# Intrinsic disorder in protein domains contributes to both organism complexity and clade-specific functions

## (Additional information)

Chao Gao[1, #], Chong Ma[1, 2, #], Huqiang Wang[1], Haolin Zhong[1], Jiayin Zang[1], Rugang Zhong[2], Fuchu He[1, *], Dong Yang[1, #, *]

[1] State Key Laboratory of Proteomics, Beijing Proteome Research Center, National Center for Protein Sciences (Beijing), Beijing Institute of Lifeomics, Beijing 102206, China

[2] Beijing Key Laboratory of Environmental and Viral Oncology, College of Life Science and Bioengineering, Beijing University of Technology, Beijing 100124, China

[#] Equal contribution.

[*] Corresponding authors.

Fuchu He, National Center for Protein Sciences Beijing, 38 Science Park Road, Changping District, Beijing 102206, China.

Tel: 8610- 61771001, E-mail: hefc@nic.bmi.ac.cn.

Dong Yang, National Center for Protein Sciences Beijing, 38 Science Park Road, Changping District, Beijing 102206, China.

Tel: 8610-61777051, Fax: 8610-61777050, E-mail: yangdongbprc@163.com.

# Contents

**Supplemental results**

- The special features of the intrinsic disordered domains (IDDs) with a consecutive disorder region (CDR) but a low disorder ratio (DSDR < 30%)

- Which physico−chemical characteristic of amino acid was correlated with values of P matrix?

- The relationship between DSDR variation and organism complexity

**Supplemental references**

**Supplemental figures**

Figure S1. Distribution of IDDs (CDRN ≥ 1) across the three superkingdoms (related to figure 1).

Figure S2. Distribution of IDDs (DSDR>30%, or CDRN ≥ 1) across 25 representative species using MobiDB-lite data.

Figure S3. Diagram of 'dominant category' approach for the disorder degree classification of domain families.

Figure S4. Calculation of the degree of domain disorder in archaea using different methods, including a corrected method.

Figure S5. Scatter-plot of the median repetition number of completely structured domains (blue triangle) and IDDs (red circle).

Figure S6. The correlation between structural disorder and organism complexity analyzed based on MobiDB-lite data.

Figure S7. Box-plot of the percentage of intrinsically disordered (CDRN $\geqslant$ 1) domains (left panel) or domain families (right panel) in each species across the different phyla of eukaryotes.

Figure S8. The distribution of IDDs in each age grade of the 25 representative species (related to figure 3C).

Figure S9. Figure S9. Comparison of the normalized DSDR values of the same domain families in evolutionarily more complex species and simpler species (related to figure 3D).

Figure S10. The number and functions of special IDDs in archaeal halobacteria.

Figure S11. The special functions of IDDs in Ascomycota and Basidiomycota revealed by the enriched mutant phenotypes.

Figure S12. Over- or under-representation analysis of BPs for IDDs with different ages and different distribution width.

Figure S13. The examples for the domains classified by domain age and disorder width (related to figure S12).

Figure S14. Comparison of the PTM site proportion in different type of domains/regions classified according to disorder degree.

Figure S15. Scatter plot of the percentage of disordered domain families in each species of the three superkingdoms (related to figure 5A-C).

Figure S16. Scatter plot of the percentage of non-domain disordered regions in each

species of the three superkingdoms (related to figure 5A-C).

Figure S17. The special features of the IDDs with one or more consecutive disorder region (CDR) but a low disorder ratio (DSDR<30%) (related to figure 5K).

Figure S18. The distribution of IDDs and young domains in different repeating domain categories.

Figure S19. The method for the normalization of DSDR values and the classification of DSDR variation.

Figure S20. Correlation between the values of physico−chemical characteristics of amino acid and the values in P matrix.

Figure S21. The correlation between DSDR variation values and domain sequence similarity and supplemental results about the distribution and functional features of domains with different DSDR variations.

Figure S22. Comparison of domain length among the protein domains with different disorder degrees.

**Supplemental tables**

**Note: The tables S1-S10 are in an additional Excel file.**

**Table S1. Domain numbers and their percentages in different disorder grades across the three superkingdoms.** DSDRs were divided into two grades according to values ($\leq 30\%$ and $> 30\%$). CDRNs were also divided into two grades ($0, \geq 1$). Domain numbers and percentages of each disorder grade were calculated in each species. SPOT-D (A), IUPred (B) were used for disorder prediction of 1870 species, and the disorder

data of 58 species from MobiDB-lite (C) were also used for validation analysis.

**Table S2. Domain family numbers and their percentages in different disorder grades across the three superkingdoms.** Each domain family was assigned with a dominant category based on dominant category approach (see Materials and Methods). The numbers and the percentages of domain families of different disorder dominant categories were calculated in each species. SPOT-D (A), IUPred (B) were used for disorder prediction of 1870 species, and the disorder data of 58 species from MobiDB-lite (C) were also used for validation analysis.

**Table S3. Colour description of each phylum and class in figure 1A and supplemental figure 1A.**

**Table S4. Median domain repetition counts in different disordered grades in each species** (related to figure S5). For each species, the median repetition counts of all the domain families at each disorder grade were calculated. SPOT-D (A) and IUPred (B) are used for disorder prediction.

**Table S5. The significantly over-represented PTM types in intrinsically disordered domains.**

**Table S6. Domain numbers and percentages of different disorder grades in different domain categories classified by domain repetitions.** SPOT-D (A) and IUPred (B) are used for disorder prediction.

**Table S7. Summary of DSDR variation in twenty-five representative species (related to figure 6B).** SPOT-D (A) and IUPred (B) are used for disorder prediction.

**Table S8. The raw data of correlation analysis between physico－chemical**

**characteristics and P matrix values.**

**Table S9. The correlation between DSDR variation and organism complexity.** A1, B1. The results of the correlation analysis between DSDR variation and organism complexity. A2, B2. The raw data for the correlation analysis between DSDR variation and organism complexity. SPOT-D (A) and IUPred (B) are used for disorder prediction.

**Table S10 Detailed information of the no-vairation and high-variation domains in twenty-five representative species.**

**Datasets**

**Dataset 1. The disorder values calculated using SPOT-D and IUPred for each amino acid residue of all proteins in the 1870 species across 3 superkingdoms.** The scores of each residue were shown according to the order of each protein sequence. The letters represent the residues in protein sequences.

**Dataset 2. DSDR and CDRN values of each domain in the 1870 species.**

**Dataset 3. The median normalized DSDR values of all protein domain families in each species across three superkingdoms** (related to Figure 3A). DSDR was normalized by the interpolation method (see *Methods* section). The normalized DSDR values of the same domain family within one species is represented by the median value. "NA" means no such domain in this species.

**Dataset 4. Sequence similarities for the domain repeats of 25 representative species.** In the title line of each table, 'PrID_DomID_SiteID-1-PrID_DomID_SiteID-2' denotes each domain pair. 'Degree_of_variation' represents the variation of DSDR.

**Dataset 5. Domain GO annotations used in the study.**

These five datasets are available at http://dis-domain-data.ncpsb.org/.

**Supplemental results**

**The special features of the IDDs with one or more consecutive disorder region (CDR) but a low disorder ratio (DSDR < 30%)**

The comparison of domain disorder type between fungi and metazoans reflected the difference of domain disorder characteristics between them. Fungi have more class II IDDs, which have a consecutive disorder region (CDR) but a low disorder ratio (DSDR<30%). From the viewpoint of domain length, it is clear that class II IDDs are significantly longer than class IV IDDs (Figure S17A). The long peptide chains of the class II IDDs tend to contain a continuous disordered region (CDR), which will increase its flexibility to form more flexible conformations via the disordered linker region. For example, fungal_trans domain (a fungi-specific domain) in the transcriptional activator SEF1 of *Saccharomyces cerevisiae* contains a CDR (40 aa, from Asp$^{523}$ to Glu$^{551}$), which is located in the middle region of fungal_trans domain (from Thr$^{357}$ to Ile$^{577}$) (Figure S17B). The structure of SEF1 was modeled using SWISS-MODEL [1-5]. The CDR in the fungal_trans domain is seated between two coils, which is similar to a flexible chain linking two compact rods (Figure S17C). It implies that the CDR in fungal_trans domain may play as an entropic chain or a linker to adjust the local conformation to bind to DNA as SEF1 is a transcription factor [6].

Another fungi-specific class II domain, cAMP phosphodiesterases class-II (PDEase_II) domain in Nadsonia fulvescens with a long CDR (45 aa, Figure S17D, E), is involved in heterocyclic metabolic process. PDEase_II catalyzes the hydrolysis of

cAMP to 5' nucleoside monophosphate. CDR may regulate this process by regulating the flexibility of the domain. Interestingly, the segment from amino acid position 209 to 258 of this protein, which highly coincides with the CDR (203-260), is predicted to be a polar region by MobiDB-lite [7]. Since this protein can catalyse the hydrolysis of cAMP, the disordered polar region may be responsible for increasing the affinity of PDEase_II to water molecules.

**Which physico−chemical characteristic of amino acid was correlated with values of P matrix?**

The degree of disorder is determined by the physical-chemical properties of amino acid residues within the domain regions and the amino acid residues close to the domain regions [8]. The relationship between the amino acid physical-chemical properties and the values of the P matrix used in the IUPred [8, 9] method was investigated (Table S8). The correlation between values in P matrix [8] and the average values of the hydrophobicity index (Table S8) was significantly negative (R = −0.47, P = 6.8e−013, Figure S20A). But there are no such significant correlations between values in P matrix and average values of molecular weight (Figure S20B) or isoelectric points (Figure S20C). This result indicated that hydrophobicity of amino acid residues plays important roles in determining the degree of disorder.

**The relationship between DSDR variation and organism complexity**

Fifty-one eukaryotes were selected for correlation analysis. We found that there was a significant positive correlation between the number of 'high-variation' domains and the

number of cell types, which indicates that variation of the disorder degree of protein domains substantially contribute to organism complexity (Table S9). We speculated the high variation of domain disorder may lead to its versatility and regulate the complex signal networks more precisely. However, after normalized by the total number of repeating domains, this correlation disappeared. We can see from table S9 that there is a stronger correlation between the number of repeating domains and organism complexity. Thus, if we calculate the correlation between the percentages of 'high-variation' domains and cell type numbers in a species, there is no significant correlation.

# Supplemental References

1. Guex, N., M.C. Peitsch, and T. Schwede, *Automated comparative protein structure modeling with SWISS-MODEL and Swiss-PdbViewer: a historical perspective.* Electrophoresis, 2009. **30 Suppl 1**: p. S162-73.

2. Benkert, P., M. Biasini, and T. Schwede, *Toward the estimation of the absolute quality of individual protein structure models.* Bioinformatics, 2011. **27**(3): p. 343-50.

3. Bienert, S., et al., *The SWISS-MODEL Repository-new features and functionality.* Nucleic Acids Res, 2017. **45**(D1): p. D313-D319.

4. Waterhouse, A., et al., *SWISS-MODEL: homology modelling of protein structures and complexes.* Nucleic Acids Res, 2018. **46**(W1): p. W296-W303.

5. Bertoni, M., et al., *Modeling protein quaternary structure of homo- and hetero-oligomers beyond binary interactions by homology.* Sci Rep, 2017. **7**(1): p. 10480.

6. van Peij, N.N., J. Visser, and L.H. de Graaff, *Isolation and analysis of xlnR, encoding a transcriptional activator co-ordinating xylanolytic expression in Aspergillus niger.* Mol Microbiol, 1998. **27**(1): p. 131-42.

7. Necci, M., et al., *MobiDB-lite: fast and highly specific consensus prediction of intrinsic disorder in proteins.* Bioinformatics, 2017. **33**(9): p. 1402-1404.

8. Dosztanyi, Z., et al., *The pairwise energy content estimated from amino acid composition discriminates between folded and intrinsically unstructured proteins.* J Mol Biol, 2005. **347**(4): p. 827-39.

9. Habchi, J., et al., *Introducing protein intrinsic disorder.* Chem Rev, 2014. **114**(13): p. 6561-88.

**Supplemental figure legends**

**Figure S1. Distribution of IDDs (CDRN ≥ 1) across the three superkingdoms (related to figure 1).** A. The percentage of IDDs (CDRN $\geqslant$ 1) in each species from archaea (light green), bacteria (lime) and eukaryotes. Each bar indicates the results of a species. All species are arranged according to taxonomy information from NCBI database. The superkingdoms are separated by solid lines and the kingdoms in eukaryotes, including protists (light blue), metazoa (aqua), fungi (yellow) and plant (pale goldenrod), are separated by dashed lines. Each phylum is filled with different colors in the middle circle. The outer circle represents different classes. The detailed information about the names and colors of each species, class and phylum is presented in supplementary table S1-S3. The abbreviations in this figure: Ha, the class of Halobacteria; CDRN, consecutive disordered region number. B. Box-plot of the percentage of intrinsically disordered (CDRN $\geqslant$ 1) domain (purple) or domain families (red) in each species across the three superkingdoms (C, excluding halobacterias). For box-plots, the values of upper and lower quartile are indicated as upper and lower edges of the box, and the median values are indicated as a bar in the box. The differences in the percentage of disordered domain distribution between different categories are examined by Mann–Whitney U test. The corrected *P* values are shown at the top of each panel. D. Cumulative probability of the percentage of intrinsically disordered domain (solid line) or domain family (dashed line) in archaea (blue), bacteria (green) and eukaryotes (red) (E, excluding halobacterias).

**Figure S2. Distribution of IDDs across 25 representative species using MobiDB-lite data.** Box-plot of the percentage of intrinsically disordered (DSDR>30% in A, B, or CDRN $\geqslant 1$ in C, D) domain (purple) or domain families (red) in each species (B, D, excluding halobacterias). For box-plots, the values of upper and lower quartile are indicated as upper and lower edges of the box, and the median values are indicated as a bar in the box. The differences in the percentage of disordered domain distribution between different categories are examined by Mann–Whitney U test. The corrected $P$ values are shown at the top of each panel.

**Figure S3. Diagram of 'dominant category' approach for the disorder degree classification of domain families.** The yellow rectangles represent one domain family and they have eight repeats among three proteins. Each repeat has a DSDR value and a CDRN value. We calculated the percentage of each DSDR and CDRN grade and the dominant grade (> 0.5) was the DSDR and CDRN dominant category of this domain family, which represented the disordered feature of this domain family.

**Figure S4. Calculation of the degree of domain disorder in archaea using different methods, including a corrected method.** Box-plots show the percentage of disordered domain in each species in different phyla of Archaea. SPOT-D (A and B), IUPred (C and D), ESpritz (E and F) and ESpritz (PSI-BLAST) (G and H) were used for intrinsic disorder prediction. Both results based on DSDR (A, C, E and G) and CDRN (B, D, F and H) are shown.

**Figure S5. Scatter-plot of the median repetition number of completely structured domains (blue triangle) and IDDs (red circle).** Each species has two corresponding

domain repetition counts. The numbers of domains are shown on the label of the X-axis (number of species in which completely structured domain repetition counts is larger than more disordered domain repetition counts/number of all species in each superkingdom).

**Figure S6. The correlation between structural disorder and organism complexity analyzed based on MobiDB-lite data. (A)** Protein disorder degree was measured by the average of all representative proteins (the longest one for each gene) in each species. **(B)** Proteins with more than 30% disordered amino acids (PSDR > 30%) were regarded as disordered proteins. Protein disorder degree was measured by the percentage of disordered proteins in each species. **(C, D)** At domain level, structural disorder of each species was measured by the percentage of IDDs (C, DSDR > 30%, or D, CDRN ≥ 1). **(E, F)** At domain family level, domain families were defined using dominant category method; structural disorder of each species was measured by the percentage of intrinsically disordered domain families (E, family_DSDR > 30%, or F, family_CDRN ≥ 1). Scatter plots was made for organism complexity and structural disorder values obtained by different methods in fifty-one eukaryotes. Spearman method was used for the correlation analysis. The correlation coefficients (R) and P values (P) are shown in the inset, among which the significant results (P<0.01) are shown as red.

**Figure S7. Box-plot of the percentage of intrinsically disordered (CDRN ≥ 1) domains (left panel) or domain families (right panel) in each species across the different phyla of eukaryotes.** The four kingdoms of eukaryotes are marked on the bottom of the figure. The values of the upper and lower quartile, and the median values

are indicated as bars in the boxes. The differences in the percentage of disordered domain distribution between different categories are examined by Mann – Whitney U test. The corrected $P$ values are marked as stars in sub-figure G (*, p < 0.05; **, p < 0.01).

**Figure S8. The distribution of IDDs in each age grade of the 25 representative species (related to figure 3C).** Domains of each representative species are divided into different age grades. **A.** The distribution of IDDs (CDRN $\geq$ 1) in each age grade of the 25 representative species using SPOTD disorder degree results. **B, C.** The same analysis to figure 3C and A in this figure, using MobiDB-lite disorder degree results.

**Figure S9. Comparison of the normalized DSDR values of the same domain families in evolutionarily more complex species and simpler species** (related to figure 3D). Scatter-plot shows the median normalized DSDR values of the same domain families in complex species (top half of the figure) and simple species (bottom half of the figure). The red dots represent the difference of the median normalized DSDR value of each domain family between complex species and simple ones. The number and percentage of domain families belonging to each differential category are shown in brackets.

**Figure S10. The number and functions of special IDDs in archaeal halobacteria.**
A. The difference between the IDDs in halobacterias and other archaea was shown in Venn diagram. B. Function analysis of the halobacterias-specific IDDs. C, D. Two examples of halobacterias-specific IDDs.

**Figure S11. The special functions of IDDs in Ascomycota and Basidiomycota**

**revealed by the enriched mutant phenotypes.** Over- or under-representation analysis of biological process (BP) and cellular component (CC) for the disordered domains. The values of ±log(P) were transformed into 14 grades (−7 to +7): −7, log(P) ≤ −8; −6, −8 < log(P) ≤ −6; −5, −6 < log(P) ≤ −4; −4, −4 < log(P) ≤ −2; −3, −2< log(P) ≤ −1.5; −2, −1.5 < log(P) ≤ −1.301; −0.25, −1.301 < log(P) ≤ 0; 0.25, 0 < −log(P) < 1.301; 2, 1.301 ≤ −log(P) <1.5; 3, 1.5≤ −log(P) < 2; 4, 2 ≤ −log(P) < 4; 5, 4 ≤ −log(P) < 6; 6, 6 ≤ −log(P) < 8; 7, −log(P) ≥ 8.

**Figure S12. Over- or under-representation analysis of BPs for IDDs with different ages and different distribution width.** The IDDs are divided into 6 categories. The over- or under-representation strengths of each class were represented by -log (P) or log (P), respectively. Heat map showing the grades of over- or under-representation strengths, scoped from -7 to 7 (See Method for detail).

**Figure S13. The examples for the domains classified by domain age and disorder width (related to figure S12).** The heatmaps show the normalized DSDR values of each domain in different clades. One protein instance for each domain was shown. The domain features, including domain age, disorder ratio, are at the bottom of each panel.

**Figure S14. Comparison of the PTM site proportion in different type of domains/regions classified according to disorder degree.** Box-plot of the PTM site proportion in protein domains with different disorder degree classified according to DSDR (A) or CDRN (B), or in CDR regions or other regions in domains (C). The values of upper and lower quartile, and the median values are indicated as bars in the boxes. The differences between the neighboring categories are examined by Mann-Whitney U

test. The corrected P values are shown. Six species were included in the analysis.

**Figure S15. Scatter plot of the percentage of disordered domain families in each species of the three superkingdoms (related to figure 5A-C)**. Each point represents a species. The dashed diagonal line represents the line where Y and X are equal. The equations in each subfigure were obtained by fitting a straight line based on the data of all species of archaea (A), bacteria (B) and fungi or metazoan (C).

**Figure S16. Scatter plot of the percentage of non-domain disordered regions in each species of the three superkingdoms (related to figure 5A-C).** The N-terminal, C-terminal and linker regions in proteins were included in the analysis. See figure S7 for detailed descriptions.

**Figure S17. The special features of the IDDs with a consecutive disorder region (CDR) but a low disorder ratio (DSDR<30%) (related to figure 5K).** A. Comparison of the length between the class II and IV domains. B. Disorder prediction of transcriptional activator SEF1. C. The structure of a part of the sequence of SEF1, modeled by SWISS-MODEL. D. Disorder prediction of 3',5'-cyclic-nucleotide phosphodiesterase. E. The structure of 3',5'-cyclic-nucleotide phosphodiesterase modeled by SWISS-MODEL. The dashed lines in subfigure B and D represent the start and end positions of the protein domains and CDRs in them.

**Figure S18. The distribution of IDDs and young domains in different repeating domain categories.** Distribution of IDDs (A, B) and young domains (C) in different repeating domain categories of the twenty-five representative species. Dark spots represent statically significant differences of IDDs in the domains encoded only once

in the genome, tested by Fisher's exact test. Dark spots above the red dashed line represent that the domains encoded only once in the genome are over-represented compared to other types, while those below the blue dashed line represent they are under-represented. The abbreviations: *Hs, Homo sapiens; Sp, Strongylocentrotus purpuratus; Dm, Drosophila melanogaster; Cel, Caenorhabditis elegans; To, Thalassiosira oceanica; Em, Eimeria mitis; Eb, Enterocytozoon bieneusi; Up, Umbilicaria pustulata; Sc, Saccharomyces cerevisiae; At, Arabidopsis thaliana; Cr, Chlamydomonas reinhardtii; Ct, Chloroherpeton thalassium; Cex, Caldisericum exile; Ai, Alistipes inops; N, Nitrospina sp. SCGC_AAA799_A02; Ks, Kytococcus sedentarius; Td, Thermosulfurimonas dismutans; Ec, Escherichia coli; Lf, Leptospirillum ferriphilum; H, Hydrogenivirga sp. 128-5-R1-1; Cp, Chlamydia psittaci; Ht, Haloterrigena turkmenica; Kc, Candidatus korarchaeum cryptofilum; My, Metallosphaera yellowstonensis; No, Candidatus nitrocosmicus oleophilus*.

**Figure S19. The method for the normalization of DSDR values and the classification of DSDR variation.** A. Normalized curve for the DSDR values. Interpolation method was used to normalize the DSDRs. Four points were used to calculate the equation: $(0, 0)$, $(10, 1)$, $(30, 2)$ and $(100, 3)$. B. Grades of DSDR variation. The values denote the difference values of each pair of normalized DSDRs. C. Diagram of DSDR variation dominant category approach. Each domain pair has a DSDR variation value. We calculated the percentage of each DSDR variation grade and the dominant grade $(> 0.5)$ was the DSDR variation dominant category of this domain family, which represented the level of variation of this domain family.

**Figure S20. Correlation between the values of physico−chemical characteristics of amino acid and the values in P matrix.** Physico−chemical characteristics of amino acids, such as hydropathy index (A), molecular weight (B) and isoelectric points (C) may influence the value of disorder score of each amino acid residue, then affect the values of DSDR. The cross represents the median of equal-sized bins. Whiskers encompass the range from a quarter to three quarters of values. The correlation coefficients (R) and $P$ values ($P$) are shown in the inset.

**Figure S21. The correlation between DSDR variation values and domain sequence similarity (A-D) and supplemental results about the distribution and functional features of domains with different DSDR variations (E-F). A-B**. Correlation between DSDR variation and domain sequence similarity of each repeating domain pair. Scatter diagrams of human DSDR variation and domain sequence similarity in the domains repeating in the same proteins (A) or the domains repeating in multiple proteins (B). Each panel was divided into four regions according to the threshold low/high variation (DSDR variation = 1) and low/high similarity (40%). The bold numbers denote the numbers of domain pairs in each region. The abbreviations: LS, low similarity; HS, high similarity; LV, low variation; HV, high variation. **C-D**. Heatmap shows the numbers and percentages of four types of domain pairs in the 25 representative species. The percentages of low-similarity domain pairs (LS/Total) and the percentages of low-DSDR variation domain pairs among the low-similarity domain pairs (LS&LV/LS) in each species are listed on the right of the heat map. Sub-figure C shows for the domains repeating in the same proteins and D shows the domains

repeating in multiple proteins. **E.** Distribution of No-variation dominant categories of DSDR variation in the domains repeating in the same proteins and the domains appearing among multiple proteins of the 25 representative species using MobiDB-lite disorder result. The DSDR variation value was calculated by the difference of normalized DSDR. Black spots represent statically significant differences, tested by the Fisher's exact test. Black spots above the dashed line represent significant over representations. The abbreviations: *Hs, Homo sapiens; Sp, Strongylocentrotus purpuratus; Dm, Drosophila melanogaster; Cel, Caenorhabditis elegans; Eca, Enterospora canceri; Lp, Lasallia pustulata; Sc, Saccharomyces cerevisiae; To, Thalassiosira oceanica; Cc, Cyclospora cayetanensis; At, Arabidopsis thaliana; Cr, Chlamydomonas reinhardtii; Ct, Chloroherpeton thalassium; Cex, Caldisericum exile; Ai, Alistipes inops; N, Nitrospina sp. SCGC_AAA799_A02; Ks, Kytococcus sedentarius; Td, Thermosulfurimonas dismutans; Ec, Escherichia coli; Lf, Leptospirillum ferriphilum; H, Hydrogenivirga sp. 128-5-R1-1; Cp, Chlamydia psittaci; Hd, Haloterrigena daqingensis; Kc, Candidatus korarchaeum cryptofilum; My, Metallosphaera yellowstonensis; No, Candidatus nitrocosmicus oleophilus.* **F.** Examples of 'no-variation' domains. DSDR was normalized by interpolation method. Heatmap shows the normalized DSDR values of the same domain family in different proteins of different species. Each column represents a different species, and each row represents the corresponding orthologs. The abbreviations of the species: *Han, Helianthus annuus; Gs, Galdieria sulphuraria; Cc, Chondrus crispus; Cm,*

*Cyanidioschyzon merolae; Dc, Daucus carota; Atr, Amborella trichopoda; At, Arabidopsis thaliana; Bn, Brassica napus; Mt, Medicago truncatula.*
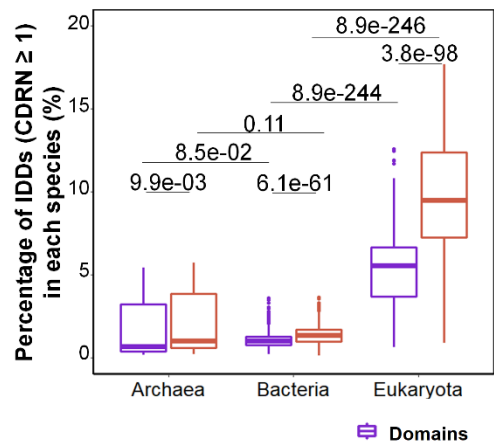
**Figure S22. Comparison of domain length among the protein domains with different disorder degrees.** Protein domains were classified according to CDRN (A-C) or DSDR (D-F) grades. The protein domains from archaea (A, D), bacteria (B, E) and eukaryotes (C, F) were analyzed respectively. The values of upper and lower quartile are indicated as upper and lower edges of the box, and the median values are indicated as a bar in the box. The differences in the domain length between different categories are examined by Mann–Whitney U test. The corrected *P* values are shown at the top of each panel.

**Figure S1**

**A**



Percentage of
intrinsically disordered domains(IDDs)
in each species

Disordered domain: CDRN≥1
Disordered domain family: CDRN≥1

**B**



8.9e-246
3.8e-98
8.9e-244
0.11
8.5e-02
9.9e-03
6.1e-61

Percentage of IDDs (CDRN ≥ 1) in each species (%)

Archaea    Bacteria    Eukaryota

**C**



8.9e-246
3.8e-98
8.9e-244
4.1e-16
1.4e-16
6.3e-04
6.1e-61

Archaea (No ha)    Bacteria    Eukaryota

Domains    Domain families

**D**



Cumulative probability

Category
Domains
Domain families

Superkingdom
Archaea
Bacteria
Eukaryota

Percentage of IDDs (CDRN ≥ 1)
in each species (%)

**E**



Category
Domains
Domain families

Superkingdom
Archaea (No ha)
Bacteria
Eukaryota

Percentage of IDDs (CDRN ≥ 1)
each species (%)

**Figure S2**

**Figure S3**

**Figure S4**

**Figure S5**

**Figure S6**

**Figure S7**

**Figure S8**



**SPOTD result**

**MobiDB-lite result**

**MobiDB-lite result**

Domain age grade:

- Cellular organism
- Archaea
- Bacteria
- Proteobacteria
- Actinobacteria
- Eukaryota
- Viridiplantae
- Streptophyta
- Fungi
- Metazoa
- Nematoda
- Insecta
- Chordata
- Mammalia

**Figure S9**

**Figure S10**

**A**

Halobacteria    Other archaea

380    82    136

**B**



Cellular amino acid metabolic process

Oxoacid metabolic process

Carboxylic acid metabolic process

Organic acid metabolic process

Oxidation-reduction process

Domain Number
- 10.0
- 12.5
- 15.0

$-\log_{10}(Pvalue)$
6
5
4
3

Species
· Halobacteria

Enrichment Score

**C**

Ubiquitinol-cytochrome C reductase Fe-S subunit TAT signal (PF10399)
Halobacterias-specific disordered domain
BP: Oxidation-reduction process; Metabolic process

ELK53625 (D320_12255)    *Haloferax sp. bab2207*
Multicopper oxidase is involved in copper homeostasis and oxidative stress response.



1    100    200    385

TAT_signal    Cu-oxidase_3    Cu-oxidase_2

DSDR (%):    52.00    63.16    42.61
Domain Age:    Cellular organisms    Cellular organisms    Cellular organisms

disorder
order

**Figure S11**

**Figure S12**

# Figure S13



**A** Androgen receptor (PF02166)
Eukaryote-specific domain, disordered in ≥50% species (50|54)
BP: Regulation of biosynthetic process; Signal transduction; Regulation of transcription

**B** Phage capsid scaffolding protein (GPO) serine peptidase (PF05929)
Prokaryote-specific domain, disordered in ≥50% species (22|43)
BP: Interspecies interaction between organisms; Symbiont process

**C** Ribosomal protein L36 (PF00444)
Common domain, disordered in ≥50% species (669|1216)
BP: Translation; Nitrogen compound metabolic process; Protein metabolic process

**D** Adaptin N terminal region (PF01602)
Eukaryote-specific domain, disordered in <50% species (0|537)
BP: Intracellular transport; Protein transport; Nitrogen compound transport

**E** XhoI (PF04555)
Prokaryote-specific domain, disordered in <50% species (0|20)
BP: DNA modification; DNA metabolic process; Nucleic acid metabolic process

**F** Cytochrome P450 (PF00067)
Common domain, disordered in <50% species (0|935)
BP: Oxidation-reduction; Metabolic process

Normalized DSDR
0  0.5  1  1.5  2  2.5  3

**Figure S14**

**Figure S15**

# Figure S16



**N_terminal**

**A** Y = 1.1928*X - 1.0823
R² = 0.981

**B** Y = 1.3894*X + 0.6085
R² = 0.948

**C** Y = 1.2793*X - 0.8614
R² = 0.9874

**Linker**

**D** Y = 1.0675*X - 0.0585
R² = 0.9126

**E** Y = 1.1342*X + 0.0812
R² = 0.9316

**F** Y = 1.4013*X + 1.0006
R² = 0.7323

**C_terminal**

**G** Y = 1.0278*X + 8.3232
R² = 0.8781

**H** Y = 1.1623*X + 6.8328
R² = 0.7264

**I** Y = 1.0255*X + 9.3662
R² = 0.8175

**Archaea**
- Halobacteria
- Other archaea

**Bacteria**
- Actinobacteria
- Bacilli
- Other bacteria

**Eukaryota**
- Plant
- Protozoa
- Fungi
- Metazoa

**Figure S17**



A

P < 2.2e-308

Domain length (aa)

DSDR>30%,
CDRN=0
(Class IV)

DSDR≤30%
CDRN≥1
(Class II)

C

Putative transcription factor SEF1 ($Asp^{523}$-$Glu^{551}$)

$Asp^{523}$

$Glu^{551}$

B

SEF1
(Putative transcription factor SEF1, P34228)
DSDR(%): 13.12

Fungal_trans (PF04082)

357    523 551 577

Disorder Tendency

Amino Acid Position

• Disordered
• Ordered

D

NADFUDRAFT_51267
(3',5'-cyclic-nucleotide phosphodiesterase, ODQ65997)
DSDR(%): 19.46

PDEase_II (PF02112)

123    203    260    420

Disorder Tendency

Amino Acid Position

• Disordered
• Ordered

E

3',5'-cyclic-nucleotide phosphodiesterase

$Pro^{203}$

$Ser^{260}$

**Figure S18**

**Figure S19**



**A**

y= x³/78750 -137x²/63000+ 253x/2100

**B**

| DSDR variation | Classification |
|----------------|----------------|
| 0 | No-variation |
| 0-1 | Low-variation |
| 1-3 | High-variation |

**C**

V1, V2, V3, V4 (V: DSDR variation) V5, V6, V7

DSDR variation dominant category:

No-variation ← d/4>0.5
Low-variation ← e/4>0.5
High-variation ← f/4>0.5
No dominant category ← others

Classification

d ← DSDR variation=0 → a
e ← 0<DSDR variation≤1 → b
f ← 1<DSDR variation≤3 → c

a/3>0.5 → No-variation
b/3>0.5 → Low-variation
c/3>0.5 → High-variation
others → No dominant category

**Figure S20**

# Figure S21



**A** Domains repeating in the same proteins

**B** Domains repeating in multiple proteins

**E**

**F**

PsaA_PsaB, Photosystem I psaA/psaB protein (PF00223)
Function: PsaA and PsaB bind the primary electron donor of photosystem I (PSI) which is a plastocyanin-ferredoxin oxidoreductase.
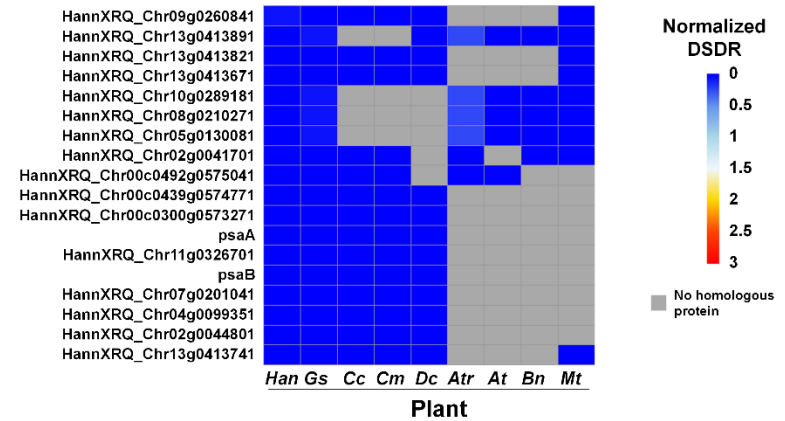
# Figure S22



**Archaea**          **Bacteria**          **Eukaryotes**