

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45

Scaling laws of graphs of 3D protein structures

Jure Pražnikar ^{*1,2}

¹ Faculty of Mathematics, Natural Sciences and Information Technologies, University of Primorska, Glagoljaška 8, Koper, Slovenia

² Department of Biochemistry, Molecular and Structural Biology, Institute Jožef Stefan, Jamova 39, Ljubljana, Slovenia

*** Corresponding author:**

E-mail: jure.praznikar@upr.si

46 **Abstract**

47 The application of graph theory in structural biology offers an alternative means of studying 3D
48 models of large macromolecules, such as proteins. However, basic structural parameters still play
49 an important role in the description of macromolecules. For example, the radius of gyration, which
50 scales with exponent ~ 0.4 , provides quantitative information about the compactness of the protein
51 structure. In this study, we combine two proven methods, the graph-theoretical and the fundamental
52 scaling laws, to study 3D protein models.

53 This study shows that the mean node degree of the protein graphs, which scales with exponent
54 0.038, is scale-invariant. In addition, proteins that differ in size have a highly similar node degree
55 distribution, which peaks at node degree 7, and additionally conforms to the same statistical
56 properties at any scale. Linear regression analysis showed that the graph parameters (radius,
57 diameter and mean eccentricity) can explain up to 90% of the total radius of gyration variance.
58 Thus, the graph parameters of radius, diameter and mean eccentricity scale with the same exponent
59 as the radius of gyration. The main advantage of graph eccentricity compared to the radius of
60 gyration is that it can be used to analyse the distribution of the central and peripheral amino
61 acids/nodes of the macromolecular structure. The central nodes are hydrophobic amino acids (Val,
62 Leu, Ile, Phe), which tend to be buried, while the peripheral nodes are more hydrophilic residues
63 (Asp, Glu, Lys). Furthermore, it has been shown that the number of central and peripheral nodes is
64 more related to the fold of the protein than to the protein length.

65

66

67

68

69

70

71

72 **Introduction**

73 Proteins, molecules that serve many critical roles in nature, consist of complex systems of amino
74 acids that have increasingly been modelled as networks over the last decade (1–3). There are many
75 ways to abstract the 3D protein structure into a graph. We can consider $C\alpha$, $C\beta$ or all heavy atoms
76 to construct an adjacency matrix. A typical method to abstract the protein model into a graph is to
77 consider $C\alpha$ atoms with 7.0 Å cut-off distance (4). We should be aware that some information is
78 lost when the 3D model is abstracted into the graph, although it still captures relevant biochemical
79 properties of the "real" 3D protein model. Once the 3D model of the protein is abstracted into a
80 graph, we have several options to analyse the 3D structure by examining different parameters of the
81 graph. Residue network models have been used to predict catalytic sites (5–8), to study protein
82 dynamics (9), to discover node-amino acids that play crucial roles in protein folding (10,11), to
83 explore allosteric pathways (12–14), and to analyse enzyme domain packing (15). In addition, graph
84 theory has also been successfully used to validate PDB entries (16), to study local errors in protein
85 models, and to discriminate decoys from the native structure (17,18). We should emphasize the
86 importance of validation and quality assessment, since only correct protein structures can answer
87 relevant biological questions (19–22).

88 However, we cannot expect that all graph parameters used and derived by mathematicians will have
89 practical implications, for example, when studying specific phenomena in structural biology.

90 However, we should try to find a connection between theoretical graph parameters and biochemical
91 phenomena that can lead to deeper insights.

92 Despite the obvious usefulness of graph theory for the analysis of protein structures, basic structural
93 parameters still serve an important role in the description of macromolecules. For example, a
94 common parameter used to describe the compactness of protein models is the radius of gyration
95 (23–25). The radius of gyration of a macromolecule describes the distribution of atoms around the
96 centre of mass. Since Flory's theory (26), the scaling law between the radius of gyration and the
97 length of the protein has been studied in detail and used to describe protein folding, and to analyse

98 the compactness of protein structures in poor or good solvent conditions. It was found that the
99 radius of gyration of globular protein structures, including monomers and oligomers, scales with
100 exponent ~ 0.4 (27). The radius of gyration can be used as a constraint when building protein models
101 or performing molecular dynamics simulations. On the other hand, when the 3D model of a protein
102 is not known, the scaling law provides qualitative information about the dimensions of the
103 macromolecule.

104 This study links two well-established approaches, the graph-theoretical and the fundamental scaling
105 laws, to study 3D protein models. In this paper, the PDB is surveyed to study scaling laws of 3D
106 protein structures using a graph-theoretical approach. This research has demonstrated that the mean
107 node degree of the protein graph is nearly scale-invariant. Additionally, the comparison of the node
108 degree distributions of proteins of different sizes exhibits marginal differences. Furthermore, this
109 study shows that the compactness of the protein, which is conventionally calculated by the radius of
110 gyration, can be estimated using graph eccentricity, which also provides insights into central
111 (buried) and peripheral (non-buried) amino acids.

112

113 **Methods**

114 **Dataset**

115 The Protein Sequence Culling Server (28) was used to obtain the PDB id list for protein structures
116 with the following characteristics: maximum mutual sequence identity of 80%, X-ray resolution
117 cut-off of 3.0 Å and minimum (maximum) chain length of 40 (10,000) residues. The PDB id list
118 was then used to retrieve 31,571 Biological Assemblies from the Protein Data Bank.

119

120 **Graph Construction**

121 From each of the 31,571 3D Biological Assemblies, the graph was constructed and analysed. $C\alpha$
122 atoms were considered as nodes, and edges between nodes were constructed if the $C\alpha$ – $C\alpha$ distance

123 between a pair of residues was less than (or equal to) 7 Å. It follows that the number of nodes was
124 equal to the number of residues (C α atoms) in the protein. Ligands, water molecules, and other
125 hetero-compounds were discarded during graph construction. Thus, if a protein has n residues, then
126 a protein graph $G = G(V, E)$ consists of a set of vertices (nodes) $V = v_1 v_2, \dots v_n$ and a set of edges
127 $E = e_1, e_2, \dots e_m$.

128

129 **Graph parameters**

130 The mean node degree (MND) of a graph G is expressed with the ratio

$$131 \quad d(G) = \frac{2e(G)}{N(G)}$$

132 where $e(G)$ represents the total number of edges in a graph G , and $N(G)$ is the number of nodes in a
133 graph G .

134 The eccentricity is a node centrality index defined as the maximum distance between a vertex to all
135 other vertices. Thus, the vertex's eccentricity is the maximal shortest path between the vertex and all
136 other vertices. Mean eccentricity is expressed as an average value of eccentricities of all vertices of
137 G . The radius of the graph is defined as a minimum eccentricity among all vertices in the graph.
138 Meanwhile, the diameter is defined as maximum eccentricity among all vertices in the graph. The
139 center of a graph or central node has eccentricity equal to the radius. A vertex is said to be a
140 peripheral node if its eccentricity is equal to the diameter.

141 Next, R script (igraph package) was used to calculate mean node degree, eccentricity, mean
142 eccentricity, radius and diameter:

```
143 MND <- mean(degree(G)) # mean node degree of graph G
```

```
144 EC <- eccentricity(G) # eccentricity of graph G
```

```
145 MEC <- mean(EC) # mean eccentricity of graph G
```

```
146 R <- min(EC) # radius of graph G, or min eccentricity of graph G
```

```
147 D <- max(EC) # diameter of graph G, or max eccentricity of graph G
```

148

149 **Radius of gyration**

150 Considering atoms as points in a three-dimensional space, the radius of gyration is defined as

$$151 \quad R_{\text{gyr}} = \sqrt{\frac{\sum m_i r_i^2}{M}}$$

152 where M is the total mass of the molecule, and m_i is the mass of the i -th atom with distance r_i from
153 the centre of mass. Radius of gyration was calculated using the *rgyr* function, which is part of the
154 Bio3D R package (29).

155

156 **Solvent accessibility of residues**

157 Secondary structure (total solvent-accessible area of proteins) was assigned according to the method
158 of Kabsch and Sander (30,31). The solvent accessible area data for protein residues were taken from
159 the work of Tien and co-workers (32).

160

161 **Table 1: Solvent accessibility of residues**

Residue	Solvent accessible area (Å²)
Alanine	121
Arginine	265
Asparagine	187
Aspartate	187
Cysteine	148
Glutamate	214
Glutamine	214
Glycine	97
Histidine	216
Isoleucine	195
Leucine	191
Lysine	230
Methionine	203
Phenylalanine	228
Proline	154
Serine	143
Threonine	163
Tryptophan	264

Tyrosine	255
Valine	165
Average	192

162

163

164 **Results and Discussion**

165

166 **Node degree-nearly scale-invariant**

167 One of the fundamental global graph parameters in graph theory is the mean node degree (MND).

168 The mean node degree shows how many edges each node has on average. In previous work,

169 Pražnikar and co-workers (33) have shown that protein models that deviate from the expected MND

170 by approximately two standard deviations or more are likely to be incorrect. Furthermore, the

171 scaling exponent calculated in the mentioned study is close to zero and indicates that the mean node

172 degree is nearly scale-invariant.

173 In this study, a large non-redundant database of biological units (31,571) was used, rather than

174 crystal asymmetric units. We can see in Figure 1A that MND is not strongly dependent on protein

175 size and that the distribution is rather narrow. Upon closer examination, however, the value

176 determined in our study (0.038) differs slightly from the value (0.024) determined in previous

177 study. The reason for the different scaling exponents is probably that the datasets are not the same.

178 An analysis performed on two large but different datasets shows that MND is nearly scale-invariant,

179 i.e., the scaling exponent is close to zero (0.024 and 0.038). Thus, MND is not strongly dependent

180 on protein length, and it can be concluded that the number of edges in the protein graph is linearly

181 related to the number of nodes (amino acids).

182

183 **Fig 1.** (A) Scaling exponent of mean node degree of protein graphs versus the number of residues.

184 (B) Probability of node degree of protein graphs for three size bins. The first size bin (black line)

185 encompasses protein structures with length between 100 and 200 residues, the second size bin (blue
186 line) encompasses protein structures with length between 500 and 600 residues, and the third size
187 bin encompasses structures with length between 900 and 1000 residues.

188

189 We could expect that larger proteins would have a higher average node degree because of a higher
190 number of core residues and a relatively lower number of surface residues, which are supposed to
191 have lower numbers of edges. Thus, to further analyse the node degree of protein nodes-residues,
192 we calculated the node degree distribution for three different size bins. The first size bin contains all
193 protein structures from the database, which have lengths between 100 and 200 residues, in the
194 second size bin are proteins with lengths between 500 and 600 residues, and in the third size bin are
195 structures with lengths between 900 and 1,000 residues. Figure 1B shows that all three distributions
196 are very similar and that there is a peak at 7 node degrees. The comparison of all three peak values
197 shows that the first size bin, which contains the smallest proteins, has the highest probability density
198 value. The lowest probability density value at 7 node degree is seen for the third size bin, which
199 contains the largest proteins among all three selected size bins. A closer look at the left (degree 2)
200 and right (degree 14) tail of the distributions shows high similarity for all three distributions. The
201 visual comparison shows the most significant differences on the left and right sides of the peak. The
202 differences can be observed at values of approximately 5 to 9 node degrees. It can be seen that the
203 first size bin has a higher probability at 3 to 6 node degrees as compared to size bins two and three.
204 The order is somehow reversed on the right side of the displayed distribution. For node degrees 8, 9
205 and 10, size bin three exhibits higher values compared with size bins two and three.

206 This analysis shows that despite the different sizes of the proteins, they have a very similar node
207 degree distribution, which peaks at node degree 7. A simple way to explain the presented results is
208 that buried residues in small or large proteins form approximately the same number of links. This is
209 a direct consequence of the fact that the amino acids are physical objects and cannot be arbitrarily

210 close to each other. The marginal difference in node degree distributions, a slight shift to higher
211 node degrees, explains the low positive scaling exponent (0.038), which is nearly scale-invariant.
212

213 **Protein graph eccentricity: an alternative method for analysis of** 214 **radius of gyration**

215 It is easy to ask a question: radius of gyration and radius as a graph parameter have a common
216 name, but do they follow the same power law? To answer this question, linear regression analysis
217 and scaling exponent were calculated for three graph parameters: radius, diameter and mean
218 eccentricity. The radius-graph parameter is defined as minimum eccentricity, whereas the
219 eccentricity of the graph is defined as the maximum distance between one node and all other nodes.
220 Notice that the diameter is defined as maximum eccentricity.

221 Figure 2A shows the scaling exponent of a radius of gyration for 31,571 selected protein structures.

222 The non-linear fitting function can be written as

$$223 R_{\text{gyr}}=R_0N^{\nu},$$

224 where R_{gyr} is the radius of gyration, R_0 is the pre-factor and ν is a scaling exponent. The pre-factor
225 R_0 can be obtained experimentally and used as a restrained value during non-linear fitting (34–36).

226 Thus, when restrained fitting was performed, the pre-factor ($R_0=2 \text{ \AA}$) was fixed. We can see that in
227 the case of restrained fitting, the scaling exponent is 0.405, which is consistent with other studies.

228 When fitted without restraint, the exponent is lower (0.351). Both values are within the range
229 reported by other studies (25,27,34,37–39).

230

231 **Fig 2.** Log-log plots of the radius of gyration (A), radius (B), diameter (C), and mean eccentricity
232 (D) versus the number of residues. The solid red line is generated by fitting without restraint; the
233 blue line is produced by restrained fitting. The legend shows unrestrained scaling factors in red. The
234 restrained scaling factors and corresponding pre-factors are in blue.

235

236 Figure 2B, C, and D show radius, diameter, and mean eccentricity plotted against protein length.
237 Similar to the analysis of the radius of gyration, the power exponent was fitted with and without
238 restraint. The pre-factor for restrained fitting was derived from linear regression analysis, as shown
239 in Figure 3. The linear regression analysis between the radius of gyration and graph parameters
240 reveals that R^2 is close to 0.90 for all three cases (Fig. 3). The highest R^2 (0.91) is observed between
241 mean eccentricity and radius of gyration (Fig. 3C). The reason for this is probably that the values of
242 radius and diameter are discrete, while mean eccentricity values are not discrete. For example, the
243 radius can be 7 or 8, but cannot be a real number between 7 and 8. Mean eccentricity is just a mean
244 value of all shortest paths to any nodes. It is seen that the distribution of mean eccentricity is
245 smoother in comparison to the discrete values of radius and diameter on the y-axis.

246

247 **Fig 3.** Scatter plot of graph parameters: radius (A), diameter (B), and mean eccentricity (C) versus
248 radius of gyration. The legend shows the regression coefficient, R-squared value, and p-value.

249

250 If we use pre-factor R_0 of a radius of gyration, which was obtained from experimental data, then we
251 can calculate the pre-factors for radius, diameter, and mean eccentricity using the slope k from a
252 regression analysis. The steepness of the linear regression model between R_{gyr} and radius was 0.42
253 \AA^{-1} , 0.77\AA^{-1} between R_{gyr} and diameter, and 0.59\AA^{-1} between R_{gyr} and mean eccentricity (see Fig.
254 3A, B and C). Using pre-factor R_0 and the steepness of linear fit k , the pre-factors for radius,
255 diameter and mean eccentricity can be calculated using the next expression:

$$256 \quad R_x = R_0 k, \quad (1)$$

257 where R_x is a new calculated pre-factor, R_0 is the pre-factor of radius of gyration and k is the
258 steepness of the linear fit. In Figure 2B, C and D are shown calculated pre-factors for radius (2.0\AA
259 $0.42 \text{\AA}^{-1} = 0.84$), diameter ($2.0 \text{\AA} 0.77 \text{\AA}^{-1} = 1.54$) and mean eccentricity ($2.0 \text{\AA} 0.59 \text{\AA}^{-1} = 1.18$).

260 We can see that the restrained scaling exponent is higher than the non-restrained scaling exponent
261 for all three cases. Furthermore, it is observed that restrained graph parameters all have very similar

262 scaling exponents (~ 0.395) which are very close to the scaling exponent of the radius of gyration
263 (0.405).
264 Thus, this study shows that the radius of gyration, which is calculated from the atomic coordinates
265 and radius of graph follow the same scaling exponent. From this we can conclude that when
266 analysing 3D models of macromolecules using a graph-theoretical theory approach, the eccentricity
267 of the graph can be used to estimate the radius of gyration. Thus, graph parameter eccentricity
268 allows us to investigate whether the 3D protein model deviates from the expected value of the
269 radius of gyration and to obtain information about the compactness of the structure. Furthermore,
270 when a scientist builds a model, and the model is still in an early phase, e.g., as an alanine chain,
271 then the $C\alpha$ only model, which can be represented as a graph, contains enough information to
272 estimate the radius of gyration.

273

274 **Central and peripheral nodes**

275 In the previous section, the relation between the radius of gyration and graph eccentricity was
276 introduced. When exploring graph eccentricity, it is ubiquitous to examine which nodes-amino
277 acids are central (close to every other node) and peripheral (distant from every other node). This
278 kind of analysis has some common points with analysis of solvent-exposed residues, which is
279 directly related to the arrangement of residues in 3D space. Buried residues constitute the core of
280 the protein; meanwhile, residues exposed to the solvent represent the outer part of the protein in 3D
281 space (40). The molecular mass of a protein is related with the total solvent exposed surface using
282 the next expression:

$$283 \quad A_{monomer} = 4.44M^{0.770}, \quad (2)$$

284 where M is the molecular mass of the protein (41). Similarly, we can introduce the relation between
285 protein length and total solvent-exposed surface. Given that the average amino acid has a total
286 solvent exposed area of 192\AA^2 (32), we can use this value as a restraint during data fitting. Figure 4
287 shows the protein length against the total exposed area. The scaling exponent of the fitted curve is

288 0.772, almost the same as the exponent in equation 2. This result was somehow expected because
289 molecular mass correlates with sequence length.

290

291 **Fig 4.** Dependence of total solvent accessible area on the number of residues in proteins. The legend
292 shows the scaling exponent and pre-factor (A_0), which was used during restrained fitting.

293

294 However, a closer examination of the relationship between the number of central (peripheral) nodes
295 and protein length demonstrated that the numbers of central and peripheral residues are not related
296 to the protein size (Figure 5A, B). To further support this finding, the comparison between the
297 radius of gyration, mean eccentricity, and numbers of central and peripheral residues for nine
298 different size bins was made. We can see (Figure 6) that the mean eccentricity and radius of
299 gyration increase according to the length of the protein. Note that the scaling exponent for both
300 mentioned parameters is approximately 0.4 (Figure 2). On the other hand, central and peripheral
301 box plots do not show such a positive trend (Figure 6C, D). It follows that the numbers of
302 peripheral and central residues are not related to the protein size.

303

304 **Fig 5.** Distribution of central/peripheral residues against the size of proteins (number of residues).

305

306 **Fig 6.** Boxplot of the radius of gyration (A), mean eccentricity (B), number of peripheral nodes (C),
307 and number of central nodes (D) for ten different size bins.

308

309 Next, we examine the case with different central to peripheral node ratio, which is shown as a graph
310 and ribbon representation of two protein structures (Figure 7 A, B, D, and C). The PDB id: 1pq7
311 structure has a low number of central nodes (4), but a high number of peripheral nodes (47): see
312 Figures 7A and B. The situation is somehow reversed for the PDB id: 3wvj structure, which has a
313 higher number of central nodes (34) and lower number of peripheral nodes (2): see Figure 7D, E.

314 This case shows that the numbers of central and peripheral nodes depend on the protein fold rather
315 than on the length of the protein chain. Furthermore, we can draw parallels between *almost-*
316 *peripheral* (42), *self-centred* (43), and protein graphs. Figure 7E shows the wheel, a graph that is
317 *almost-peripheral*, containing only one central node and 6 peripheral nodes. Meanwhile, the graph
318 shown in Figure 7F is an almost *self-centred* graph that contains 5 central nodes and 2 peripheral
319 nodes. We could say that the graph abstracted from the PDB id: 3wvj structure is centred, while
320 peripheral nodes-amino acids dominate the PDB id: 1pq7 structure.

321

322 **Fig 7.** Examples of central and peripheral protein graphs. (A) Graph and (B) ribbon presentation of
323 Biological Assembly 1 of PDB entry 1pq7. (D) Graph and (C) ribbon presentation of Biological
324 Assembly 1 of PDB entry 3wvj. (E) Presents an almost peripheral graph, while (F) represents an
325 almost self-centred graph.

326

327 Further analysis shows the frequency distribution of central and peripheral amino acids (Figure 8),
328 which are subdivided into three groups: (i) charged, (ii) polar, and (iii) non-polar side chains. It can
329 be seen that amino acids histidine, cysteine, methionine, and tryptophan have the lowest probability
330 values; it follows that they are neither central nor peripheral nodes. Further, we can see that central
331 nodes are hydrophobic amino acids (Val, Leu, Ile, Phe), which tend to be buried, while peripheral
332 nodes are more likely hydrophilic residues (Asp, Glu, Lys) which form hydrogen bonds with
333 solvent.

334

335 **Fig 8.** Probability density of central and peripheral nodes for different types of amino acids.

336

337 Thus, when the 3D protein structure is analysed using graph theory, the eccentricity of the graph
338 can be very strong for evaluating the radius of gyration when studying central nodes, which tend to
339 be hydrophobic amino acids, and peripheral nodes, which tend to be hydrophilic amino acids. It is

340 remarkable that this graph parameter (eccentricity) allowed us to study the compactness (R_{gyr}) of the
341 protein and the arrangement of the residues (buried/not buried), making it versatile and useful for
342 the analysis of 3D macromolecules.

343

344 **Conclusion**

345 This study showed that the mean node degree of the protein graph is nearly scale-invariant. In other
346 words, a small scaling exponent (0.038) indicates that the mean node degree is scale-free. Scale
347 invariance was further supported by the analysis of node degree distribution, which showed very
348 similar node degree distributions for proteins that differ according to size. This scale invariably
349 offers a valuable tool for validating structures by simply counting the number of edges.

350 Furthermore, an additional comparison between the expected node degree and node degree of a
351 candidate could be used to explore and interpret large deviations. For example, intrinsically
352 disordered proteins are expected to have a considerably lower mean node degree than globular
353 proteins of the same size.

354 The comparison between the mean eccentricity of the graph and radius of gyration revealed a high
355 R^2 . In other words, the mean eccentricity and radius of gyration follow the same scaling exponent
356 (~ 0.4). The eccentricity of the graph, in addition to the estimation of the radius of gyration, also
357 allows us to study the distribution of central (buried) and peripheral amino acids (non-buried). We
358 should be aware that the mean eccentricity alone (or radius of gyration), which is used as a
359 constraint when running molecular dynamics simulations or manually building a model, does not
360 provide the correctness of the protein model. It is also crucial to determine how the amino acids are
361 distributed in real space, and this can be elucidated by studying peripheral and central nodes. Thus,
362 a single graph parameter (eccentricity) can be used to control the compactness of the
363 macromolecule and the distribution of amino acids in 3D space, which makes it a valuable tool for
364 analysing protein models.

365

366 **Acknowledgement**

367 This work was supported by Structural Biology grant P1-0048 and Infrastructure programme grant
368 I0-0035-2790, provided by the Slovenian Research Agency.

369

370 **References**

- 371 1. Estrada E. Universality in protein residue networks. *Biophys J* [Internet]. 2010;98(5):890–
372 900. Available from: <http://dx.doi.org/10.1016/j.bpj.2009.11.017>
- 373 2. Greene LH. Protein structure networks. *Brief Funct Genomics*. 2012;11(6):469–78.
- 374 3. Vishveshwara S, Brinda K V., Kannan N. Protein Structure: Insights From Graph Theory. *J*
375 *Theor Comput Chem*. 2002;01(01):187–211.
- 376 4. da Silveira CH, Pires DE V., Minardi RC, Ribeiro C, Veloso CJM, Lopes JCD, et al. Protein
377 cutoff scanning: A comparative analysis of cutoff dependent and cutoff free methods for
378 prospecting contacts in proteins. *Proteins Struct Funct Bioinforma* [Internet]. 2009 Feb 15
379 [cited 2017 Mar 15];74(3):727–43. Available from: <http://doi.wiley.com/10.1002/prot.22187>
- 380 5. del Sol A, Fujihashi H, Amoros D, Nussinov R. Residue centrality, functionally important
381 residues, and active site shape: Analysis of enzyme and non-enzyme families. *Protein Sci*
382 [Internet]. 2006 Sep 28;15(9):2120–8. Available from:
383 <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2242611/>
- 384 6. Thibert B, Bredesen DE, del Rio G. Improved prediction of critical residues for protein
385 function based on network and phylogenetic analyses. *BMC Bioinformatics* [Internet]. 2005
386 Aug 26;6:213. Available from: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1208857/>
- 387 7. Wangikar PP, Tendulkar A V, Ramya S, Mali DN, Sarawagi S. Functional Sites in Protein
388 Families Uncovered via an Objective and Automated Graph Theoretic Approach. *J Mol Biol*
389 [Internet]. 2003 Feb 21;326(3):955–78. Available from:
390 <http://www.sciencedirect.com/science/article/pii/S0022283602013840>

- 391 8. Chea E, Livesay DR. How accurate and statistically robust are catalytic site predictions based
392 on closeness centrality? *BMC Bioinformatics*. 2007;8:153.
- 393 9. Atilgan AR, Akan P, Baysal C. Small-world communication of residues and significance for
394 protein dynamics. *Biophys J [Internet]*. 2004;86(1 Pt 1):85–91. Available from:
395 [http://dx.doi.org/10.1016/S0006-3495\(04\)74086-2](http://dx.doi.org/10.1016/S0006-3495(04)74086-2)
- 396 10. Vendruscolo M, Dokholyan N V., Paci E, Karplus M. Small-world view of the amino acids
397 that play a key role in protein folding. *Phys Rev E Stat Nonlin Soft Matter Phys*.
398 2002;65(6):4.
- 399 11. Vendruscolo M, Paci E, Dobson CM, Karplus M. Three key residues form a critical contact
400 network in a protein folding transition state. *Nature [Internet]*. 2001 Feb 1;409(6820):641–5.
401 Available from: <http://dx.doi.org/10.1038/35054591>
- 402 12. Negre CFA, Morzan UN, Hendrickson HP, Pal R, Lisi GP, Loria JP, et al. Eigenvector
403 centrality for characterization of protein allosteric pathways. *Proc Natl Acad Sci [Internet]*.
404 2018 Dec 26;115(52):E12201 LP-E12208. Available from:
405 <http://www.pnas.org/content/115/52/E12201.abstract>
- 406 13. Daily MD, Upadhyaya TJ, Gray JJ. Contact rearrangements form coupled networks from
407 local motions in allosteric proteins. *Proteins*. 2008 Apr;71(1):455–66.
- 408 14. Daily MD, Gray JJ. Allosteric communication occurs via networks of tertiary and quaternary
409 motions in proteins. *PLoS Comput Biol*. 2009;5(2).
- 410 15. Pintar S, Borišek J, Usenik A, Perdih A, Turk D. Domain sliding of two *Staphylococcus*
411 *aureus* N-acetylglucosaminidases enables their substrate-binding prior to its catalysis.
412 *Commun Biol*. 2020;3(1):1–9.
- 413 16. Pražnikar J, Tomić M, Turk D. Validation and quality assessment of macromolecular
414 structures using complex network analysis. *Sci Rep*. 2019 Dec 1;9(1).
- 415 17. Chatterjee S, Bhattacharyya M, Vishveshwara S. Network properties of protein-decoy
416 structures. *J Biomol Struct Dyn [Internet]*. 2012;29(6):1110–26. Available from:

- 417 <https://doi.org/10.1080/07391102.2011.672625>
- 418 18. Chatterjee S, Ghosh S, Vishveshwara S. Network properties of decoys and CASP predicted
419 models: A comparison with native protein structures. *Mol Biosyst.* 2013;9(7):1774–88.
- 420 19. Kleywegt GJ. Validation of protein crystal structures. *Acta Crystallogr Sect D Biol*
421 *Crystallogr.* 2000;56(3):249–65.
- 422 20. Read RJ, Adams PD, Arendall WB, Brunger AT, Emsley P, Joosten RP, et al. A new
423 generation of crystallographic validation tools for the Protein Data Bank. *Structure.*
424 2011;19(10):1395–412.
- 425 21. Wlodawer A, Minor W, Dauter Z, Jaskolski M. Protein crystallography for non-
426 crystallographers, or how to get the best (but not more) from published macromolecular
427 structures. *FEBS J.* 2008;275(1):1–21.
- 428 22. Benkert P, Biasini M, Schwede T. Toward the estimation of the absolute quality of individual
429 protein structure models. *Bioinformatics.* 2011;27(3):343–50.
- 430 23. Lobanov MY, Bogatyreva NS, Galzitskaya O V. Radius of gyration as an indicator of protein
431 structure compactness. *Mol Biol.* 2008;42(4):623–8.
- 432 24. Enright MB, Leitner DM. Mass fractal dimension and the compactness of proteins. *Phys Rev*
433 *E - Stat Nonlinear, Soft Matter Phys.* 2005;71(1):1–9.
- 434 25. Hong L, Lei J. Scaling law for the radius of gyration of proteins and its dependence on
435 hydrophobicity. *J Polym Sci Part B Polym Phys [Internet].* 2009;47(2):207–14. Available
436 from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/polb.21634>
- 437 26. Flory PJ. The Configuration of Real Polymer Chains. *J Chem Phys [Internet].*
438 1949;17(3):303–10. Available from: <https://doi.org/10.1063/1.1747243>
- 439 27. Tanner JJ. Empirical power laws for the radii of gyration of protein oligomers. *Acta*
440 *Crystallogr Sect D Struct Biol.* 2016;72:1119–29.
- 441 28. Wang G, Dunbrack Jr RL. PISCES: a protein sequence culling server. *Bioinformatics*
442 [Internet]. 2003 Aug 12;19(12):1589–91. Available from:

- 443 <https://doi.org/10.1093/bioinformatics/btg224>
- 444 29. Grant BJ, Rodrigues APC, ElSawy KM, McCammon JA, Caves LSD. Bio3d: An R package
445 for the comparative analysis of protein structures. *Bioinformatics*. 2006;22(21):2695–6.
- 446 30. Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of
447 hydrogen-bonded and geometrical features. *Biopolymers*. 1983 Dec;22(12):2577–637.
- 448 31. Touw WG, Baakman C, Black J, Te Beek TAH, Krieger E, Joosten RP, et al. A series of
449 PDB-related databanks for everyday needs. *Nucleic Acids Res*. 2015;43(D1):D364–8.
- 450 32. Tien MZ, Meyer AG, Sydykova DK, Spielman SJ, Wilke CO. Maximum allowed solvent
451 accessibilities of residues in proteins. *PLoS One*. 2013;8(11).
- 452 33. Pražnikar J, Tomić M, Turk D. Validation and quality assessment of macromolecular
453 structures using complex network analysis. *Sci Rep [Internet]*. 2019;9(1):1678. Available
454 from: <https://doi.org/10.1038/s41598-019-38658-9>
- 455 34. Hofmann H, Soranno A, Borgia A, Gast K, Nettels D, Schuler B. Polymer scaling laws of
456 unfolded and intrinsically disordered proteins quantified with single-molecule spectroscopy.
457 *Proc Natl Acad Sci U S A*. 2012;109(40):16155–60.
- 458 35. Kohn JE, Millett IS, Jacob J, Zagrovic B, Dillon TM, Cingel N, et al. Random-coil behavior
459 and the dimensions of chemically unfolded proteins. *Proc Natl Acad Sci U S A*. 2004
460 Aug;101(34):12491–6.
- 461 36. Wilkins DK, Grimshaw SB, Receveur V, Dobson CM, Jones JA, Smith LJ. Hydrodynamic
462 Radii of Native and Denatured Proteins Measured by Pulse Field Gradient NMR Techniques.
463 *Biochemistry [Internet]*. 1999 Dec 1;38(50):16424–31. Available from:
464 <https://doi.org/10.1021/bi991765q>
- 465 37. Dima RI, Thirumalai D. Asymmetry in the Shapes of Folded and Denatured States of
466 Proteins. *J Phys Chem B [Internet]*. 2004 May 1;108(21):6564–70. Available from:
467 <https://doi.org/10.1021/jp037128y>
- 468 38. Gong H, Fleming PJ, Rose GD. Building native protein conformation from highly

- 469 approximate backbone torsion angles. Proc Natl Acad Sci U S A [Internet]. 2005 Nov
470 8;102(45):16227 LP – 16232. Available from:
471 <http://www.pnas.org/content/102/45/16227.abstract>
- 472 39. Hinsen K, Hu S, Kneller GR, Niemi AJ. A comparison of reduced coordinate sets for
473 describing protein structure. J Chem Phys [Internet]. 2013;139(12):124115. Available from:
474 <https://doi.org/10.1063/1.4821598>
- 475 40. Miller S, Janin J, Lesk AM, Chothia C. Interior and surface of monomeric proteins. J Mol
476 Biol. 1987 Aug;196(3):641–56.
- 477 41. Marsh JA. Buried and accessible surface area control intrinsic protein flexibility. J Mol Biol
478 [Internet]. 2013;425(17):3250–63. Available from:
479 <http://dx.doi.org/10.1016/j.jmb.2013.06.019>
- 480 42. Klavzar S, Narayankar KP, Walikar HB, Lokesh SB. Almost-Peripheral Graphs. Taiwan J
481 Math. 2014;18:463–71.
- 482 43. Klavžar S, Narayankar KP, Walikar HB. Almost self-centered graphs. Acta Math Sin Engl
483 Ser [Internet]. 2011;27(12):2343–50. Available from: [https://doi.org/10.1007/s10114-011-](https://doi.org/10.1007/s10114-011-9628-3)
484 9628-3

485

486 **Data availability**

487 The data is freely available on github at <https://github.com/jure-praznikar/Scaling-laws>.

488

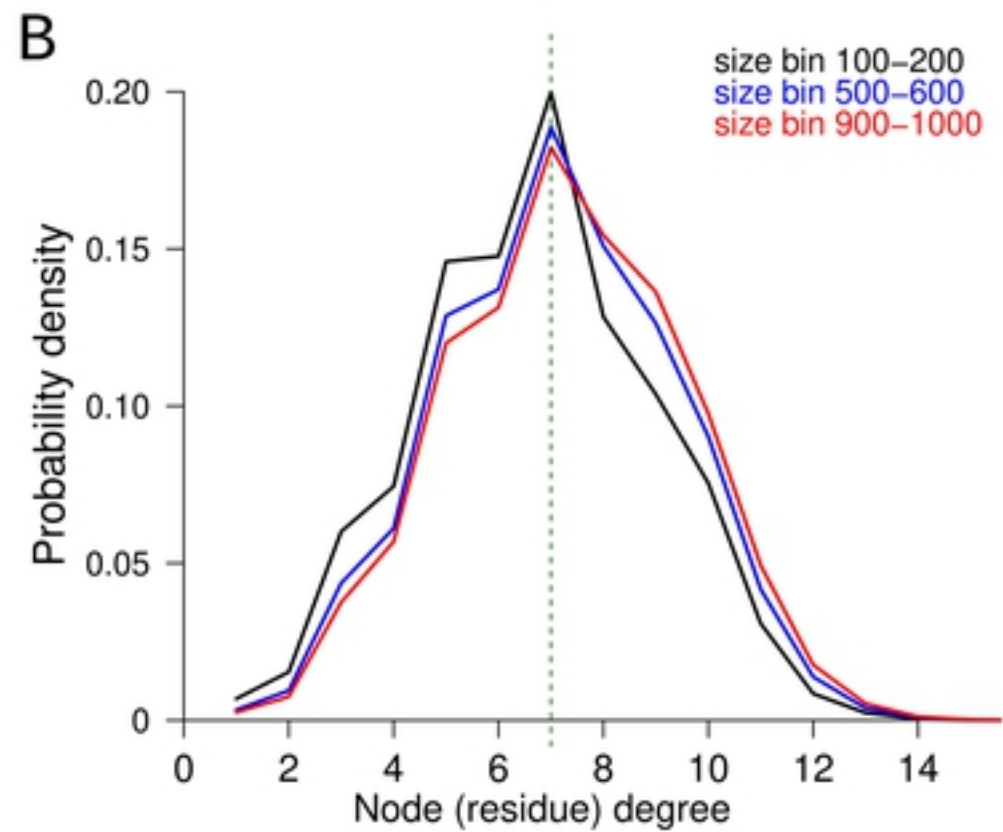
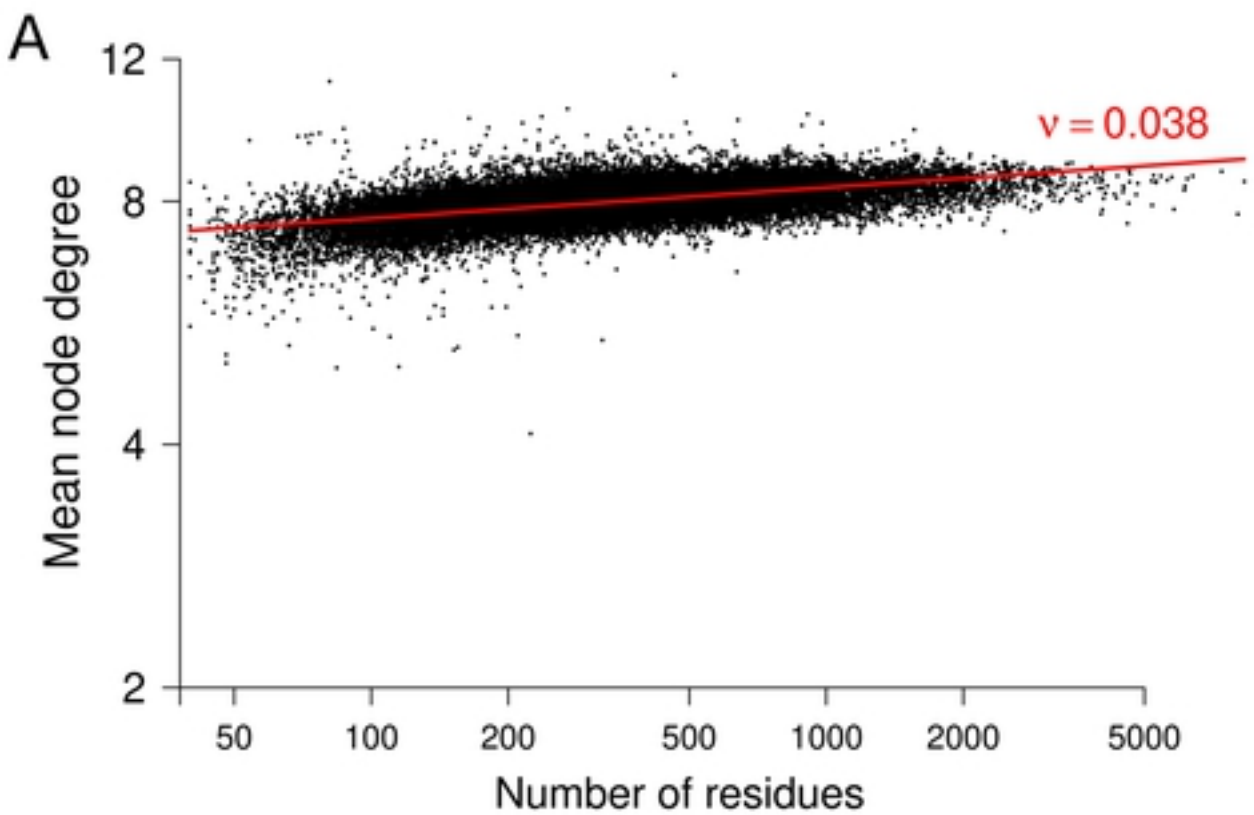


Figure 1

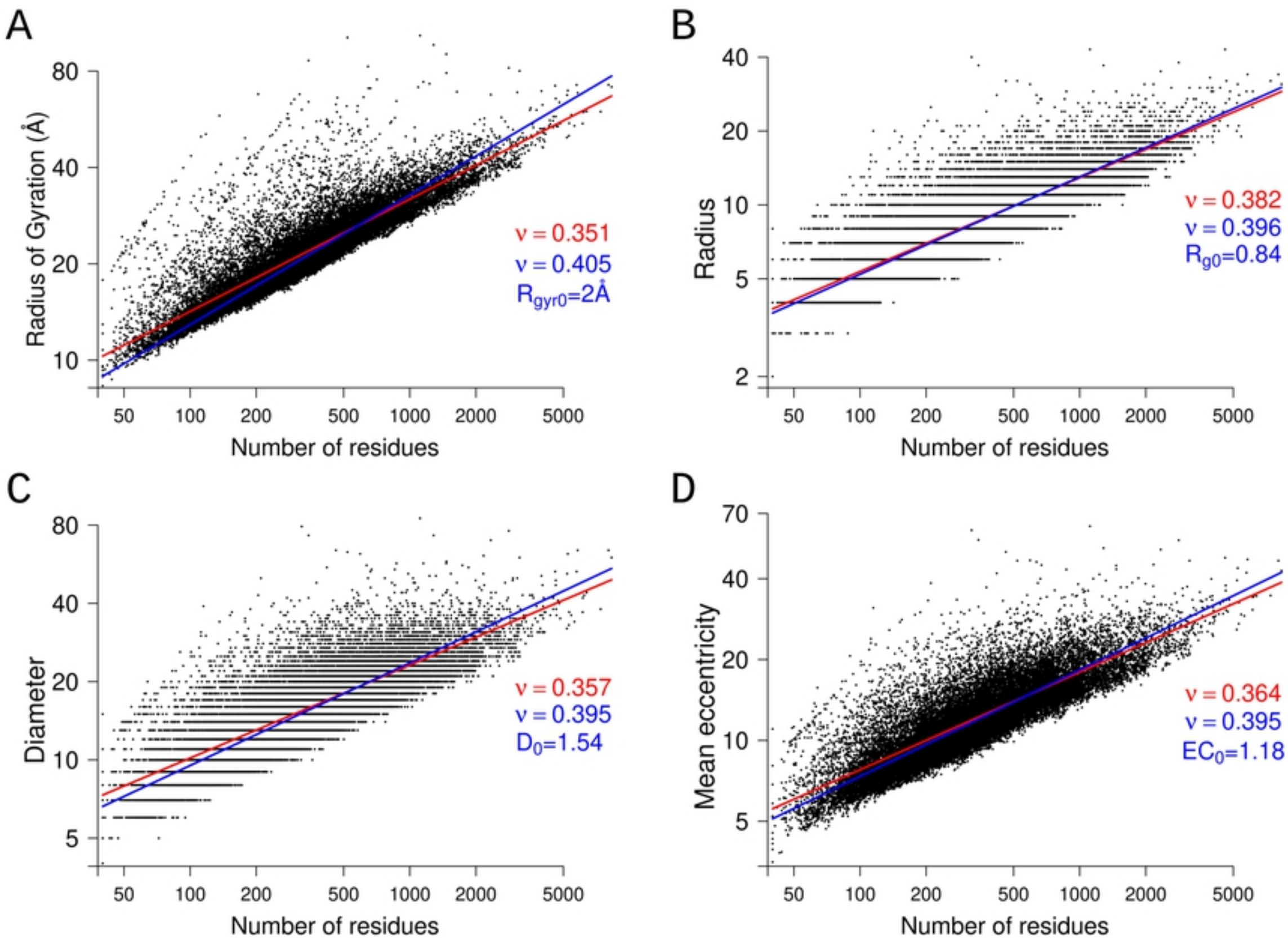


Figure 2

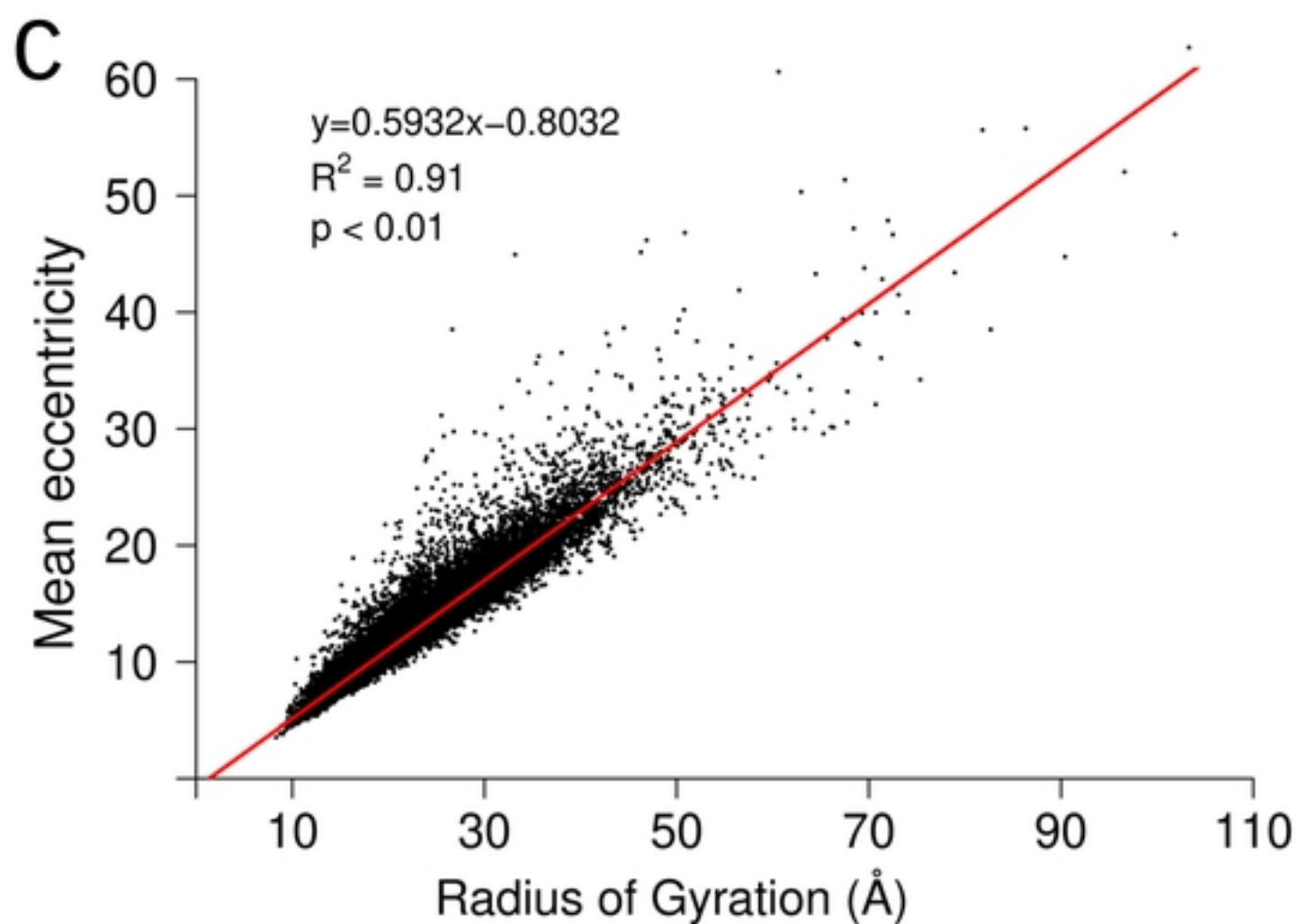
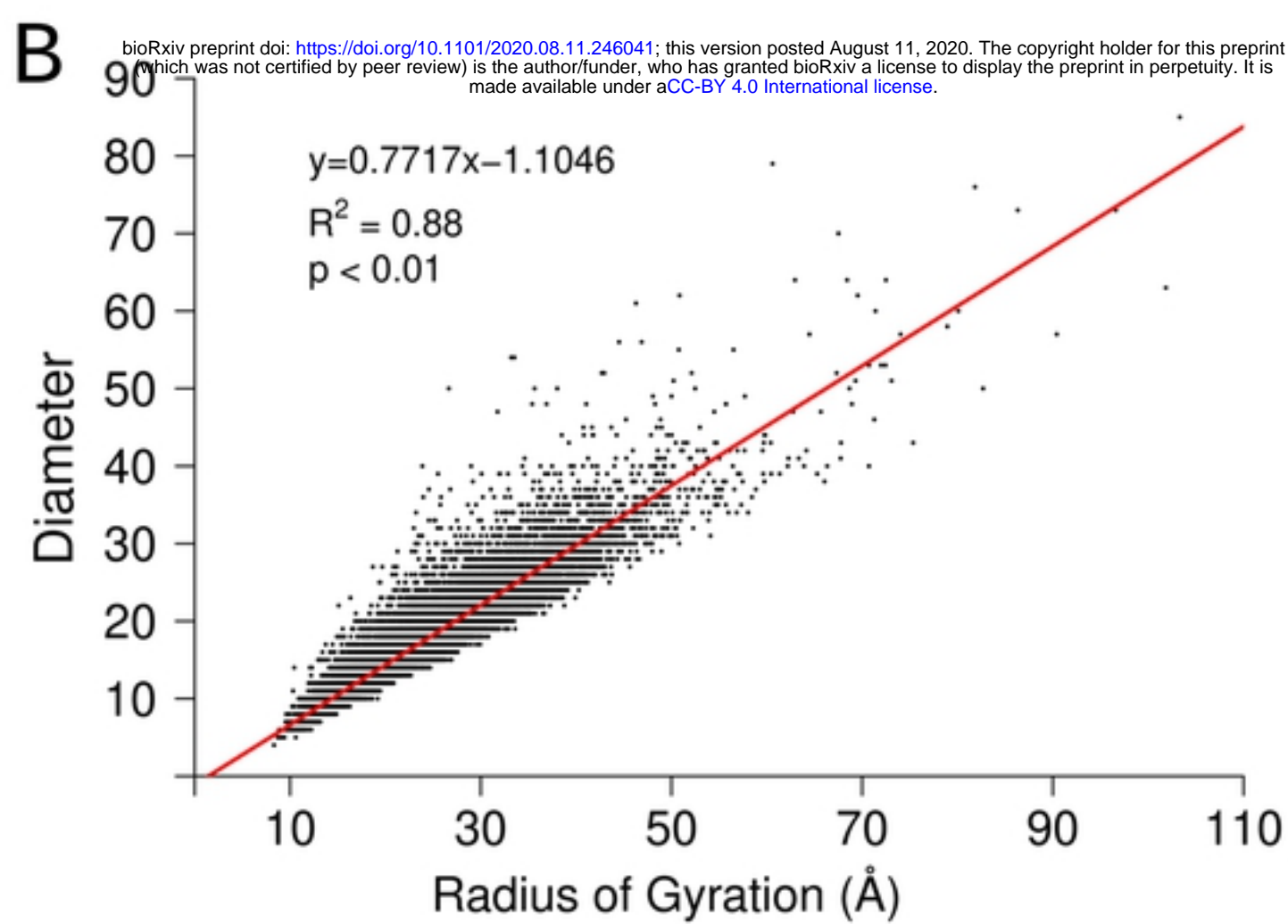
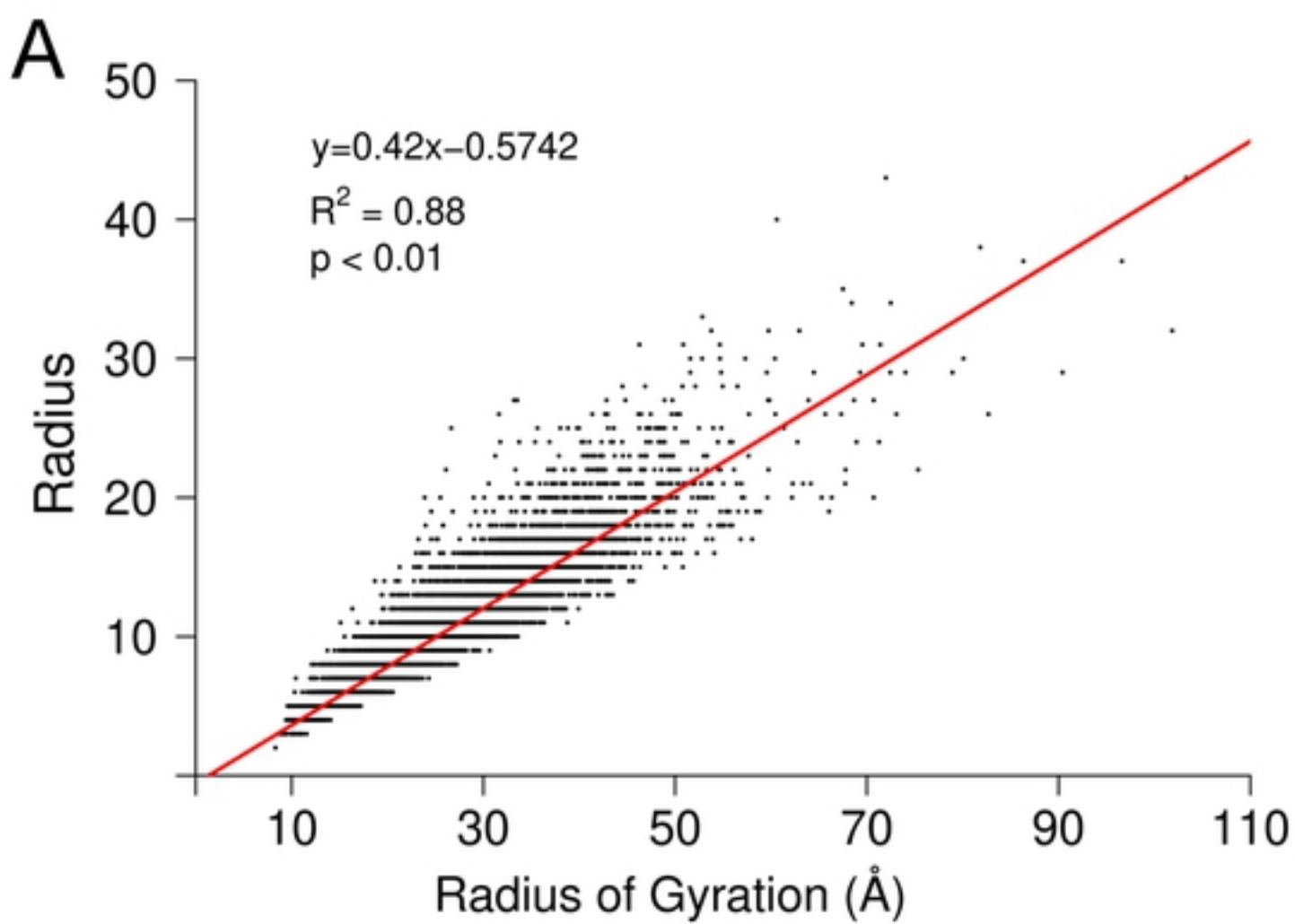


Figure 3

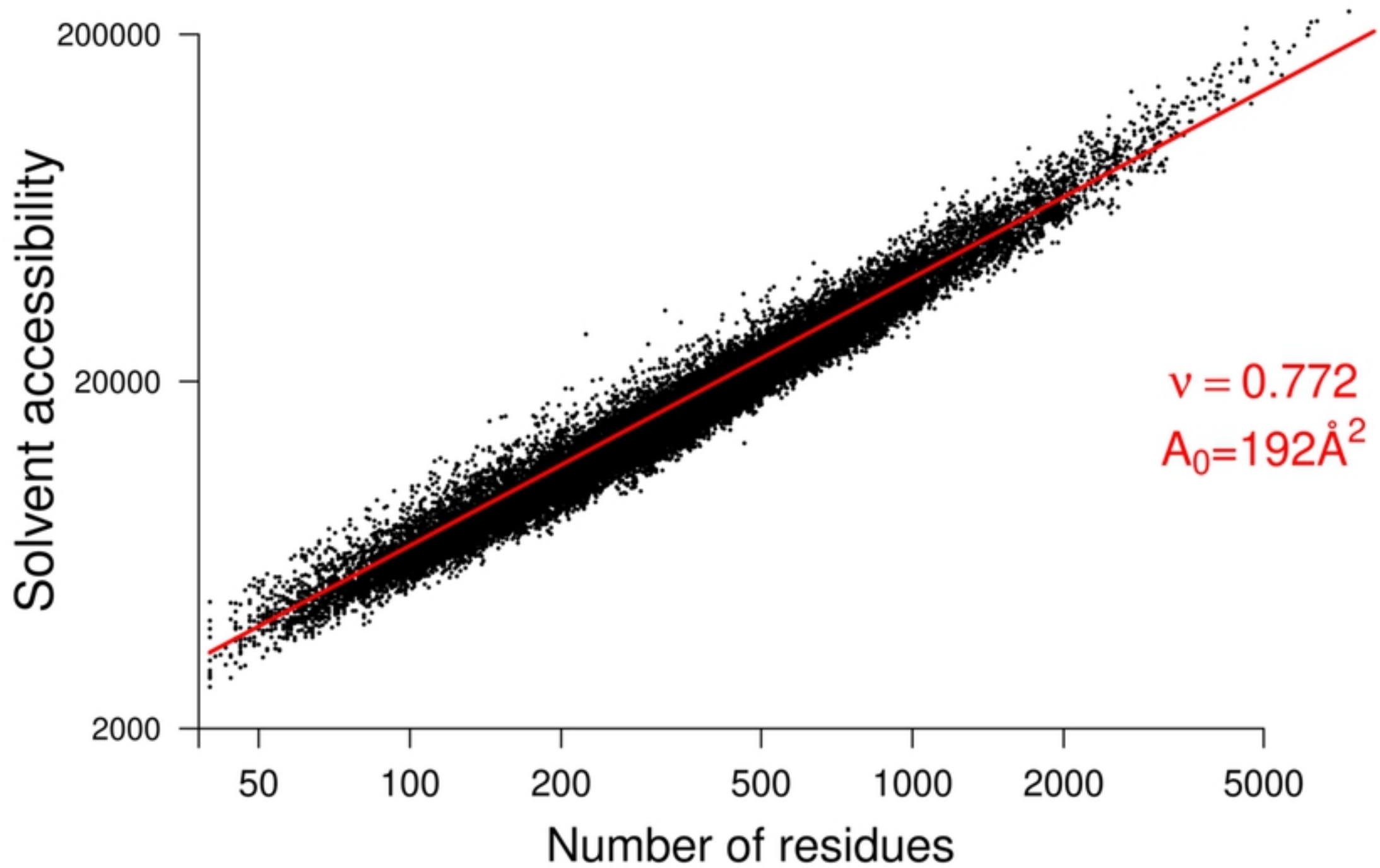


Figure 4

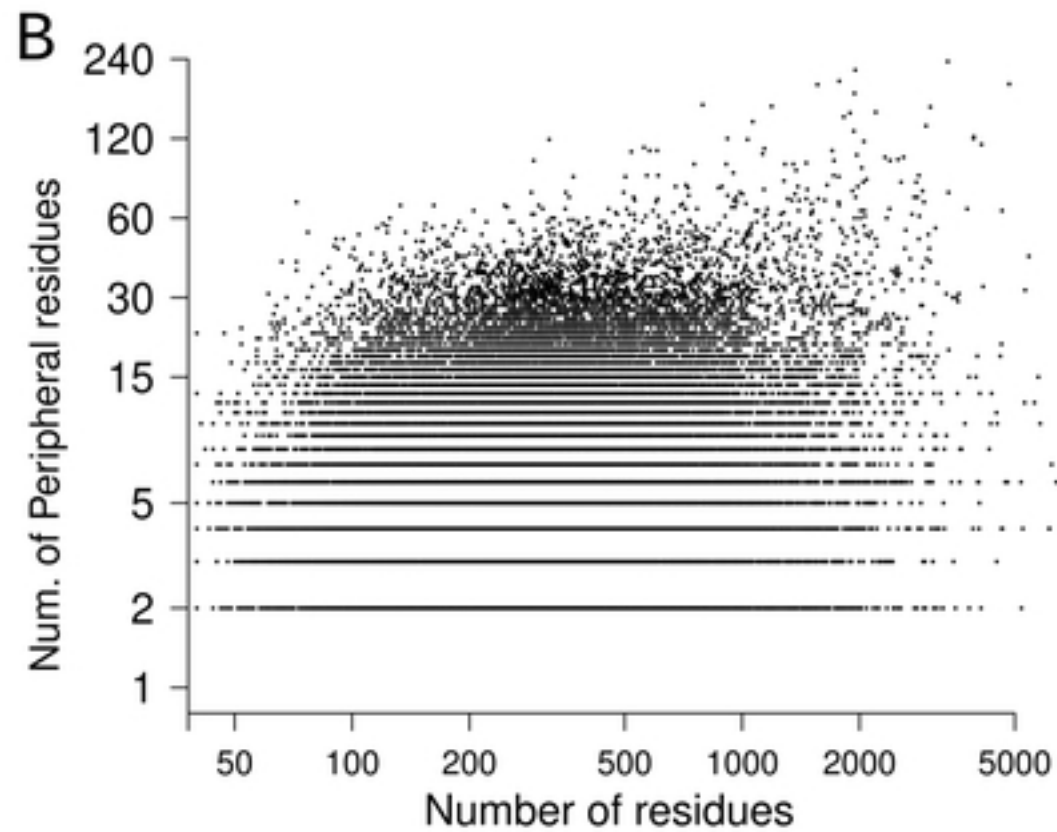
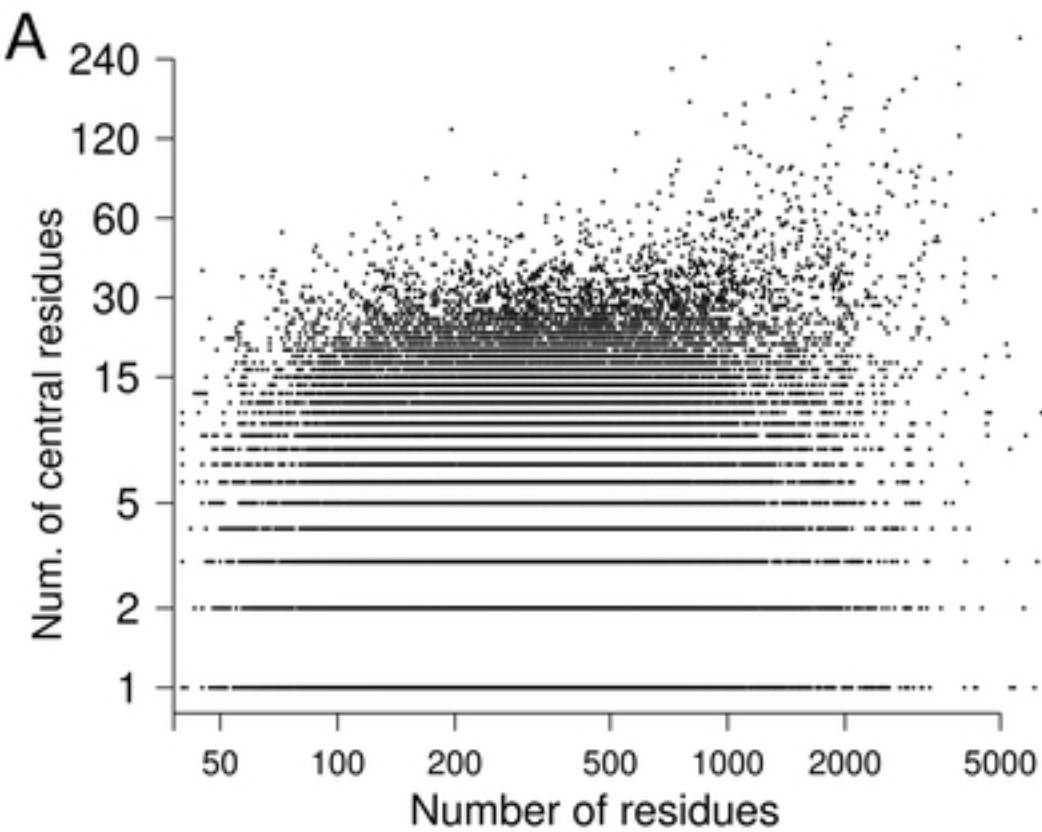
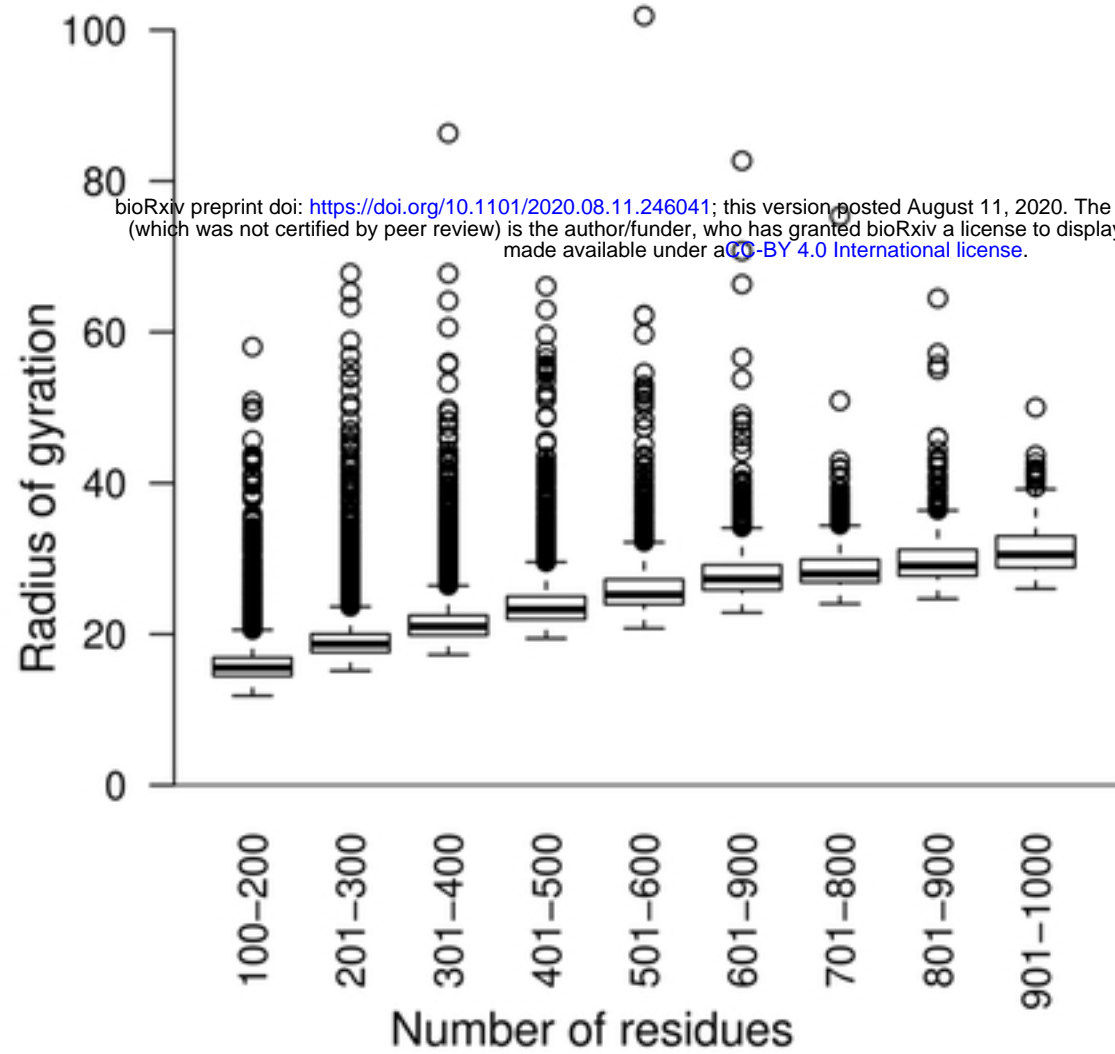
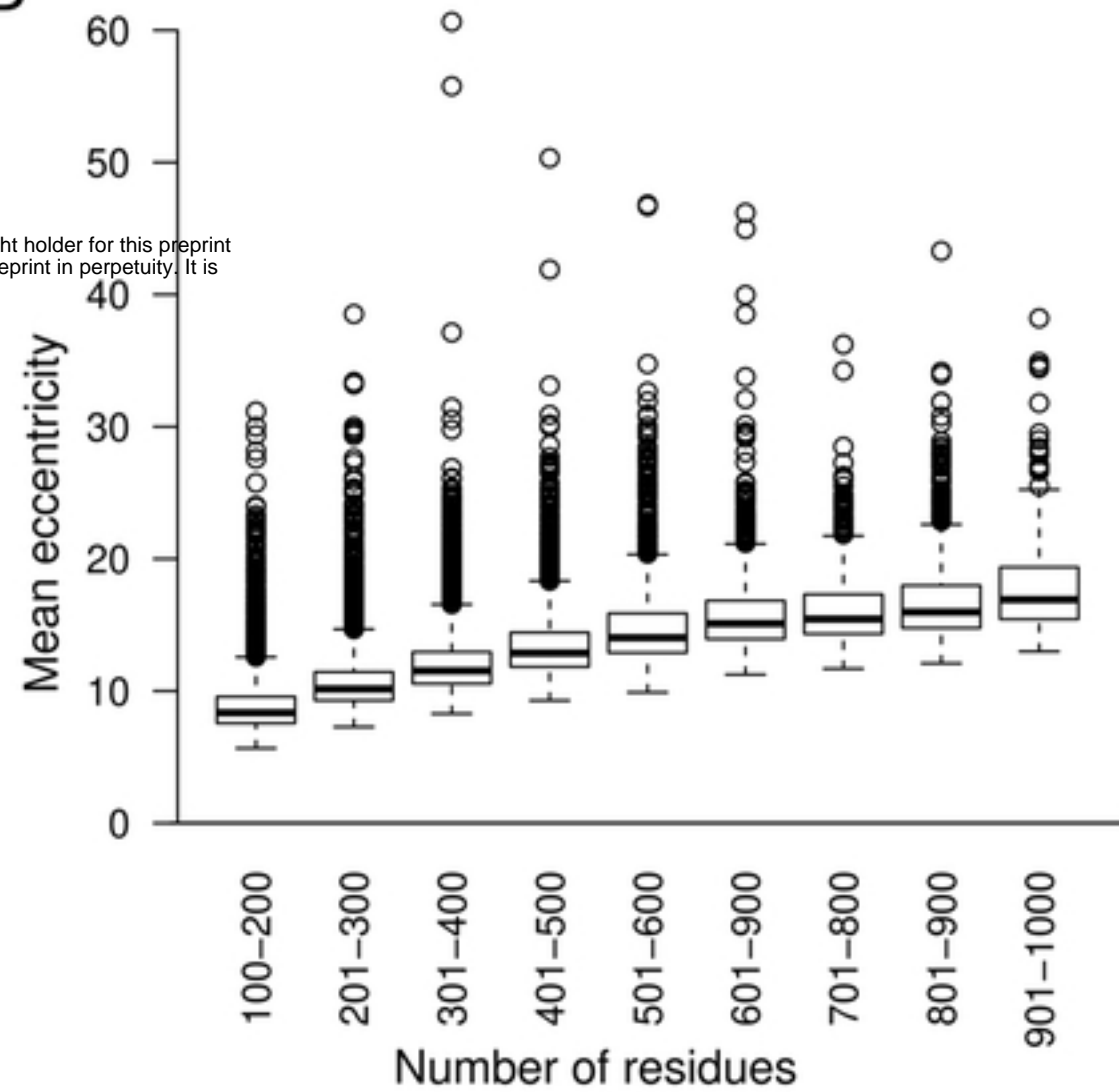


Figure 5

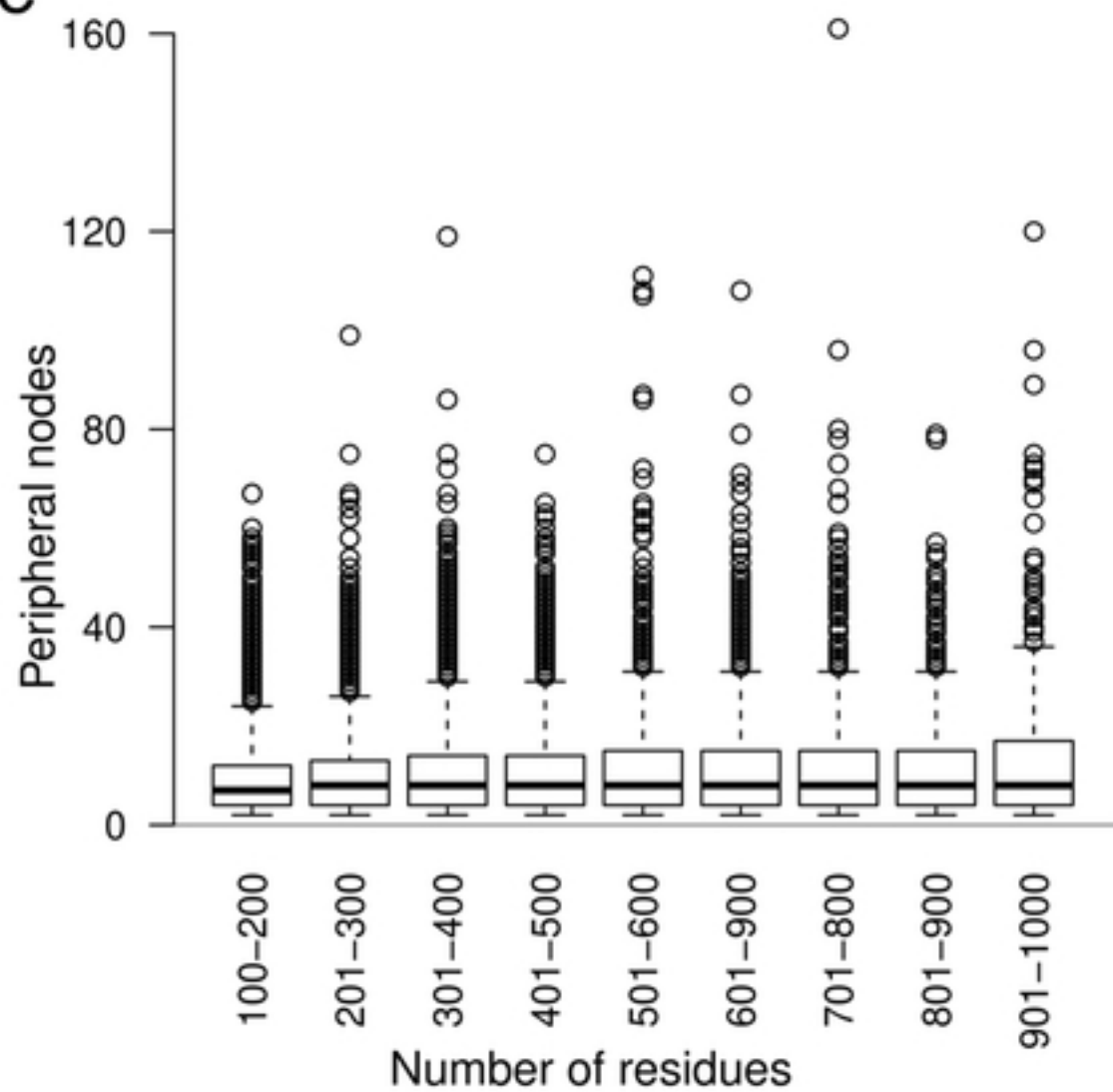
A



B



C



D

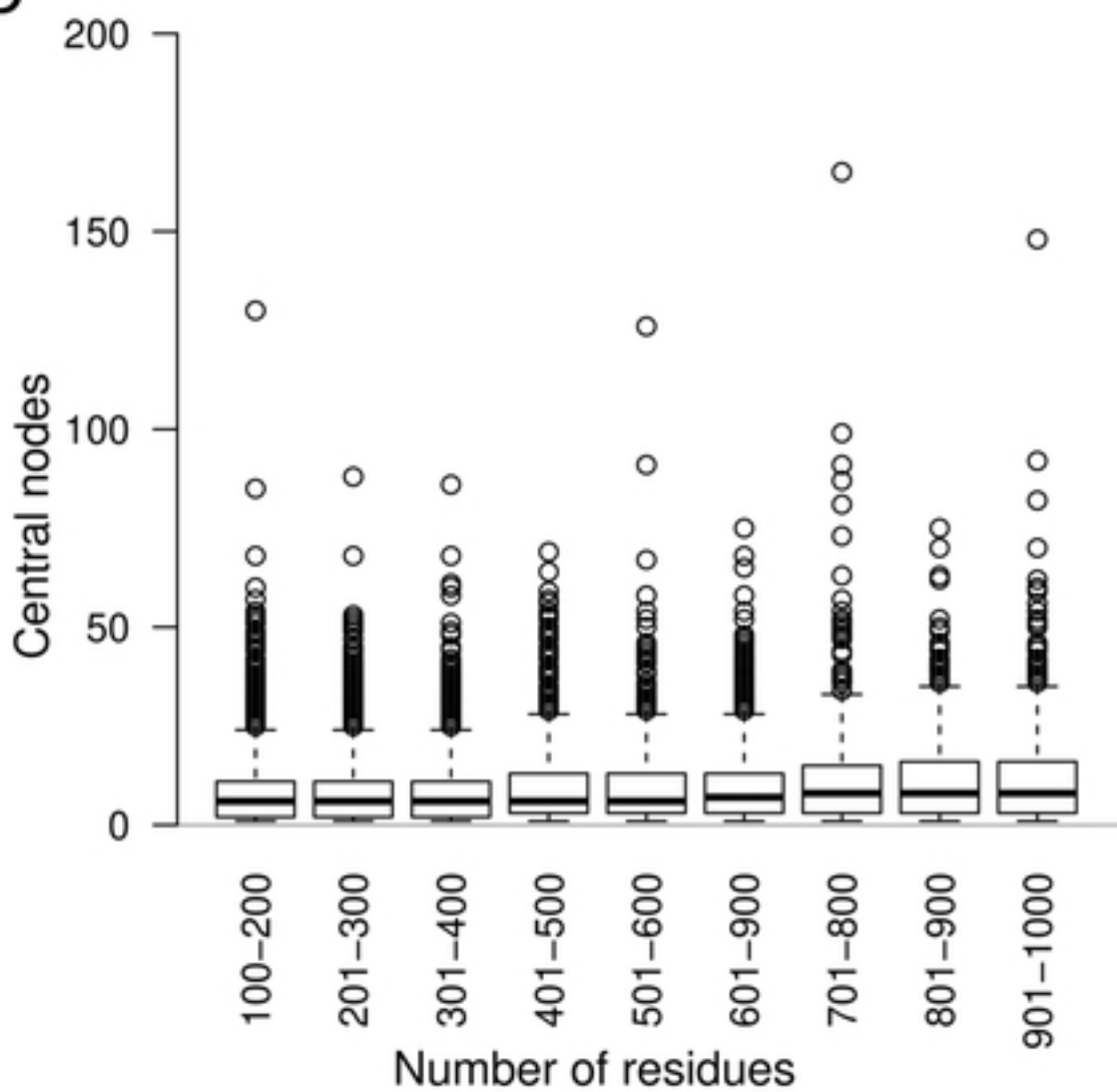


Figure 6

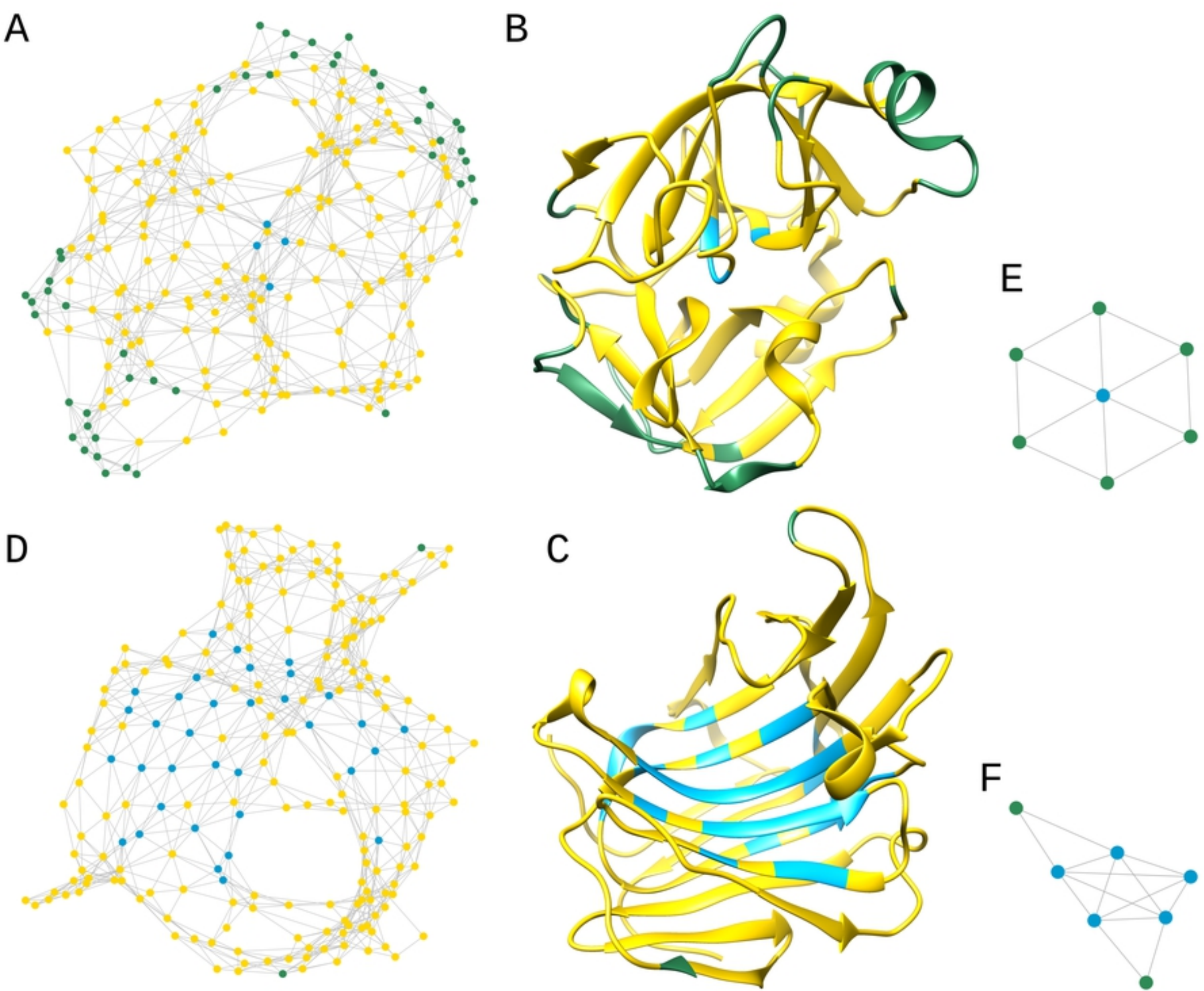


Figure 7

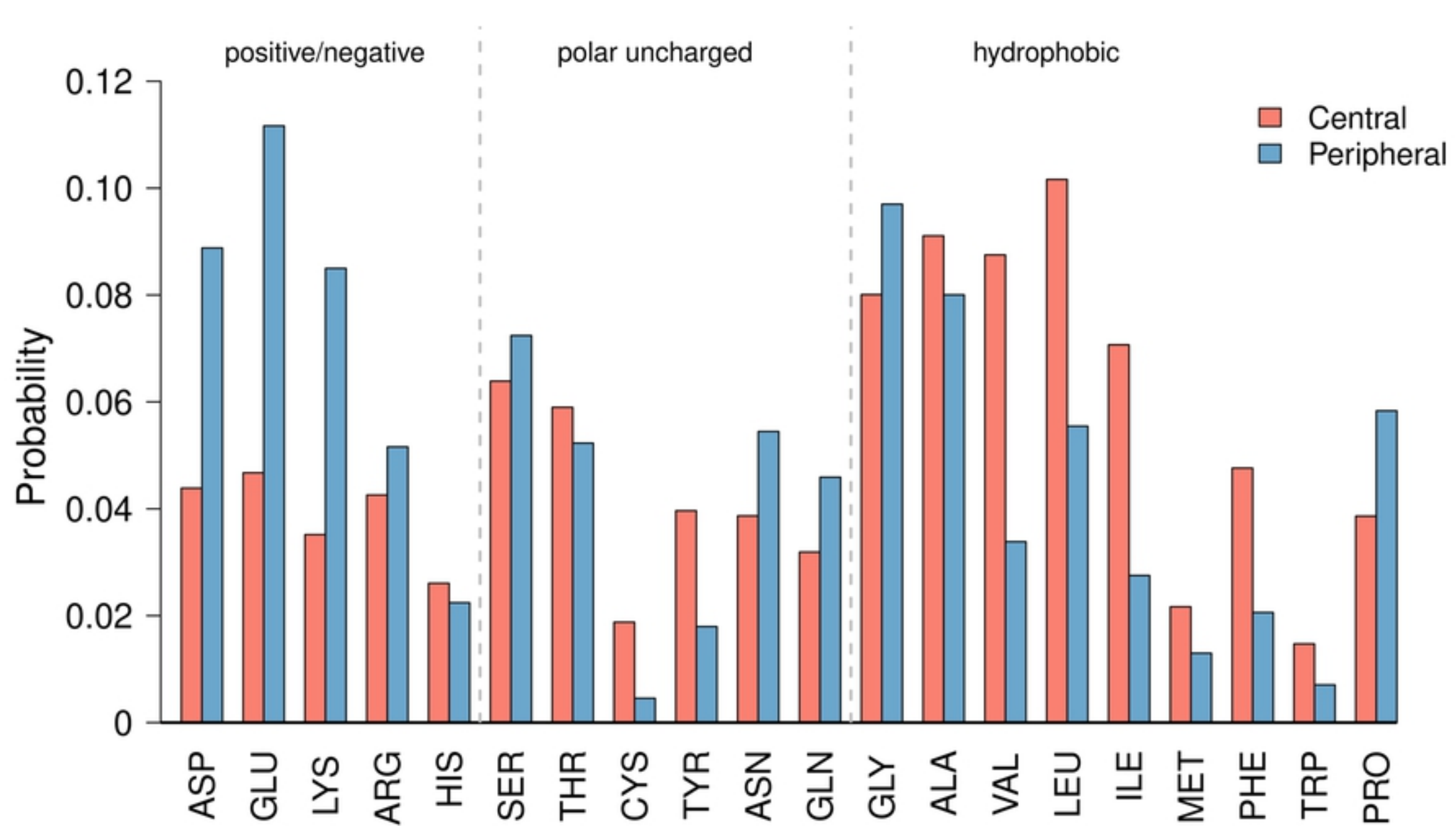


Figure 8