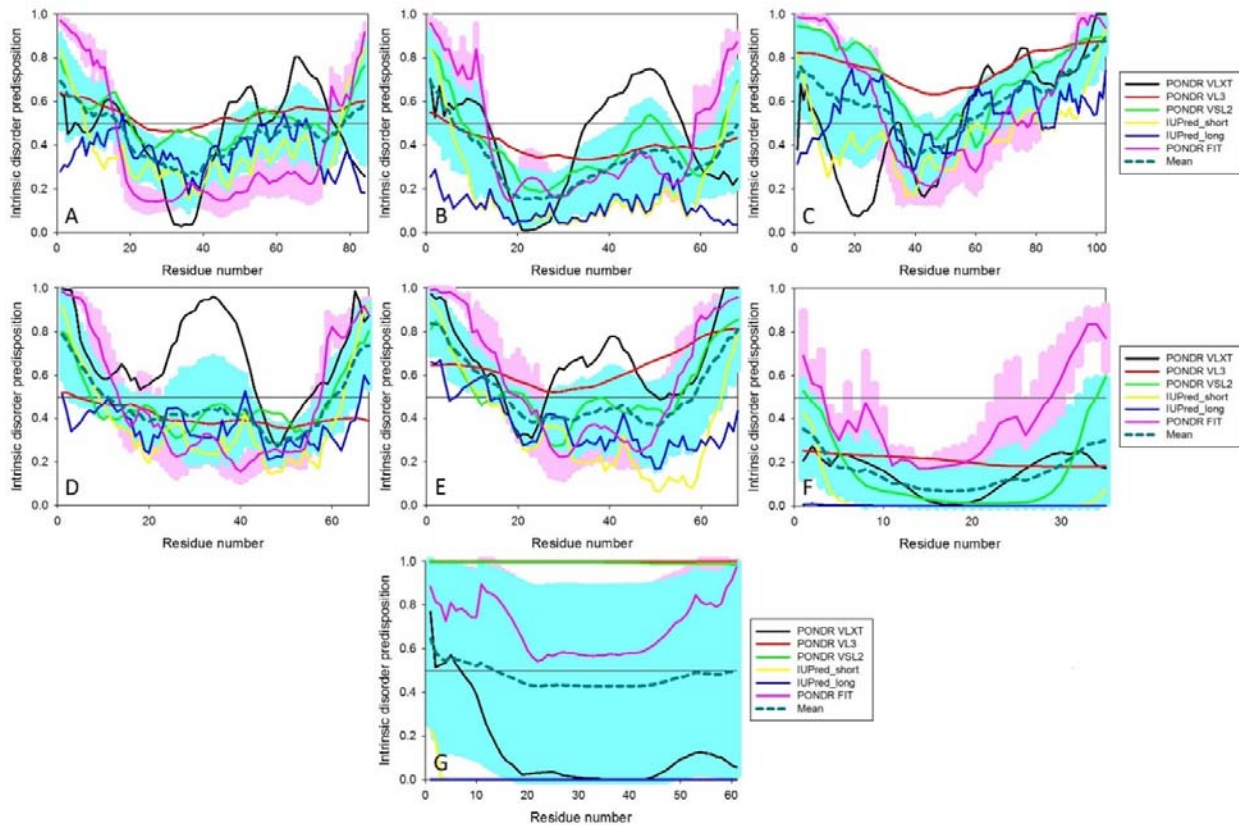
114. Toth-Petroczy A, Oldfield CJ, Simon I, Takagi Y, Dunker AK, Uversky VN, Fuxreiter M (2008) Malleable machines in transcription regulation: the mediator complex. *PLoS Comput Biol* 4 (12):e1000243. doi:10.1371/journal.pcbi.1000243
115. Di Mauro E, Dunker AK, Trifonov EN (2012) Disorder to order, non-life to life: In the beginning there was a mistake. In: Seckbach J (ed) *Genesis - In the beginning. Precursors of Life, Chemical Models and Early Biological Evolution*. Springer, Dordrecht, Heidelberg, New York, London,
116. Uversky VN (2013) A decade and a half of protein intrinsic disorder: biology still waits for physics. *Protein Sci* 22 (6):693-724. doi:10.1002/pro.2261
117. Kulkarni P, Uversky VN (2018) Intrinsically Disordered Proteins: The Dark Horse of the Dark Proteome. *Proteomics* 18 (21-22):e1800061. doi:10.1002/pmic.201800061
118. Longo LM, Blaber M (2012) Protein design at the interface of the pre-biotic and biotic worlds. *Arch Biochem Biophys* 526 (1):16-21. doi:10.1016/j.abb.2012.06.009
119. Longo LM, Blaber M (2014) Prebiotic protein design supports a halophile origin of foldable proteins. *Front Microbiol* 4. doi:ARTN 418  
10.3389/fmicb.2013.00418
120. Trifonov EN (2000) Consensus temporal order of amino acids and evolution of the triplet code. *Gene* 261 (1):139-151
121. Brooks DJ, Fresco JR, Lesk AM, Singh M (2002) Evolution of amino acid frequencies in proteins over deep time: inferred order of introduction of amino acids into the genetic code. *Mol Biol Evol* 19 (10):1645-1655. doi:10.1093/oxfordjournals.molbev.a003988
122. Longo LM, Tenorio CA, Kumru OS, Middaugh CR, Blaber M (2015) A single aromatic core mutation converts a designed "primitive" protein from halophile to mesophile folding. *Protein Sci* 24 (1):27-37. doi:10.1002/pro.2580
123. Poole AM, Jeffares DC, Penny D (1998) The path from the RNA world. *J Mol Evol* 46 (1):1-17
124. Tompa P, Csermely P (2004) The role of structural disorder in the function of RNA and protein chaperones. *Faseb J* 18 (11):1169-1175. doi:10.1096/fj.04-1584rev  
18/11/1169 [pii]
125. Treiber DK, Williamson JR (2001) Beyond kinetic traps in RNA folding. *Curr Opin Struct Biol* 11 (3):309-314. doi:S0959-440X(00)00206-2 [pii]
126. Cristofari G, Darlix JL (2002) The ubiquitous nature of RNA chaperone proteins. *Prog Nucleic Acid Res Mol Biol* 72:223-268
127. Gilbert W (1986) Origin of life - the RNA world. *Nature* 319 (6055):618-618. doi:Doi 10.1038/319618a0
128. Shatsky M, Nussinov R, Wolfson HJ (2004) A method for simultaneous alignment of multiple protein structures. *Proteins* 56 (1):143-156. doi:10.1002/prot.10628
129. Humphrey W, Dalke A, Schulten K (1996) VMD: visual molecular dynamics. *J Mol Graph* 14 (1):33-38, 27-38



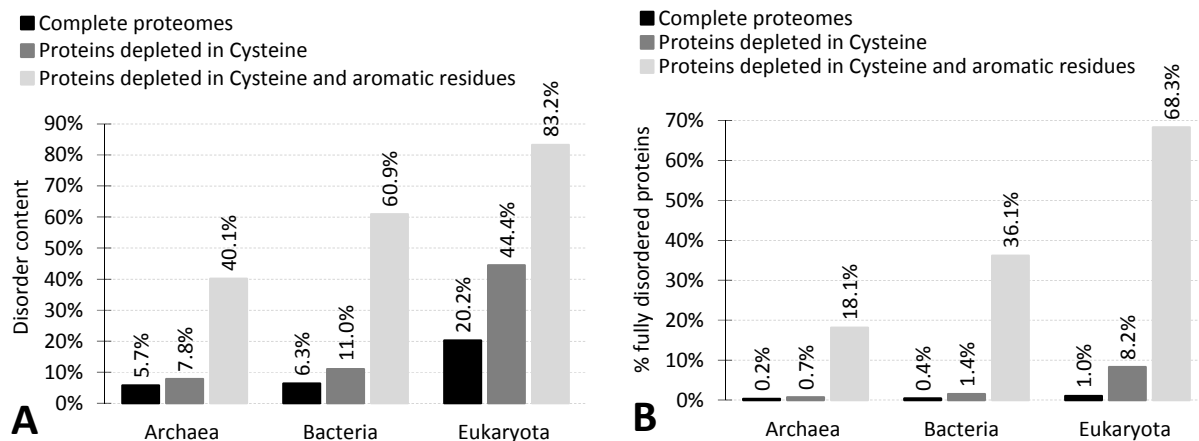
**Figure 1.** Abundance of C-depleted (panel A) and CFYW-H-depleted proteins (panel B) in the three domains of life and among the structurally solved proteins from PDB. The blue, red and green bars show the median per-proteome fraction of C-/CFYW-H-depleted proteins among the 64 Archaeal, 552 Bacterial and 201 Eukaryotic organisms, respectively. The whiskers denote the first and third quartiles of these per-proteome fractions. Statistical significance of the differences for the per-proteome values between domains of life was assessed with the Wilcoxon test for unpaired data; distributions of the measured values are not normal. The black, dark gray and light gray bars show the fraction of the C-/CFYW-H-depleted proteins among wild-type proteins chains from PDB (PDB dataset), wild-type PDB proteins that interact with nucleic acids (PDB NA dataset) and wild-type PDB monomers (PDB monomers dataset), respectively.



**Figure 2.** Structural characterization of the CFYW-H- and FYWH-depleted proteins found in PDB: **A.** NMR solution structure of a CFYW-H-depleted 10-48 deletion mutant (residues 1–9/49–123) of the transcarboxylase biotin carrier subunit (also known as biotin carboxyl carrier protein, BCCP) from *Propionibacterium freudenreichii* subsp. *shermanii* (PDB ID: 1O78); **B.** Crystal structure of the homohexameric molybdenum-pterin-binding protein 2 (molbindin-2 or MopI) from *Clostridium pasteurianum* (PDB ID: 1GUT); **C.** Aligned structures of the molbindin-2 protomers. Structures were aligned using a MultiProt server [128]. Plot was created using VMD platform [129]; **D.** High-resolution cryo-electron microscopy structure of a 40S ribosomal protein S33 from *Trypanosoma brucei brucei* (strain 927/4 GUTat10.1). Structure of this protein was extracted from of the cryo-EM structure of bacterial ribosome (PDB ID: 4V8M-AZ). Plot was created using VMD platform [129]; **E.** Solution NMR structure of a 30S bacterial ribosomal protein S28E from *Methanobacterium thermoautotrophicum* (PDB ID: 1NE3); **F.** Aligned structures of a eukaryotic 40S ribosomal protein rpS28e from *Tetrahymena thermophila* (PDB ID: 4V5O-A1/B1 and PDB ID: 4BTS-A1/B1/C1/D1). Corresponding structures were extracted from the crystal structures of the eukaryotic 40S ribosomal subunit in complex with initiation factor 1 (PDB ID: 4V5O-A1/B1) and the crystal structure of the eukaryotic 40S ribosomal subunit in complex with eIF1 and eIF1A (PDB ID: 4BTS-A1/B1/C1/D1). Structures were aligned using a MultiProt server [128]. Plot was created using VMD platform [129]; **G.** Solution NMR structure of the pulmonary surfactant-associated polypeptide C (SP-C) solved in apolar solvent (a mixed solvent of C<sub>2</sub>H<sub>5</sub>Cl/C<sub>2</sub>H<sub>5</sub>OH/ 1 M HCl 32:64:5 (v/v)) (PDB ID: 1SPF); and **H.** Crystal structure of the metal-bound dimer rat metallothionein-2 (PDB ID: 4MT2).

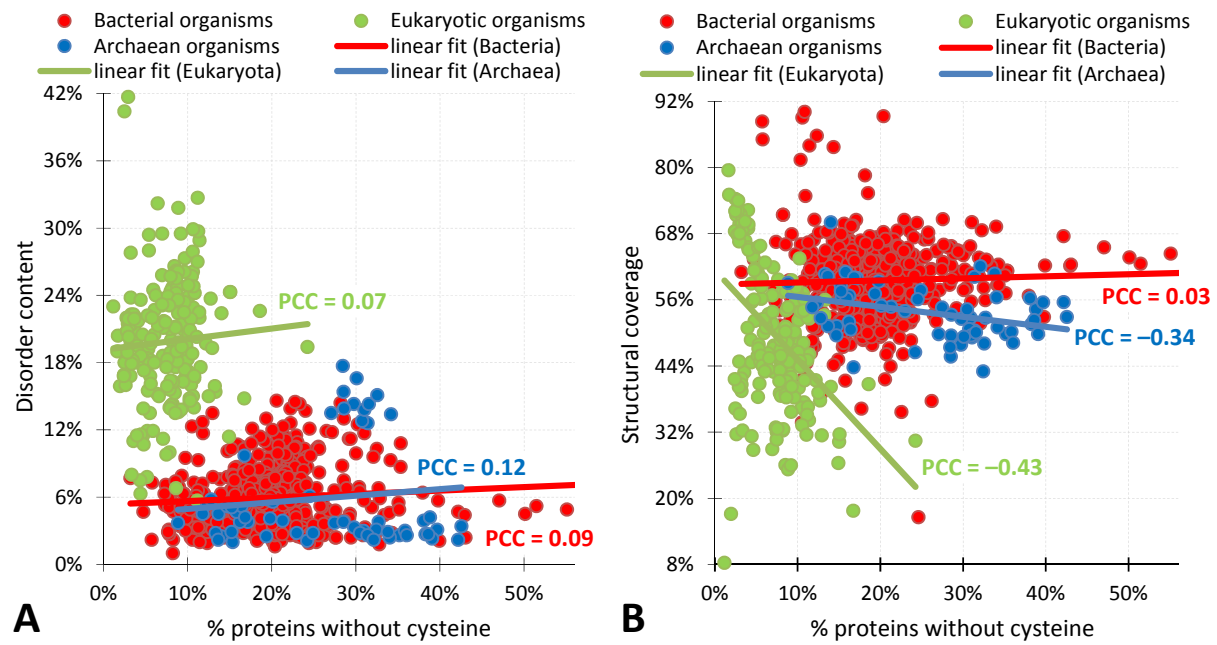


**Figure 3.** Multiparametric analysis of the intrinsic disorder predisposition of the CFYWH- and FYWH-depleted proteins found in PDB by several common predictors of intrinsic disorder: PONDRL VLXT [32] (black curves), PONDRL VL3 [95] (red curves), PONDRL VSL2 [96] (green curves), IUPred\_short [97] (yellow curves), IUPred\_long [97] (blue curves), and PONDRL FIT [98] (pink curves). Dark cyan dashed line shows the mean disorder propensity calculated by averaging disorder profiles of individual predictors. Light pink shadow around the PONDRL FIT shows error distribution for this predictor, whereas light cyan shadow around the mean disorder curve shows error distribution for evaluation of mean disorder. In these analyses, the predicted intrinsic disorder scores above 0.5 are considered to correspond to the disordered residues/regions, whereas regions with the disorder scores between 0.2 and 0.5 are considered flexible. Analyzed proteins were **A.** BCCP (residues 1–9/49–123 of UniProt ID: Q2904); **B.** Molbindin-2 (UniProt ID: P08854); **C.** 40S ribosomal protein S33 from *Trypanosoma brucei* (UniProt ID: Q57U30); **D.** 30S ribosomal protein S28E from *Methanobacterium thermoautotrophicum* (UniProt ID: O26356); **E.** 40S ribosomal protein rpS28e from *Tetrahymena thermophila* (UniProt ID: Q234G5); **F.** pulmonary surfactant-associated polypeptide C (UniProt ID: P15785); and **G.** metallothionein-2 (UniProt ID: P04355).

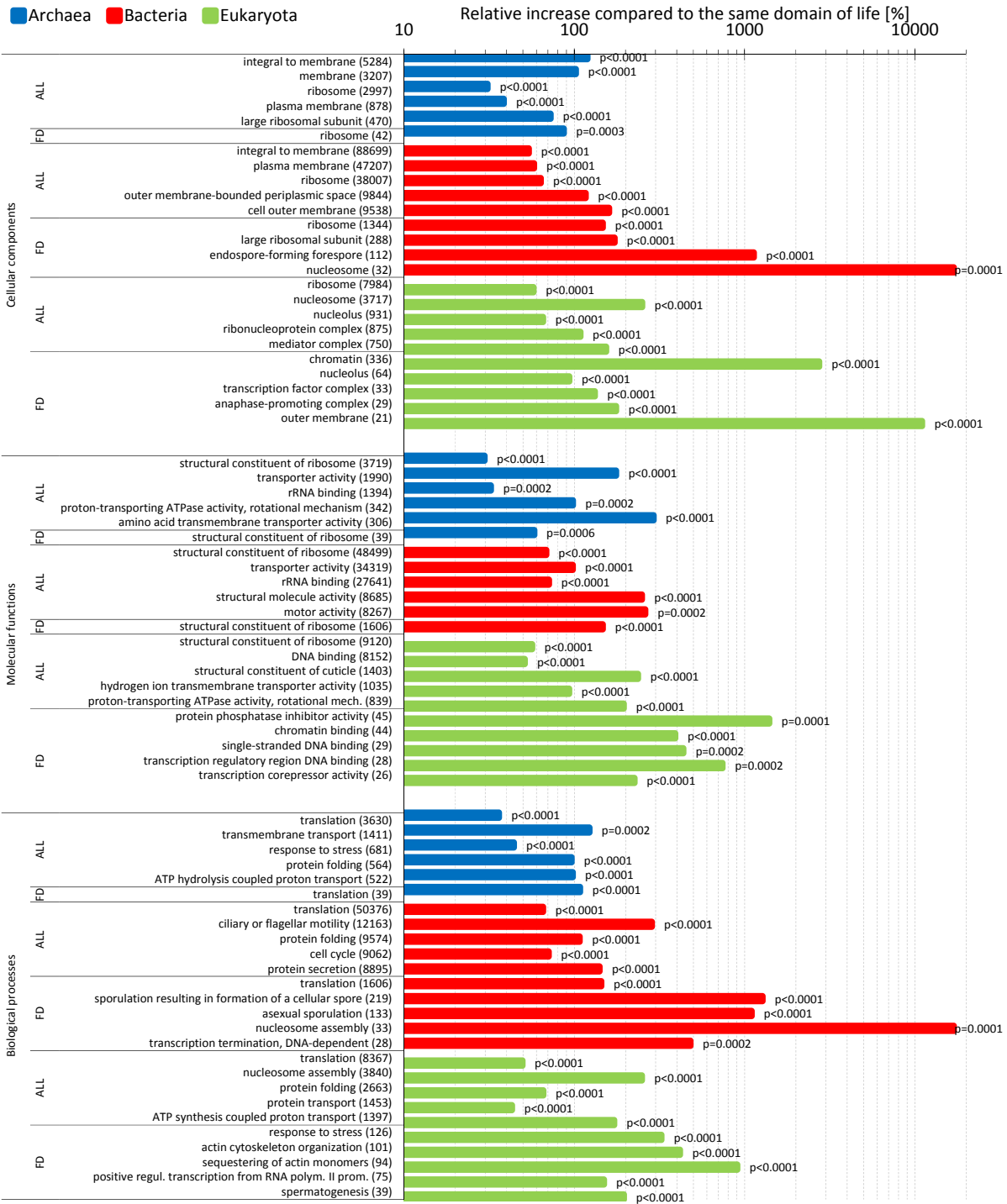


**Figure 4.** Comparison of the disorder content (panel A) and fraction of fully disordered proteins (panel B) between complete proteomes, C-depleted proteins and CFYWH-depleted proteins in the three domains of life.

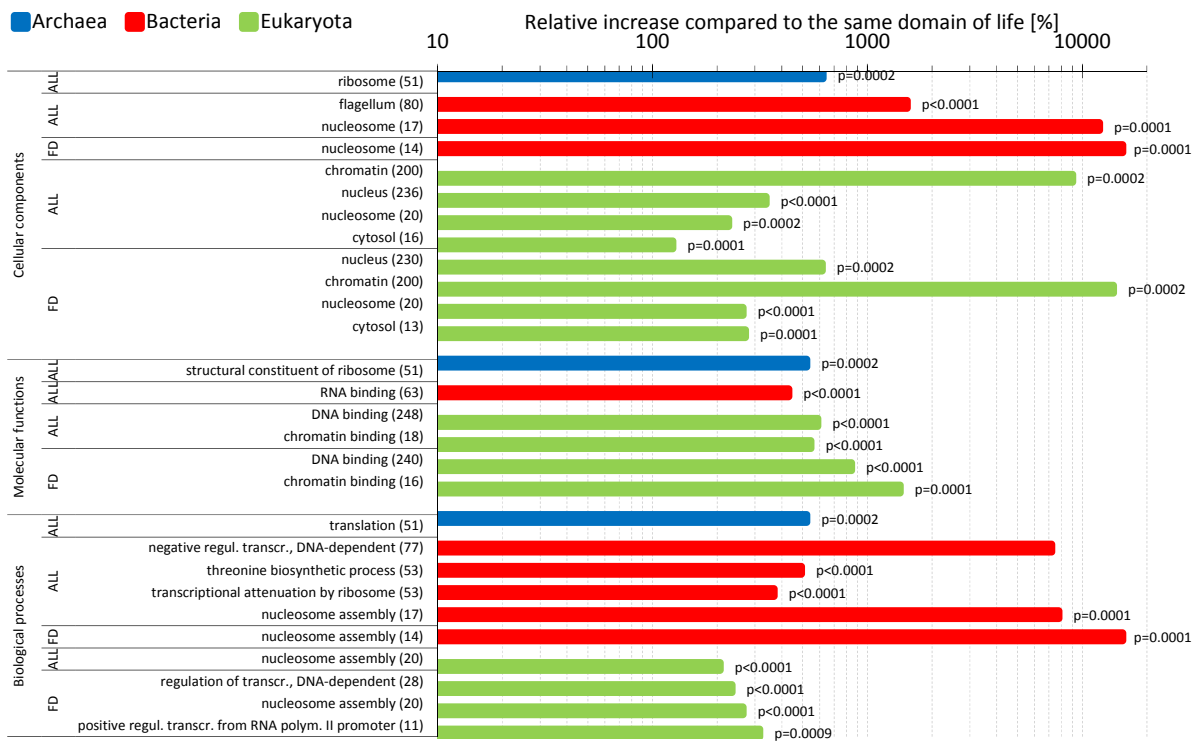




**Figure 5.** Relation between proteome-level abundance of the C-depleted proteins and structural coverage (panel A) and disorder content (panel B). Each point represents a single proteome. Lines represent linear fit into the data for a specific domain of life, which is accompanied with the corresponding values of the Person correlation coefficient (PCC).



**Figure 6.** Cellular components (top of the figure), molecular functions (middle of the figure) and biological processes (bottom of the figure) that are significantly enriched in all C-depleted proteins (ALL lines) and among fully disordered C-depleted proteins (FD lines). The analysis is performed separately for eukaryotic species (green bars), bacterial species (red bars) archaeal species (blue bars). The y-axis lists five most frequent (the corresponding number of annotated proteins is shown inside the brackets) and significantly enriched functions/components ( $p$ -value < 0.001 and enrichment > 30%). The x-axis shows the enrichment measured as relative increase in frequency when compared to the size and chain length matched set of randomly chosen proteins from the same domain of life. Details of the calculation are explained in the Materials and Methods section. The functions/cellular components are sorted, within each group, by the number of annotated proteins.



**Figure 7.** Cellular components (top of the figure), molecular functions (middle of the figure) and biological processes (bottom of the figure) that are significantly enriched in all CFYWH-depleted proteins (ALL lines) and among fully disordered CFYWH-depleted proteins (FD lines). The analysis is performed separately for eukaryotic species (green bars), bacterial species (red bars) and archaeal species (blue bars). The y-axis lists five most frequent (the corresponding number of annotated proteins is shown inside the brackets) and significantly enriched functions/components ( $p$ -value < 0.001 and enrichment > 30%). The x-axis shows the enrichment measured as relative increase in frequency when compared to the size and chain length matched set of randomly chosen proteins from the same domain of life. Details of the calculation are explained in the Materials and Methods section. The functions/cellular components are sorted, within each group, by the number of annotated proteins.

# Supplementary materials

The supplementary materials include the list and taxonomic classification of the 817 considered proteomes.

DomainOfLife	TaxonomicClassification
Archaea	Crenarchaeota Thermoprotei Acidianus
Archaea	Crenarchaeota Thermoprotei Acidilobus
Archaea	Crenarchaeota Thermoprotei Aeropyrum
Archaea	Crenarchaeota Thermoprotei Caldivirga
Archaea	Crenarchaeota Thermoprotei Desulfurococcus
Archaea	Crenarchaeota Thermoprotei Hyperthermus
Archaea	Crenarchaeota Thermoprotei Ignicoccus
Archaea	Crenarchaeota Thermoprotei Ignisphaera
Archaea	Crenarchaeota Thermoprotei Metallosphaera
Archaea	Crenarchaeota Thermoprotei Pyroaculum
Archaea	Crenarchaeota Thermoprotei Pyrolobus
Archaea	Crenarchaeota Thermoprotei Staphylothermus
Archaea	Crenarchaeota Thermoprotei Sulfolobus
Archaea	Crenarchaeota Thermoprotei Thermofilum
Archaea	Crenarchaeota Thermoprotei Thermoproteus
Archaea	Crenarchaeota Thermoprotei Thermosphaera
Archaea	Crenarchaeota Thermoprotei Vulcanisaeta
Archaea	Euryarchaeota Acidilipifundum
Archaea	Euryarchaeota Archaeoglobi Archaeoglobus
Archaea	Euryarchaeota Archaeoglobi Ferroglobus
Archaea	Euryarchaeota Halobacteria Halalkalicoccus
Archaea	Euryarchaeota Halobacteria Haloarcula
Archaea	Euryarchaeota Halobacteria Halobacterium
Archaea	Euryarchaeota Halobacteria Haloferax
Archaea	Euryarchaeota Halobacteria Halogeometricum
Archaea	Euryarchaeota Halobacteria Halomicrobium
Archaea	Euryarchaeota Halobacteria Halopiger
Archaea	Euryarchaeota Halobacteria Haloquadratum
Archaea	Euryarchaeota Halobacteria Halorhabdus
Archaea	Euryarchaeota Halobacteria Halorubrum
Archaea	Euryarchaeota Halobacteria Haloterrigena
Archaea	Euryarchaeota Halobacteria Natrialba
Archaea	Euryarchaeota Halobacteria Natronomonas
Archaea	Euryarchaeota Methanobacteria Methanobacterium
Archaea	Euryarchaeota Methanobacteria Methanobrevibacter
Archaea	Euryarchaeota Methanobacteria Methanosphaera
Archaea	Euryarchaeota Methanobacteria Methanothermobacter
Archaea	Euryarchaeota Methanobacteria Methanothermus
Archaea	Euryarchaeota Methanococci Methanocaldococcus
Archaea	Euryarchaeota Methanococci Methanococcus
Archaea	Euryarchaeota Methanococci Methanothermococcus
Archaea	Euryarchaeota Methanococci Methanoterris
Archaea	Euryarchaeota Methanomicrobia Methanocella
Archaea	Euryarchaeota Methanomicrobia Methanococcoides
Archaea	Euryarchaeota Methanomicrobia Methanocorpusculum
Archaea	Euryarchaeota Methanomicrobia Methanoculleus
Archaea	Euryarchaeota Methanomicrobia Methanohalobium
Archaea	Euryarchaeota Methanomicrobia Methanohalophilus
Archaea	Euryarchaeota Methanomicrobia Methanoplanus
Archaea	Euryarchaeota Methanomicrobia Methanoregula
Archaea	Euryarchaeota Methanomicrobia Methanoseta
Archaea	Euryarchaeota Methanomicrobia Methanosalsum
Archaea	Euryarchaeota Methanomicrobia Methanosarcina
Archaea	Euryarchaeota Methanomicrobia Methanosphaerula
Archaea	Euryarchaeota Methanomicrobia Methanospirillum
Archaea	Euryarchaeota Methanopyri Methanopyrus
Archaea	Euryarchaeota Thermococci Pyrococcus
Archaea	Euryarchaeota Thermococci Thermococcus
Archaea	Euryarchaeota Thermoplasmata Picrophilus
Archaea	Euryarchaeota Thermoplasmata Thermoplasma
Archaea	Korarchaeota Candidatus Korarchaeum
Archaea	Nanoarchaeota Nanoarchaeum
Archaea	Thaumarchaeota Cenarchaeales Cenarchaeum
Archaea	Thaumarchaeota Nitrosopumilales Nitrosopumilus
Bacteria	Acidobacteria Acidobacteriales Acidobacterium
Bacteria	Acidobacteria Acidobacteriales Granulicella
Bacteria	Acidobacteria Acidobacteriales Terriglobus
Bacteria	Acidobacteria Candidatus Chloracidobacterium

Bacteria Acidobacteria Candidatus Koribacter  
 Bacteria Acidobacteria Solibacteres Candidatus Solibacter  
 Bacteria Actinobacteria Acidimicrobidae Acidimicrobium  
 Bacteria Actinobacteria Actinobacteridae Acidothermus  
 Bacteria Actinobacteria Actinobacteridae Actinoplanes  
 Bacteria Actinobacteria Actinobacteridae Actinosynnema  
 Bacteria Actinobacteria Actinobacteridae Amycolatopsis  
 Bacteria Actinobacteria Actinobacteridae Amycolicoccus  
 Bacteria Actinobacteria Actinobacteridae Arcanobacterium  
 Bacteria Actinobacteria Actinobacteridae Arthrobacter  
 Bacteria Actinobacteria Actinobacteridae Beutenbergia  
 Bacteria Actinobacteria Actinobacteridae Bifidobacterium  
 Bacteria Actinobacteria Actinobacteridae Blastococcus  
 Bacteria Actinobacteria Actinobacteridae Brachybacterium  
 Bacteria Actinobacteria Actinobacteridae Catenuispora  
 Bacteria Actinobacteria Actinobacteridae Cellulomonas  
 Bacteria Actinobacteria Actinobacteridae Clavibacter  
 Bacteria Actinobacteria Actinobacteridae Corynebacterium  
 Bacteria Actinobacteria Actinobacteridae Frankia  
 Bacteria Actinobacteria Actinobacteridae Gardnerella  
 Bacteria Actinobacteria Actinobacteridae Geodermatophilus  
 Bacteria Actinobacteria Actinobacteridae Gordonia  
 Bacteria Actinobacteria Actinobacteridae Intrasporangium  
 Bacteria Actinobacteria Actinobacteridae Isoptericola  
 Bacteria Actinobacteria Actinobacteridae Jonesia  
 Bacteria Actinobacteria Actinobacteridae Kineococcus  
 Bacteria Actinobacteria Actinobacteridae Kitasatospora  
 Bacteria Actinobacteria Actinobacteridae Kocuria  
 Bacteria Actinobacteria Actinobacteridae Kribbella  
 Bacteria Actinobacteria Actinobacteridae Kytococcus  
 Bacteria Actinobacteria Actinobacteridae Leifsonia  
 Bacteria Actinobacteria Actinobacteridae Microbacterium  
 Bacteria Actinobacteria Actinobacteridae Micrococcus  
 Bacteria Actinobacteria Actinobacteridae Microlunatus  
 Bacteria Actinobacteria Actinobacteridae Micromonospora  
 Bacteria Actinobacteria Actinobacteridae Mobiluncus  
 Bacteria Actinobacteria Actinobacteridae Mycobacterium  
 Bacteria Actinobacteria Actinobacteridae Mycobacterium abscessus  
 Bacteria Actinobacteria Actinobacteridae Mycobacterium avium complex (MAC)  
 Bacteria Actinobacteria Actinobacteridae Mycobacterium tuberculosis complex  
 Bacteria Actinobacteria Actinobacteridae Nakamurella  
 Bacteria Actinobacteria Actinobacteridae Nocardia  
 Bacteria Actinobacteria Actinobacteridae Nocardioides  
 Bacteria Actinobacteria Actinobacteridae Nocardiosis  
 Bacteria Actinobacteria Actinobacteridae Propionibacterium  
 Bacteria Actinobacteria Actinobacteridae Pseudonocardia  
 Bacteria Actinobacteria Actinobacteridae Renibacterium  
 Bacteria Actinobacteria Actinobacteridae Rhodococcus  
 Bacteria Actinobacteria Actinobacteridae Rothia  
 Bacteria Actinobacteria Actinobacteridae Saccharomonospora  
 Bacteria Actinobacteria Actinobacteridae Saccharopolyspora  
 Bacteria Actinobacteria Actinobacteridae Salinispora  
 Bacteria Actinobacteria Actinobacteridae Sanguibacter  
 Bacteria Actinobacteria Actinobacteridae Segniliparus  
 Bacteria Actinobacteria Actinobacteridae Stackebrandtia  
 Bacteria Actinobacteria Actinobacteridae Streptomyces  
 Bacteria Actinobacteria Actinobacteridae Streptosporangium  
 Bacteria Actinobacteria Actinobacteridae Thermobifida  
 Bacteria Actinobacteria Actinobacteridae Thermobispora  
 Bacteria Actinobacteria Actinobacteridae Thermomonospora  
 Bacteria Actinobacteria Actinobacteridae Tropheryma  
 Bacteria Actinobacteria Actinobacteridae Tsukamurella  
 Bacteria Actinobacteria Actinobacteridae Verrucosipora  
 Bacteria Actinobacteria Actinobacteridae Xylanimonas  
 Bacteria Actinobacteria Coriobacteridae Atopobium  
 Bacteria Actinobacteria Coriobacteridae Coriobacterium  
 Bacteria Actinobacteria Coriobacteridae Cryptobacterium  
 Bacteria Actinobacteria Coriobacteridae Eggerthella  
 Bacteria Actinobacteria Coriobacteridae Olsenella  
 Bacteria Actinobacteria Coriobacteridae Slackia  
 Bacteria Actinobacteria Rubrobacteridae Conexibacter  
 Bacteria Actinobacteria Rubrobacteridae Rubrobacter  
 Bacteria Aquificae Aquificales Aquifex  
 Bacteria Aquificae Aquificales Desulfurobacterium  
 Bacteria Aquificae Aquificales Hydrogenobacter  
 Bacteria Aquificae Aquificales Hydrogenobaculum  
 Bacteria Aquificae Aquificales Persephonella  
 Bacteria Aquificae Aquificales Sulfurhydrogenibium

Bacteria Aquificae Aquificales Thermocrinis  
 Bacteria Aquificae Aquificales Thermovibrio  
 Bacteria Bacteroidetes Bacteroidetes Order II. Incertae sedis Rhodothermus  
 Bacteria Bacteroidetes Bacteroidetes Order II. Incertae sedis Salinibacter  
 Bacteria Bacteroidetes Bacteroidia Bacteroides  
 Bacteria Bacteroidetes Bacteroidia Candidatus Azobacteroides  
 Bacteria Bacteroidetes Bacteroidia Odoribacter  
 Bacteria Bacteroidetes Bacteroidia Paludibacter  
 Bacteria Bacteroidetes Bacteroidia Parabacteroides  
 Bacteria Bacteroidetes Bacteroidia Porphyromonas  
 Bacteria Bacteroidetes Bacteroidia Prevotella  
 Bacteria Bacteroidetes Bacteroidia Tannerella  
 Bacteria Bacteroidetes Candidatus Amoebophilus  
 Bacteria Bacteroidetes Cytophagia Cyclobacterium  
 Bacteria Bacteroidetes Cytophagia Cytophaga  
 Bacteria Bacteroidetes Cytophagia Dyadobacter  
 Bacteria Bacteroidetes Cytophagia Leadbetterella  
 Bacteria Bacteroidetes Cytophagia Marivirga  
 Bacteria Bacteroidetes Cytophagia Runella  
 Bacteria Bacteroidetes Cytophagia Spirosoma  
 Bacteria Bacteroidetes Flavobacteriia Blattabacterium  
 Bacteria Bacteroidetes Flavobacteriia Candidatus Sulcia  
 Bacteria Bacteroidetes Flavobacteriia Capnocytophaga  
 Bacteria Bacteroidetes Flavobacteriia Cellulophaga  
 Bacteria Bacteroidetes Flavobacteriia Croceibacter  
 Bacteria Bacteroidetes Flavobacteriia Flavobacteriaceae  
 Bacteria Bacteroidetes Flavobacteriia Flavobacterium  
 Bacteria Bacteroidetes Flavobacteriia Fluvicola  
 Bacteria Bacteroidetes Flavobacteriia Gramella  
 Bacteria Bacteroidetes Flavobacteriia Krokinobacter  
 Bacteria Bacteroidetes Flavobacteriia Lacinutrix  
 Bacteria Bacteroidetes Flavobacteriia Maribacter  
 Bacteria Bacteroidetes Flavobacteriia Muricauda  
 Bacteria Bacteroidetes Flavobacteriia Owenweeksia  
 Bacteria Bacteroidetes Flavobacteriia Riemerella  
 Bacteria Bacteroidetes Flavobacteriia Robiginitalea  
 Bacteria Bacteroidetes Flavobacteriia Weeksella  
 Bacteria Bacteroidetes Flavobacteriia Zobellia  
 Bacteria Bacteroidetes Flavobacteriia Zunongwangia  
 Bacteria Bacteroidetes Sphingobacteriia Chitinophaga  
 Bacteria Bacteroidetes Sphingobacteriia Haliscomenobacter  
 Bacteria Bacteroidetes Sphingobacteriia Niasella  
 Bacteria Bacteroidetes Sphingobacteriia Pedobacter  
 Bacteria Bacteroidetes Sphingobacteriia Saprospira  
 Bacteria Bacteroidetes Sphingobacteriia Solitalea  
 Bacteria Bacteroidetes Sphingobacteriia Sphingobacterium  
 Bacteria Caldiseptica Caldiseptica Caldisepticum  
 Bacteria candidate division WWE1 Candidatus Cloacamonas  
 Bacteria Chlamydiae Chlamydiales Candidatus Protochlamydia  
 Bacteria Chlamydiae Chlamydiales Chlamydia  
 Bacteria Chlamydiae Chlamydiales Chlamydomphila  
 Bacteria Chlamydiae Chlamydiales Parachlamydia  
 Bacteria Chlamydiae Chlamydiales Simkania  
 Bacteria Chlamydiae Chlamydiales Waddlia  
 Bacteria Chlorobi Chlorobia Chlorobaculum  
 Bacteria Chlorobi Chlorobia Chlorobium  
 Bacteria Chlorobi Chlorobia Chloroherpeton  
 Bacteria Chlorobi Chlorobia Pelodictyon  
 Bacteria Chlorobi Chlorobia Prosthecochloris  
 Bacteria Chloroflexi Anaerolineae Anaerolinea  
 Bacteria Chloroflexi Caldilineae Caldilinea  
 Bacteria Chloroflexi Chloroflexales Chloroflexus  
 Bacteria Chloroflexi Chloroflexales Roseiflexus  
 Bacteria Chloroflexi Dehalococcoidetes Dehalococcoides  
 Bacteria Chloroflexi Dehalococcoidetes Dehalogenimonas  
 Bacteria Chloroflexi Herpetosiphonales Herpetosiphon  
 Bacteria Chloroflexi Sphaerobacteridae Sphaerobacter  
 Bacteria Chloroflexi Thermomicrobiales Thermomicrobium  
 Bacteria Chrysiogenetes Chrysiogenales Desulfurispirillum  
 Bacteria Cyanobacteria Chroococcales  
 Bacteria Cyanobacteria Chroococcales Acaryochloris  
 Bacteria Cyanobacteria Chroococcales Cyanothece  
 Bacteria Cyanobacteria Chroococcales Microcystis  
 Bacteria Cyanobacteria Chroococcales Synechococcus  
 Bacteria Cyanobacteria Chroococcales Synechocystis  
 Bacteria Cyanobacteria Chroococcales Thermosynechococcus  
 Bacteria Cyanobacteria Gloeobacteria Gloeobacter  
 Bacteria Cyanobacteria Nostocales Anabaena

Bacteria Cyanobacteria Nostocales Nostoc  
 Bacteria Cyanobacteria Nostocales Trichormus  
 Bacteria Cyanobacteria Oscillatoriales Trichodesmium  
 Bacteria Cyanobacteria Prochlorophytes Prochlorococcus  
 Bacteria Deferribacteres Deferribacterales Deferribacter  
 Bacteria Deferribacteres Deferribacterales Deferribacteraceae  
 Bacteria Deferribacteres Deferribacterales Denitrovibrio  
 Bacteria Deferribacteres Deferribacterales Flexistipes  
 Bacteria Deinococcus-Thermus Deinococci Deinococcus  
 Bacteria Deinococcus-Thermus Deinococci Marinithermus  
 Bacteria Deinococcus-Thermus Deinococci Meiothermus  
 Bacteria Deinococcus-Thermus Deinococci Oceanithermus  
 Bacteria Deinococcus-Thermus Deinococci Thermus  
 Bacteria Deinococcus-Thermus Deinococci Truepera  
 Bacteria Dictyoglomi Dictyoglomales Dictyoglomus  
 Bacteria Elusimicrobia Elusimicrobia Elusimicrobium  
 Bacteria Elusimicrobia environmental samples  
 Bacteria Fibrobacteres Fibrobacterales Fibrobacter  
 Bacteria Firmicutes Bacillales Alicyclobacillus  
 Bacteria Firmicutes Bacillales Anoxybacillus  
 Bacteria Firmicutes Bacillales Bacillus  
 Bacteria Firmicutes Bacillales Bacillus cereus group  
 Bacteria Firmicutes Bacillales Brevibacillus  
 Bacteria Firmicutes Bacillales Exiguobacterium  
 Bacteria Firmicutes Bacillales Geobacillus  
 Bacteria Firmicutes Bacillales Halobacillus  
 Bacteria Firmicutes Bacillales Kyrpidia  
 Bacteria Firmicutes Bacillales Listeria  
 Bacteria Firmicutes Bacillales Lysinibacillus  
 Bacteria Firmicutes Bacillales Macrococcus  
 Bacteria Firmicutes Bacillales Oceanobacillus  
 Bacteria Firmicutes Bacillales Paenibacillus  
 Bacteria Firmicutes Bacillales Solibacillus  
 Bacteria Firmicutes Bacillales Staphylococcus  
 Bacteria Firmicutes Clostridia Acetobacterium  
 Bacteria Firmicutes Clostridia Acetohalobium  
 Bacteria Firmicutes Clostridia Alkaliphilus  
 Bacteria Firmicutes Clostridia Ammonifex  
 Bacteria Firmicutes Clostridia Anaerococcus  
 Bacteria Firmicutes Clostridia Butyrivibrio  
 Bacteria Firmicutes Clostridia Caldanaerobacter  
 Bacteria Firmicutes Clostridia Caldicellulosiruptor  
 Bacteria Firmicutes Clostridia Candidatus Arthromitus  
 Bacteria Firmicutes Clostridia Candidatus Desulforudis  
 Bacteria Firmicutes Clostridia Carboxydothermus  
 Bacteria Firmicutes Clostridia Cellulosilyticum  
 Bacteria Firmicutes Clostridia Clostridiales  
 Bacteria Firmicutes Clostridia Clostridium  
 Bacteria Firmicutes Clostridia Coprothermobacter  
 Bacteria Firmicutes Clostridia Desulfitobacterium  
 Bacteria Firmicutes Clostridia Desulfosporosinus  
 Bacteria Firmicutes Clostridia Desulfotomaculum  
 Bacteria Firmicutes Clostridia Ethanoligenens  
 Bacteria Firmicutes Clostridia Eubacterium  
 Bacteria Firmicutes Clostridia Filifactor  
 Bacteria Firmicutes Clostridia Finegoldia  
 Bacteria Firmicutes Clostridia Halanaerobium  
 Bacteria Firmicutes Clostridia Halothermothrix  
 Bacteria Firmicutes Clostridia Heliobacterium  
 Bacteria Firmicutes Clostridia Mahella  
 Bacteria Firmicutes Clostridia Moorella  
 Bacteria Firmicutes Clostridia Natranaerobius  
 Bacteria Firmicutes Clostridia Oscillibacter  
 Bacteria Firmicutes Clostridia Pelotomaculum  
 Bacteria Firmicutes Clostridia Peptostreptococcaceae  
 Bacteria Firmicutes Clostridia Roseburia  
 Bacteria Firmicutes Clostridia Ruminococcus  
 Bacteria Firmicutes Clostridia Sulfobacillus  
 Bacteria Firmicutes Clostridia Symbiobacterium  
 Bacteria Firmicutes Clostridia Syntrophobotulus  
 Bacteria Firmicutes Clostridia Syntrophomonas  
 Bacteria Firmicutes Clostridia Syntrophothermus  
 Bacteria Firmicutes Clostridia Tepidanaerobacter  
 Bacteria Firmicutes Clostridia Thermaerobacter  
 Bacteria Firmicutes Clostridia Thermicola  
 Bacteria Firmicutes Clostridia Thermoanaerobacter  
 Bacteria Firmicutes Clostridia Thermoanaerobacterium  
 Bacteria Firmicutes Clostridia Thermosediminibacter

Bacteria Firmicutes Erysipelotrichi Erysipelothrix  
 Bacteria Firmicutes Lactobacillales Aerococcus  
 Bacteria Firmicutes Lactobacillales Carnobacterium  
 Bacteria Firmicutes Lactobacillales Enterococcus  
 Bacteria Firmicutes Lactobacillales Lactobacillus  
 Bacteria Firmicutes Lactobacillales Lactococcus  
 Bacteria Firmicutes Lactobacillales Leuconostoc  
 Bacteria Firmicutes Lactobacillales Melissococcus  
 Bacteria Firmicutes Lactobacillales Oenococcus  
 Bacteria Firmicutes Lactobacillales Pediococcus  
 Bacteria Firmicutes Lactobacillales Streptococcus  
 Bacteria Firmicutes Lactobacillales Tetragenococcus  
 Bacteria Firmicutes Lactobacillales Weissella  
 Bacteria Firmicutes Negativicutes Acidaminococcus  
 Bacteria Firmicutes Negativicutes Selenomonas  
 Bacteria Firmicutes Negativicutes Veillonella  
 Bacteria Fusobacteria Fusobacterales Fusobacterium  
 Bacteria Fusobacteria Fusobacterales Ilyobacter  
 Bacteria Fusobacteria Fusobacterales Leptotrichia  
 Bacteria Fusobacteria Fusobacterales Sebaldella  
 Bacteria Fusobacteria Fusobacterales Streptobacillus  
 Bacteria Gemmatimonadetes Gemmatimonadales Gemmatimonas  
 Bacteria Ignavibacteria Ignavibacteria Ignavibacterium  
 Bacteria Nitrospirae Nitrospirales Leptospirillum  
 Bacteria Nitrospirae Nitrospirales Thermodesulfovibrio  
 Bacteria Planctomycetes Phycisphaerae Phycisphaera  
 Bacteria Planctomycetes Planctomycetia Isosphaera  
 Bacteria Planctomycetes Planctomycetia Pirellula  
 Bacteria Planctomycetes Planctomycetia Planctomyces  
 Bacteria Planctomycetes Planctomycetia Rhodopirellula  
 Bacteria Proteobacteria Alphaproteobacteria Acetobacter  
 Bacteria Proteobacteria Alphaproteobacteria Acidiphilium  
 Bacteria Proteobacteria Alphaproteobacteria Agrobacterium  
 Bacteria Proteobacteria Alphaproteobacteria Agrobacterium tumefaciens complex  
 Bacteria Proteobacteria Alphaproteobacteria Anaplasma  
 Bacteria Proteobacteria Alphaproteobacteria Asticcacaulis  
 Bacteria Proteobacteria Alphaproteobacteria Azorhizobium  
 Bacteria Proteobacteria Alphaproteobacteria Azospirillum  
 Bacteria Proteobacteria Alphaproteobacteria Bartonella  
 Bacteria Proteobacteria Alphaproteobacteria Beijerinckia  
 Bacteria Proteobacteria Alphaproteobacteria belli group  
 Bacteria Proteobacteria Alphaproteobacteria Bradyrhizobium  
 Bacteria Proteobacteria Alphaproteobacteria Brevundimonas  
 Bacteria Proteobacteria Alphaproteobacteria Brucella  
 Bacteria Proteobacteria Alphaproteobacteria Candidatus Hodgkinia  
 Bacteria Proteobacteria Alphaproteobacteria Candidatus Liberibacter  
 Bacteria Proteobacteria Alphaproteobacteria Candidatus Midichloria  
 Bacteria Proteobacteria Alphaproteobacteria Candidatus Pelagibacter  
 Bacteria Proteobacteria Alphaproteobacteria Candidatus Puniceispirillum  
 Bacteria Proteobacteria Alphaproteobacteria Caulobacter  
 Bacteria Proteobacteria Alphaproteobacteria Chelativorans  
 Bacteria Proteobacteria Alphaproteobacteria Dinoroseobacter  
 Bacteria Proteobacteria Alphaproteobacteria Ehrlichia  
 Bacteria Proteobacteria Alphaproteobacteria Erythrobacter  
 Bacteria Proteobacteria Alphaproteobacteria Gluconacetobacter  
 Bacteria Proteobacteria Alphaproteobacteria Gluconobacter  
 Bacteria Proteobacteria Alphaproteobacteria Granulibacter  
 Bacteria Proteobacteria Alphaproteobacteria Hirschia  
 Bacteria Proteobacteria Alphaproteobacteria Hyphomicrobium  
 Bacteria Proteobacteria Alphaproteobacteria Hyphomonas  
 Bacteria Proteobacteria Alphaproteobacteria Jannaschia  
 Bacteria Proteobacteria Alphaproteobacteria Ketogulonicigenium  
 Bacteria Proteobacteria Alphaproteobacteria Magnetococcus  
 Bacteria Proteobacteria Alphaproteobacteria Magnetospirillum  
 Bacteria Proteobacteria Alphaproteobacteria Maricaulis  
 Bacteria Proteobacteria Alphaproteobacteria Mesorhizobium  
 Bacteria Proteobacteria Alphaproteobacteria Methylobacterium  
 Bacteria Proteobacteria Alphaproteobacteria Methylocella  
 Bacteria Proteobacteria Alphaproteobacteria Micavibrio  
 Bacteria Proteobacteria Alphaproteobacteria Neorickettsia  
 Bacteria Proteobacteria Alphaproteobacteria Nitrobacter  
 Bacteria Proteobacteria Alphaproteobacteria Novosphingobium  
 Bacteria Proteobacteria Alphaproteobacteria Ochrobactrum  
 Bacteria Proteobacteria Alphaproteobacteria Oligotropha  
 Bacteria Proteobacteria Alphaproteobacteria Orientia  
 Bacteria Proteobacteria Alphaproteobacteria Paracoccus  
 Bacteria Proteobacteria Alphaproteobacteria Parvibaculum  
 Bacteria Proteobacteria Alphaproteobacteria Parvularcula



Bacteria Proteobacteria Alphaproteobacteria Pelagibacterium  
 Bacteria Proteobacteria Alphaproteobacteria phagocytophilum group  
 Bacteria Proteobacteria Alphaproteobacteria Phenyllobacterium  
 Bacteria Proteobacteria Alphaproteobacteria Polymorphum  
 Bacteria Proteobacteria Alphaproteobacteria Pseudovibrio  
 Bacteria Proteobacteria Alphaproteobacteria Rhizobium  
 Bacteria Proteobacteria Alphaproteobacteria Rhodobacter  
 Bacteria Proteobacteria Alphaproteobacteria Rhodomicrobium  
 Bacteria Proteobacteria Alphaproteobacteria Rhodopseudomonas  
 Bacteria Proteobacteria Alphaproteobacteria Rhodospirillum  
 Bacteria Proteobacteria Alphaproteobacteria Roseobacter  
 Bacteria Proteobacteria Alphaproteobacteria Ruegeria  
 Bacteria Proteobacteria Alphaproteobacteria Sinorhizobium  
 Bacteria Proteobacteria Alphaproteobacteria Sphingobium  
 Bacteria Proteobacteria Alphaproteobacteria Sphingomonas  
 Bacteria Proteobacteria Alphaproteobacteria Sphingopyxis  
 Bacteria Proteobacteria Alphaproteobacteria spotted fever group  
 Bacteria Proteobacteria Alphaproteobacteria Starkeya  
 Bacteria Proteobacteria Alphaproteobacteria typhus group  
 Bacteria Proteobacteria Alphaproteobacteria Wolbachia  
 Bacteria Proteobacteria Alphaproteobacteria Xanthobacter  
 Bacteria Proteobacteria Alphaproteobacteria Zymomonas  
 Bacteria Proteobacteria Betaproteobacteria Achromobacter  
 Bacteria Proteobacteria Betaproteobacteria Acidovorax  
 Bacteria Proteobacteria Betaproteobacteria Albidiferax  
 Bacteria Proteobacteria Betaproteobacteria Alicyclophilus  
 Bacteria Proteobacteria Betaproteobacteria Aromatoleum  
 Bacteria Proteobacteria Betaproteobacteria Azoarcus  
 Bacteria Proteobacteria Betaproteobacteria Azospira  
 Bacteria Proteobacteria Betaproteobacteria Bordetella  
 Bacteria Proteobacteria Betaproteobacteria Burkholderia  
 Bacteria Proteobacteria Betaproteobacteria Burkholderia cepacia complex  
 Bacteria Proteobacteria Betaproteobacteria Candidatus Accumulibacter  
 Bacteria Proteobacteria Betaproteobacteria Candidatus Tremblaya  
 Bacteria Proteobacteria Betaproteobacteria Candidatus Zinderia  
 Bacteria Proteobacteria Betaproteobacteria Chromobacterium  
 Bacteria Proteobacteria Betaproteobacteria Collimonas  
 Bacteria Proteobacteria Betaproteobacteria Comamonas  
 Bacteria Proteobacteria Betaproteobacteria Cupriavidus  
 Bacteria Proteobacteria Betaproteobacteria Dechloromonas  
 Bacteria Proteobacteria Betaproteobacteria Delfia  
 Bacteria Proteobacteria Betaproteobacteria Gallionella  
 Bacteria Proteobacteria Betaproteobacteria Herbaspirillum  
 Bacteria Proteobacteria Betaproteobacteria Herminiimonas  
 Bacteria Proteobacteria Betaproteobacteria Janthinobacterium  
 Bacteria Proteobacteria Betaproteobacteria Laribacter  
 Bacteria Proteobacteria Betaproteobacteria Leptothrix  
 Bacteria Proteobacteria Betaproteobacteria Methylibium  
 Bacteria Proteobacteria Betaproteobacteria Methylobacillus  
 Bacteria Proteobacteria Betaproteobacteria Methylothena  
 Bacteria Proteobacteria Betaproteobacteria Methylovorus  
 Bacteria Proteobacteria Betaproteobacteria Neisseria  
 Bacteria Proteobacteria Betaproteobacteria Nitrosomonas  
 Bacteria Proteobacteria Betaproteobacteria Nitrospira  
 Bacteria Proteobacteria Betaproteobacteria Polaromonas  
 Bacteria Proteobacteria Betaproteobacteria Polynucleobacter  
 Bacteria Proteobacteria Betaproteobacteria Pseudogulbenkiana  
 Bacteria Proteobacteria Betaproteobacteria pseudomallei group  
 Bacteria Proteobacteria Betaproteobacteria Pusillimonas  
 Bacteria Proteobacteria Betaproteobacteria Ralstonia  
 Bacteria Proteobacteria Betaproteobacteria Ramlibacter  
 Bacteria Proteobacteria Betaproteobacteria Rubrivivax  
 Bacteria Proteobacteria Betaproteobacteria Sideroxydans  
 Bacteria Proteobacteria Betaproteobacteria Taylorella  
 Bacteria Proteobacteria Betaproteobacteria Thauera  
 Bacteria Proteobacteria Betaproteobacteria Thiobacillus  
 Bacteria Proteobacteria Betaproteobacteria Thiomonas  
 Bacteria Proteobacteria Betaproteobacteria Variovorax  
 Bacteria Proteobacteria Betaproteobacteria Verminephrobacter  
 Bacteria Proteobacteria Deltaproteobacteria Anaeromyxobacter  
 Bacteria Proteobacteria Deltaproteobacteria Bacteriovorax  
 Bacteria Proteobacteria Deltaproteobacteria Bdellovibrio  
 Bacteria Proteobacteria Deltaproteobacteria Coralloccoccus  
 Bacteria Proteobacteria Deltaproteobacteria Desulfarculus  
 Bacteria Proteobacteria Deltaproteobacteria Desulfatibacillum  
 Bacteria Proteobacteria Deltaproteobacteria Desulfobacca  
 Bacteria Proteobacteria Deltaproteobacteria Desulfobacterium  
 Bacteria Proteobacteria Deltaproteobacteria Desulfobulbus

Bacteria Proteobacteria Deltaproteobacteria Desulfococcus  
 Bacteria Proteobacteria Deltaproteobacteria Desulfohalobium  
 Bacteria Proteobacteria Deltaproteobacteria Desulfomicrobium  
 Bacteria Proteobacteria Deltaproteobacteria Desulfotalea  
 Bacteria Proteobacteria Deltaproteobacteria Desulfovibrio  
 Bacteria Proteobacteria Deltaproteobacteria Desulfurivibrio  
 Bacteria Proteobacteria Deltaproteobacteria Geobacter  
 Bacteria Proteobacteria Deltaproteobacteria Haliangium  
 Bacteria Proteobacteria Deltaproteobacteria Hippea  
 Bacteria Proteobacteria Deltaproteobacteria Lawsonia  
 Bacteria Proteobacteria Deltaproteobacteria Myxococcus  
 Bacteria Proteobacteria Deltaproteobacteria Pelobacter  
 Bacteria Proteobacteria Deltaproteobacteria Sorangium  
 Bacteria Proteobacteria Deltaproteobacteria Stigmatella  
 Bacteria Proteobacteria Deltaproteobacteria Syntrophobacter  
 Bacteria Proteobacteria Deltaproteobacteria Syntrophus  
 Bacteria Proteobacteria Epsilonproteobacteria Arcobacter  
 Bacteria Proteobacteria Epsilonproteobacteria Campylobacter  
 Bacteria Proteobacteria Epsilonproteobacteria Helicobacter  
 Bacteria Proteobacteria Epsilonproteobacteria Nautilia  
 Bacteria Proteobacteria Epsilonproteobacteria Nitratifactor  
 Bacteria Proteobacteria Epsilonproteobacteria Nitratiruptor  
 Bacteria Proteobacteria Epsilonproteobacteria Sulfuricurvum  
 Bacteria Proteobacteria Epsilonproteobacteria Sulfurimonas  
 Bacteria Proteobacteria Epsilonproteobacteria Sulfurospirillum  
 Bacteria Proteobacteria Epsilonproteobacteria Sulfurovum  
 Bacteria Proteobacteria Epsilonproteobacteria Wolinella  
 Bacteria Proteobacteria Gammaproteobacteria Acidithiobacillus  
 Bacteria Proteobacteria Gammaproteobacteria Acinetobacter  
 Bacteria Proteobacteria Gammaproteobacteria Acinetobacter calcoaceticus or baumannii complex  
 Bacteria Proteobacteria Gammaproteobacteria Actinobacillus  
 Bacteria Proteobacteria Gammaproteobacteria Aeromonas  
 Bacteria Proteobacteria Gammaproteobacteria Aggregatibacter  
 Bacteria Proteobacteria Gammaproteobacteria Alcanivorax  
 Bacteria Proteobacteria Gammaproteobacteria Allivibrio  
 Bacteria Proteobacteria Gammaproteobacteria Alkalilimnicola  
 Bacteria Proteobacteria Gammaproteobacteria Allochromatium  
 Bacteria Proteobacteria Gammaproteobacteria Alteromonas  
 Bacteria Proteobacteria Gammaproteobacteria Azotobacter  
 Bacteria Proteobacteria Gammaproteobacteria Basfia  
 Bacteria Proteobacteria Gammaproteobacteria Buchnera  
 Bacteria Proteobacteria Gammaproteobacteria Candidatus Baumannia  
 Bacteria Proteobacteria Gammaproteobacteria Candidatus Blochmannia  
 Bacteria Proteobacteria Gammaproteobacteria Candidatus Carsonella  
 Bacteria Proteobacteria Gammaproteobacteria Candidatus Hamiltonella  
 Bacteria Proteobacteria Gammaproteobacteria Candidatus Moranella  
 Bacteria Proteobacteria Gammaproteobacteria Candidatus Riesia  
 Bacteria Proteobacteria Gammaproteobacteria Cellvibrio  
 Bacteria Proteobacteria Gammaproteobacteria Chromohalobacter  
 Bacteria Proteobacteria Gammaproteobacteria Citrobacter  
 Bacteria Proteobacteria Gammaproteobacteria Colwellia  
 Bacteria Proteobacteria Gammaproteobacteria Coxiella  
 Bacteria Proteobacteria Gammaproteobacteria Cronobacter  
 Bacteria Proteobacteria Gammaproteobacteria Dichelobacter  
 Bacteria Proteobacteria Gammaproteobacteria Dickeya  
 Bacteria Proteobacteria Gammaproteobacteria Edwardsiella  
 Bacteria Proteobacteria Gammaproteobacteria Enterobacter  
 Bacteria Proteobacteria Gammaproteobacteria Enterobacter cloacae complex  
 Bacteria Proteobacteria Gammaproteobacteria Erwinia  
 Bacteria Proteobacteria Gammaproteobacteria Escherichia  
 Bacteria Proteobacteria Gammaproteobacteria Ferrimonas  
 Bacteria Proteobacteria Gammaproteobacteria Francisella  
 Bacteria Proteobacteria Gammaproteobacteria Frateuria  
 Bacteria Proteobacteria Gammaproteobacteria Gallibacterium  
 Bacteria Proteobacteria Gammaproteobacteria Glaciecola  
 Bacteria Proteobacteria Gammaproteobacteria Haemophilus  
 Bacteria Proteobacteria Gammaproteobacteria Hahella  
 Bacteria Proteobacteria Gammaproteobacteria Halomonas  
 Bacteria Proteobacteria Gammaproteobacteria Halorhodospira  
 Bacteria Proteobacteria Gammaproteobacteria Halotheobacillus  
 Bacteria Proteobacteria Gammaproteobacteria Histophilus  
 Bacteria Proteobacteria Gammaproteobacteria Idiomarina  
 Bacteria Proteobacteria Gammaproteobacteria Kangiella  
 Bacteria Proteobacteria Gammaproteobacteria Klebsiella  
 Bacteria Proteobacteria Gammaproteobacteria Legionella  
 Bacteria Proteobacteria Gammaproteobacteria Listonella  
 Bacteria Proteobacteria Gammaproteobacteria Marinobacter  
 Bacteria Proteobacteria Gammaproteobacteria Marinomonas

Bacteria Proteobacteria Gammaproteobacteria Methylococcus  
 Bacteria Proteobacteria Gammaproteobacteria Methylochromium  
 Bacteria Proteobacteria Gammaproteobacteria Methylomonas  
 Bacteria Proteobacteria Gammaproteobacteria Moraxella  
 Bacteria Proteobacteria Gammaproteobacteria Nitrosococcus  
 Bacteria Proteobacteria Gammaproteobacteria Oceanimonas  
 Bacteria Proteobacteria Gammaproteobacteria Pantoea  
 Bacteria Proteobacteria Gammaproteobacteria Pasteurella  
 Bacteria Proteobacteria Gammaproteobacteria Pectobacterium  
 Bacteria Proteobacteria Gammaproteobacteria Photobacterium  
 Bacteria Proteobacteria Gammaproteobacteria Photorhabdus  
 Bacteria Proteobacteria Gammaproteobacteria Proteus  
 Bacteria Proteobacteria Gammaproteobacteria Providencia  
 Bacteria Proteobacteria Gammaproteobacteria Pseudoalteromonas  
 Bacteria Proteobacteria Gammaproteobacteria Pseudomonas  
 Bacteria Proteobacteria Gammaproteobacteria Pseudomonas syringae  
 Bacteria Proteobacteria Gammaproteobacteria Pseudoxanthomonas  
 Bacteria Proteobacteria Gammaproteobacteria Psychrobacter  
 Bacteria Proteobacteria Gammaproteobacteria Psychromonas  
 Bacteria Proteobacteria Gammaproteobacteria Rahnella  
 Bacteria Proteobacteria Gammaproteobacteria Saccharophagus  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella dublin  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella enterica subsp. enterica serovar Agona  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella enterica subsp. enterica serovar Choleraesuis  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella enterica subsp. enterica serovar Gallinarum or pullorum  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella enterica subsp. enterica serovar Heidelberg  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella enterica subsp. enterica serovar Paratyphi B  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella enterica subsp. enterica serovar Paratyphi C  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella enteritidis  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella gallinarum  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella newport  
 Bacteria Proteobacteria Gammaproteobacteria Salmonella typhimurium  
 Bacteria Proteobacteria Gammaproteobacteria Serratia  
 Bacteria Proteobacteria Gammaproteobacteria Shewanella  
 Bacteria Proteobacteria Gammaproteobacteria Shigella  
 Bacteria Proteobacteria Gammaproteobacteria Sodalis  
 Bacteria Proteobacteria Gammaproteobacteria Stenotrophomonas maltophilia group  
 Bacteria Proteobacteria Gammaproteobacteria sulfur-oxidizing symbionts  
 Bacteria Proteobacteria Gammaproteobacteria Teredinibacter  
 Bacteria Proteobacteria Gammaproteobacteria Thioalkalimicrobium  
 Bacteria Proteobacteria Gammaproteobacteria Thioalkalivibrio  
 Bacteria Proteobacteria Gammaproteobacteria Thiomicrospira  
 Bacteria Proteobacteria Gammaproteobacteria Tolomonas  
 Bacteria Proteobacteria Gammaproteobacteria Vibrio  
 Bacteria Proteobacteria Gammaproteobacteria Wigglesworthia  
 Bacteria Proteobacteria Gammaproteobacteria Xanthomonas  
 Bacteria Proteobacteria Gammaproteobacteria Xenorhabdus  
 Bacteria Proteobacteria Gammaproteobacteria Xylella  
 Bacteria Proteobacteria Gammaproteobacteria Yersinia  
 Bacteria Spirochaetes Spirochaetales Borrelia  
 Bacteria Spirochaetes Spirochaetales Borrelia burgdorferi group  
 Bacteria Spirochaetes Spirochaetales Brachyspira  
 Bacteria Spirochaetes Spirochaetales Leptospira  
 Bacteria Spirochaetes Spirochaetales Sphaerochaeta  
 Bacteria Spirochaetes Spirochaetales Spirochaeta  
 Bacteria Spirochaetes Spirochaetales Treponema  
 Bacteria Synergistetes Synergistia Aminobacterium  
 Bacteria Synergistetes Synergistia Thermanaerovibrio  
 Bacteria Synergistetes Synergistia Thermovirga  
 Bacteria Tenericutes Mollicutes 16SrX (Apple proliferation group)  
 Bacteria Tenericutes Mollicutes 16SrXII (Stolbur group)  
 Bacteria Tenericutes Mollicutes Acholeplasma  
 Bacteria Tenericutes Mollicutes Candidatus Phytoplasma asteris  
 Bacteria Tenericutes Mollicutes Mesoplasma  
 Bacteria Tenericutes Mollicutes Mycoplasma  
 Bacteria Tenericutes Mollicutes Ureaplasma  
 Bacteria Thermobaculum  
 Bacteria Thermodesulfobacteria Thermodesulfobacteriales Thermodesulfator  
 Bacteria Thermodesulfobacteria Thermodesulfobacteriales Thermodesulfobacterium  
 Bacteria Thermotogae Thermotogales Fervidobacterium  
 Bacteria Thermotogae Thermotogales Kosmotoga  
 Bacteria Thermotogae Thermotogales Marinitoga  
 Bacteria Thermotogae Thermotogales Petrotoga  
 Bacteria Thermotogae Thermotogales Thermosiphon  
 Bacteria Thermotogae Thermotogales Thermotoga  
 Bacteria Verrucomicrobia Opiritae Coraliomargarita  
 Bacteria Verrucomicrobia Opiritae Opiritus

Bacteria Verrucomicrobia unclassified Verrucomicrobia Methylocyathum  
 Bacteria Verrucomicrobia Verrucomicrobiae Akkermansia  
  
 Eukaryota Alveolata Apicomplexa Babesia  
 Eukaryota Alveolata Apicomplexa Cryptosporidium  
 Eukaryota Alveolata Apicomplexa Neospora  
 Eukaryota Alveolata Apicomplexa Plasmodium (Laverania)  
 Eukaryota Alveolata Apicomplexa Plasmodium (Plasmodium)  
 Eukaryota Alveolata Apicomplexa Plasmodium (Vinckeia)  
 Eukaryota Alveolata Apicomplexa Theileria  
 Eukaryota Alveolata Apicomplexa Toxoplasma  
 Eukaryota Alveolata Ciliophora Ichthyophthirius  
 Eukaryota Alveolata Ciliophora Paramecium  
 Eukaryota Alveolata Ciliophora Tetrahymena  
 Eukaryota Alveolata Perkinsea Perkinsus  
 Eukaryota Amoebozoa Archamoebae Entamoeba  
 Eukaryota Amoebozoa Mycetozoa Dictyostelium  
 Eukaryota Amoebozoa Mycetozoa Polysphondylium  
 Eukaryota Choanoflagellida Codonosigidae Monosiga  
 Eukaryota Choanoflagellida Salpingoecidae Salpingoeca  
 Eukaryota Diplomonadida Hexamitidae Giardia  
 Eukaryota Euglenozoa Kinetoplastida Duttonella  
 Eukaryota Euglenozoa Kinetoplastida Leishmania  
 Eukaryota Euglenozoa Kinetoplastida Leishmania braziliensis species complex  
 Eukaryota Euglenozoa Kinetoplastida Nannomonas  
 Eukaryota Euglenozoa Kinetoplastida Schizotrypanum  
 Eukaryota Euglenozoa Kinetoplastida Trypanosoma  
 Eukaryota Fungi Chytridiomycota Batrachochytrium  
 Eukaryota Fungi Dikarya Ajellomyces  
 Eukaryota Fungi Dikarya Arthroderma  
 Eukaryota Fungi Dikarya Aspergillus  
 Eukaryota Fungi Dikarya Botryotinia  
 Eukaryota Fungi Dikarya Candida  
 Eukaryota Fungi Dikarya Chaetomium  
 Eukaryota Fungi Dikarya Clavispora  
 Eukaryota Fungi Dikarya Coccidioides  
 Eukaryota Fungi Dikarya Colletotrichum  
 Eukaryota Fungi Dikarya Coprinopsis  
 Eukaryota Fungi Dikarya Cordyceps  
 Eukaryota Fungi Dikarya Debaryomyces  
 Eukaryota Fungi Dikarya Emericella  
 Eukaryota Fungi Dikarya Eremothecium  
 Eukaryota Fungi Dikarya Exophiala  
 Eukaryota Fungi Dikarya Filobasidiella or Cryptococcus neoformans species complex  
 Eukaryota Fungi Dikarya Fusarium oxysporum species complex  
 Eukaryota Fungi Dikarya Gibberella  
 Eukaryota Fungi Dikarya Glarea  
 Eukaryota Fungi Dikarya Glomerella  
 Eukaryota Fungi Dikarya Grosmannia  
 Eukaryota Fungi Dikarya Hypocrea  
 Eukaryota Fungi Dikarya Kazachstania  
 Eukaryota Fungi Dikarya Kluyveromyces  
 Eukaryota Fungi Dikarya Komagataella  
 Eukaryota Fungi Dikarya Laccaria  
 Eukaryota Fungi Dikarya Lachancea  
 Eukaryota Fungi Dikarya Leptosphaeria maculans complex  
 Eukaryota Fungi Dikarya Lodderomyces  
 Eukaryota Fungi Dikarya Magnaporthe  
 Eukaryota Fungi Dikarya Malassezia  
 Eukaryota Fungi Dikarya Melampsora  
 Eukaryota Fungi Dikarya Metarhizium  
 Eukaryota Fungi Dikarya Meyerozyma  
 Eukaryota Fungi Dikarya Milleroyzma  
 Eukaryota Fungi Dikarya mitosporic Nakaseomyces  
 Eukaryota Fungi Dikarya Mixia  
 Eukaryota Fungi Dikarya Moniliophthora  
 Eukaryota Fungi Dikarya Myceliophthora  
 Eukaryota Fungi Dikarya Naumovozyma  
 Eukaryota Fungi Dikarya Nectria haematococca complex  
 Eukaryota Fungi Dikarya Neosartorya  
 Eukaryota Fungi Dikarya Neurospora  
 Eukaryota Fungi Dikarya Ogataea  
 Eukaryota Fungi Dikarya Orbilia  
 Eukaryota Fungi Dikarya Paracoccidioides  
 Eukaryota Fungi Dikarya Penicillium  
 Eukaryota Fungi Dikarya Penicillium chrysogenum complex  
 Eukaryota Fungi Dikarya Phaeosphaeria  
 Eukaryota Fungi Dikarya Piriformospora

Eukaryota Fungi Dikarya Podospora  
 Eukaryota Fungi Dikarya Postia  
 Eukaryota Fungi Dikarya Puccinia  
 Eukaryota Fungi Dikarya Pyrenophora  
 Eukaryota Fungi Dikarya Rhodotorula  
 Eukaryota Fungi Dikarya Saccharomyces  
 Eukaryota Fungi Dikarya Scheffersomyces  
 Eukaryota Fungi Dikarya Schizophyllum  
 Eukaryota Fungi Dikarya Schizosaccharomyces  
 Eukaryota Fungi Dikarya Sclerotinia  
 Eukaryota Fungi Dikarya Serpula  
 Eukaryota Fungi Dikarya Sordaria  
 Eukaryota Fungi Dikarya Spathaspora  
 Eukaryota Fungi Dikarya Sporisorium  
 Eukaryota Fungi Dikarya Talaromyces  
 Eukaryota Fungi Dikarya Tetrapisispora  
 Eukaryota Fungi Dikarya Thielavia  
 Eukaryota Fungi Dikarya Torulaspora  
 Eukaryota Fungi Dikarya Trichoderma  
 Eukaryota Fungi Dikarya Trichophyton  
 Eukaryota Fungi Dikarya Tuber  
 Eukaryota Fungi Dikarya Uncinocarpus  
 Eukaryota Fungi Dikarya Ustilago  
 Eukaryota Fungi Dikarya Vanderwaltozyma  
 Eukaryota Fungi Dikarya Verticillium  
 Eukaryota Fungi Dikarya Yarrowia  
 Eukaryota Fungi Dikarya Zygosaccharomyces  
 Eukaryota Fungi Dikarya Zymoseptoria  
 Eukaryota Fungi Fungi incertae sedis Rhizopus  
 Eukaryota Fungi Microsporidia Encephalitozoon  
 Eukaryota Fungi Microsporidia Enterocytozoon  
 Eukaryota Fungi Microsporidia Nematocida  
 Eukaryota Fungi Microsporidia Nosema  
 Eukaryota Heterolobosea Schizopyrenida Naegleria  
 Eukaryota Ichthyosporea Capsaspora  
 Eukaryota Metazoa Arthropoda Acromyrmex  
 Eukaryota Metazoa Arthropoda Anopheles  
 Eukaryota Metazoa Arthropoda Apis  
 Eukaryota Metazoa Arthropoda Atta  
 Eukaryota Metazoa Arthropoda Bombyx  
 Eukaryota Metazoa Arthropoda Camponotus  
 Eukaryota Metazoa Arthropoda Culex  
 Eukaryota Metazoa Arthropoda Danaus  
 Eukaryota Metazoa Arthropoda Daphnia  
 Eukaryota Metazoa Arthropoda Drosophila  
 Eukaryota Metazoa Arthropoda Harpegnathos  
 Eukaryota Metazoa Arthropoda Hawaiian Drosophila  
 Eukaryota Metazoa Arthropoda Ixodes  
 Eukaryota Metazoa Arthropoda Pediculus  
 Eukaryota Metazoa Arthropoda Solenopsis  
 Eukaryota Metazoa Arthropoda Sophophora  
 Eukaryota Metazoa Arthropoda Stegomyia  
 Eukaryota Metazoa Arthropoda Tribolium  
 Eukaryota Metazoa Chordata Ailuropoda  
 Eukaryota Metazoa Chordata Anolis  
 Eukaryota Metazoa Chordata Bos  
 Eukaryota Metazoa Chordata Branchiostoma  
 Eukaryota Metazoa Chordata Callithrix  
 Eukaryota Metazoa Chordata Canis  
 Eukaryota Metazoa Chordata Cavia  
 Eukaryota Metazoa Chordata Ciona  
 Eukaryota Metazoa Chordata Cricetulus  
 Eukaryota Metazoa Chordata Danio  
 Eukaryota Metazoa Chordata Equus  
 Eukaryota Metazoa Chordata Gallus  
 Eukaryota Metazoa Chordata Gasterosteus  
 Eukaryota Metazoa Chordata Gorilla  
 Eukaryota Metazoa Chordata Heterocephalus  
 Eukaryota Metazoa Chordata Homo  
 Eukaryota Metazoa Chordata Ictidomys  
 Eukaryota Metazoa Chordata Latimeria  
 Eukaryota Metazoa Chordata Loxodonta  
 Eukaryota Metazoa Chordata Macaca  
 Eukaryota Metazoa Chordata Meleagris  
 Eukaryota Metazoa Chordata Monodelphis  
 Eukaryota Metazoa Chordata Mus  
 Eukaryota Metazoa Chordata Myotis  
 Eukaryota Metazoa Chordata Nomascus

Eukaryota Metazoa Chordata Oikopleura  
 Eukaryota Metazoa Chordata Oreochromis  
 Eukaryota Metazoa Chordata Ornithorhynchus  
 Eukaryota Metazoa Chordata Oryctolagus  
 Eukaryota Metazoa Chordata Oryzias  
 Eukaryota Metazoa Chordata Otolemur  
 Eukaryota Metazoa Chordata Pan  
 Eukaryota Metazoa Chordata Pongo  
 Eukaryota Metazoa Chordata Rattus  
 Eukaryota Metazoa Chordata Sarcophilus  
 Eukaryota Metazoa Chordata Silurana  
 Eukaryota Metazoa Chordata Sus  
 Eukaryota Metazoa Chordata Taeniopygia  
 Eukaryota Metazoa Chordata Takifugu  
 Eukaryota Metazoa Chordata Tetraodon  
 Eukaryota Metazoa Cnidaria Nematostella  
 Eukaryota Metazoa Echinodermata Strongylocentrotus  
 Eukaryota Metazoa Nematoda Brugia  
 Eukaryota Metazoa Nematoda Caenorhabditis  
 Eukaryota Metazoa Nematoda Loa  
 Eukaryota Metazoa Nematoda Pristionchus  
 Eukaryota Metazoa Nematoda Trichinella  
 Eukaryota Metazoa Placozoa Trichoplax  
 Eukaryota Metazoa Platyhelminthes Clonorchis  
 Eukaryota Metazoa Platyhelminthes Schistosoma  
 Eukaryota Metazoa Porifera Amphimedon  
 Eukaryota Parabasalia Trichomonadida Trichomonas  
 Eukaryota stramenopiles Bacillariophyta Phaeodactylum  
 Eukaryota stramenopiles Bacillariophyta Thalassiosira  
 Eukaryota stramenopiles Blastocystis  
 Eukaryota stramenopiles Oomycetes Phytophthora  
 Eukaryota stramenopiles Pelagophyceae Aureococcus  
 Eukaryota stramenopiles PX clade Ectocarpus  
 Eukaryota Viridiplantae Chlorophyta Chlamydomonas  
 Eukaryota Viridiplantae Chlorophyta Chlorella  
 Eukaryota Viridiplantae Chlorophyta Micromonas  
 Eukaryota Viridiplantae Chlorophyta Ostreococcus  
 Eukaryota Viridiplantae Chlorophyta Volvox  
 Eukaryota Viridiplantae Streptophyta Arabidopsis  
 Eukaryota Viridiplantae Streptophyta Brachypodium  
 Eukaryota Viridiplantae Streptophyta Glycine  
 Eukaryota Viridiplantae Streptophyta Oryza  
 Eukaryota Viridiplantae Streptophyta Physcomitrella  
 Eukaryota Viridiplantae Streptophyta Populus  
 Eukaryota Viridiplantae Streptophyta Ricinus  
 Eukaryota Viridiplantae Streptophyta Selaginella  
 Eukaryota Viridiplantae Streptophyta Sorghum  
 Eukaryota Viridiplantae Streptophyta Vitis