

1 **The length scale of multivalent interactions is evolutionarily conserved in fungal**  
2 **and vertebrate phase-separating proteins**

3  
4 Pouria Dasmeh<sup>1,3,4</sup>, Roman Doronin<sup>1,4</sup> & Andreas Wagner<sup>1,2,4</sup>

5 <sup>1</sup>Institute for Evolutionary Biology and Environmental Studies, University of Zurich, Zurich, Switzerland.

6 <sup>2</sup>The Santa Fe Institute, Santa Fe, New Mexico, <sup>3</sup>Department of Chemistry and Chemical Biology, Harvard  
7 University, Cambridge, MA USA 02139. <sup>4</sup>Swiss Institute of Bioinformatics (SIB).

8

9 **Abstract**

10 One key feature of proteins that form liquid droplets by phase separation inside a cell is  
11 the presence of multiple sites – multivalency – that mediate interactions with other  
12 proteins. We know little about the variation of multivalency on evolutionary time scales.  
13 Here, we investigated the long-term evolution (~600 million years) of multivalency in  
14 fungal mRNA decapping subunit 2 protein (Dcp2), and in the FET protein family. We found  
15 that multivalency varies substantially among the orthologs of these proteins. However,  
16 evolution has maintained the length scale at which sequence motifs that enable protein-  
17 protein interactions occur. That is, the total number of such motifs per hundred amino  
18 acids is higher and less variable than expected by neutral evolution. To help explain this  
19 evolutionary conservation, we developed a conformation classifier using machine-  
20 learning algorithms. This classifier demonstrates that disordered segments in Dcp2 and  
21 FET proteins tend to adopt compact conformations, which is necessary for phase  
22 separation. Thus, the evolutionary conservation we detected may help proteins preserve  
23 the ability to undergo phase separation. Altogether, our study reveals that the length scale  
24 of multivalent interactions is an evolutionarily conserved feature of two classes of phase-  
25 separating proteins in fungi and vertebrates.

26

## 27 **Body**

28 Proteins that undergo liquid-liquid phase separation in a cell have various features that  
29 facilitate their condensation into liquid droplets. The presence of multiple interaction sites  
30 (multivalency) is one of these features.<sup>1-3</sup> Despite the pivotal role of multivalency, we  
31 know little about its evolution. The reason is that multivalency can take different forms,  
32 including the presence of interacting patches on a protein surface, short linear amino acid  
33 motifs, and specific amino acids within the intrinsically disordered regions of phase-  
34 separating proteins<sup>4</sup>.

35 We investigated the evolution of multivalency in two well-known classes of  
36 multivalent phase-separating proteins during ~ 600 million years of evolution. The first  
37 class comprises orthologs of the fungal mRNA decapping subunit 2 protein (Dcp2). Dcp2  
38 is one of the scaffold proteins that help RNA processing bodies (P-bodies) self-assemble  
39 by liquid-liquid phase separation<sup>5-8</sup>. P-bodies are conserved membrane-less eukaryotic  
40 organelles that contribute to the regulation of gene expression by participating in RNA  
41 decay and degradation<sup>9</sup>. They also serve as mRNA storage depots when cells are  
42 stressed<sup>10</sup>. Dcp2 undergoes multivalent interactions using short helical leucine-rich motifs  
43 (HLMs) in its disordered C-terminal domain<sup>11</sup>. HLMs form eight out of 12 identified  
44 interactions between Dcp2 and other core proteins in P-bodies<sup>5</sup>.

45 The second class of proteins comprises orthologs of six members of the FET family  
46 of RNA-binding proteins, including FUS, EWS, HNRNPA1, HNRNPA3, HNRNPR, and  
47 TAF15. These proteins have a common domain architecture that consists of a prion-like  
48 domain (PLD), and other domains with RNA/DNA binding affinities<sup>10,12-14</sup>. They contribute

49 to DNA damage repair, transcriptional control, and the regulation of the life-time of RNAs  
50 in metazoan species<sup>13</sup>. The prion-like domain of these proteins has low sequence  
51 complexity, and is enriched in few amino acids, such as asparagine, glutamine, tyrosine,  
52 and glycine<sup>15,16</sup>. The aromatic residues within the prion-like domain, particularly tyrosine,  
53 and arginine in RNA binding domains, are responsible for the multivalency of these  
54 proteins<sup>17-19</sup>. Interactions between these residues drive the phase-separation of FET  
55 proteins<sup>20</sup>.

56 We use the stickers-and-spacers representation<sup>21</sup> of phase-separating proteins  
57 throughout this work. Stickers are specific amino acids, motifs or protein domains whose  
58 interactions derive phase separation. Spacers are the sequences that separate the  
59 stickers. The helical leucine-rich motifs in Dcp2 and the aromatic residues in FET proteins  
60 are such stickers. For the FET proteins, we particularly focus on tyrosine residues,  
61 because their number and patterning in the sequence modulates the phase-separation  
62 propensity of these proteins<sup>20,22,23</sup>.

63 We first investigated the evolution of helical leucine-rich motifs in 48 Dcp2 proteins  
64 of the phylum *Ascomycota* (Dataset S1). HLMs lie within the disordered C-terminal  
65 domain of Dcp2 (Figure 1A, residues 229 - 930 in *S. cerevisiae*), and take the form LL-  
66 x $\phi$ -L, where L stands for leucine,  $\phi$  is a hydrophobic residue, and x represents any amino  
67 acid. We identified 347 motifs in these sequences that exactly matched the LL-x $\phi$ -L  
68 pattern (Figure S1). As shown in Figure 1B, HLMs are the most conserved sequence  
69 segments within the intrinsically disordered C-terminal domain of Dcp2. However, their  
70 number substantially varies from a minimum of three in *L. elongisporus* to a maximum of

71 16 in *K. capsulata* (Table S1). Importantly, the number of HLMS increases with the length  
72 of the disordered C-terminal domain of a Dcp2 sequence (Figure 1C; Spearman  
73 correlation,  $R=0.44$ ,  $p=0.0017$ ). The average and median length of the spacer segments  
74 that separate HLMS are  $\sim 70$  and  $\sim 51$  amino acids. Based on these observations, we  
75 hypothesized that the scaling between the number of HLMS and the length of the  
76 disordered domain (1 HLM in  $\sim 70$  residues) reflects a requirement for a characteristic  
77 sequence length that separates sticker motifs. We tested this hypothesis by asking  
78 whether this characteristic sequence length may be subject to natural selection.

79 To understand the evolutionary forces that shape the scaling between the number  
80 of HLMS and the length of the C-terminal domain in Dcp2, we determined the likelihood  
81 that HLMS arise by chance through neutral evolution. To this end, we simulated neutral  
82 protein sequence evolution, using realistic divergence times of real Dcp2 sequences (see  
83 Method for details). We found that neutral evolution can indeed create motifs that exactly  
84 match known HLMS (Figure S2; Dataset S2), but the fraction of these neutrally-evolved  
85 HLMS per unit sequence length was much lower than that of HLMS in real sequences.  
86 Specifically, neutral evolution creates only one HLM per  $\sim 1500$  amino acids. In other  
87 words, HLMS in neutrally evolving sequences are  $\sim 35$  times less frequent than in real  
88 Dcp2 sequences (Figure 1D). We recalculated the fraction of HLMS per unit of sequence  
89 length for various codon frequencies, nonsynonymous substitution rates, and values of  
90 transition/transversion bias (Dataset S2). In all these calculations, we found a  
91 substantially higher incidence of HLMS per unit of sequence length in biological  
92 sequences compared to sequences evolved by neutral evolution (Table S2).

93 We also compared the distribution of spacer lengths (segments that separate  
94 HLMS) in the C-terminal domain of Dcp2 orthologs with that of neutrally evolved  
95 sequences. The median length of spacers is 81 amino acids in neutrally evolved  
96 sequences, significantly higher than the 51 amino acids in biological Dcp2 sequences ( $p$   
97  $\sim 10^{-6}$ ; Wilcoxon ran-sum test; Figure 1E). In addition, spacer lengths are significantly  
98 more variable in neutrally evolved sequences compared to the biological Dcp2 proteins  
99 ( $p \sim 10^{-7}$ , one-sided F-test for the equality of variances), and the length distributions are  
100 significantly different ( $p \sim 10^{-8}$ ; Kolmogorov-Smirnov test). Altogether, these results show  
101 that evolution has not only increased the incident of HLMS in Dcp2 sequences, but also  
102 has stabilized the lengths of sequences that separate HLMS.

103 To find out whether the scaling of sticker number with the length of a disordered  
104 region is a more general property, we next studied the FET family of proteins in  
105 vertebrates. We identified  $\sim 200$ -300 orthologs for each of the six FET proteins, and  
106 compiled a set of 1480 sequences of these proteins (Dataset S3). Analogous to Dcp2  
107 and its HLMS, we observed that longer FET proteins have more arginine (R) and tyrosine  
108 (Y) sticker residues in their prion like domain (Figure 2A, Spearman correlation,  $R=0.8$ ,  
109  $p < 10^{-16}$ ). Importantly, among all 20 amino acids, the number of Rs and Ys showed the  
110 highest correlation with sequence length (Figure 2B; adjusted  $R^2 \sim 0.81$ , and 0.70 in a  
111 linear model with 5-fold cross-validation and 10 replicates). In sum, the scaling of  
112 multivalency with the lengths of disordered domains is not unique to Dcp2 in fungi. It also  
113 exists in the FET protein family of vertebrates.

114 We further examined the spacer lengths that separate Ys and Rs in the sequence  
115 of FET proteins to find out whether natural selection has influenced the number of stickers  
116 per unit sequence length. We compared the distance distribution of both R and Y residues  
117 in FET proteins with that of neutrally-evolved sequences (see Methods for details). For  
118 both amino acids, the distribution of distances between tyrosine residues in FET proteins  
119 is significantly less variable than that of neutrally evolved sequences (See Figure 2C for  
120 spacers between tyrosine residues;  $p < 10^{-16}$ ; Kolmogorov-Smirnov test, and Figure S3  
121 for spacers between arginine residues). The median distance between tyrosine residues  
122 is seven amino acids for FET proteins, which is significantly less than the corresponding  
123 distance of 11 amino acids in neutrally evolving proteins ( $p \sim 10^{-6}$ ; Wilcoxon rank-sum  
124 test). In addition, the distance distribution of FET proteins is much more sharply peaked  
125 (leptokurtic, Figure 2C) and significantly differed from neutrally-evolved sequences ( $p \sim$   
126  $10^{-8}$ ; Kolmogorov-Smirnov test). This suggests that natural selection has likely stabilized  
127 this distance distribution in FET proteins.

128 Next we asked why the scaling of the number of stickers may be conserved,  
129 focusing on the hypothesis that it helps maintain a network of protein interactions that is  
130 necessary for condensation and phase separation<sup>24</sup>. To maintain this interaction network,  
131 disordered sequences should be able to adopt compact conformations. The reason is  
132 that only this type of conformation substantially increases the chance of interactions  
133 between stickers<sup>24</sup>. We thus wanted to find out whether this ability exists in our proteins.  
134 First, we calculated the fraction of charged residues in the spacers that separate HLMs  
135 in Dcp2 and aromatic residues in the prion-like domain of FET proteins. This fraction of

136 charged residues is a proxy for the effective solvation and hence the conformation of  
137 disordered spacers<sup>24</sup>. Previous studies have suggested that spacers whose fraction of  
138 charged residues is less than 0.5 can self-associate and drive the formation of a  
139 condensation-promoting network of interactions (Figure 3A). As shown in Figures 3B-C,  
140 we found that almost all spacers in Dcp2 and FET proteins have a fraction of charged  
141 residues between 0.2 and 0.4, indicating that they can adopt compact conformations.

142         Second, we predicted a structural feature of disordered sequences known as the  
143  $\Delta$ -parameter. This parameter is the average difference between inter-residue distances  
144 of a disordered sequence and the corresponding distances of a typical Flory random  
145 coil<sup>24</sup>. Flory random coils are an idealized kind of disordered sequences in which the  
146 attractive and repulsive forces between residues and solvent molecules are at balance.  
147 Spacers that self-associate and promote phase-separation are characterized by  $\Delta \leq 0.1$   
148 nm. As  $\Delta$  increases beyond 0.1 nm, spacers adopt more extended conformations,  
149 resembling another type of idealized sequence known as a self-avoiding random coil  
150 (Figure 3A).

151         We developed a sequence-based classifier of  $\Delta$  using a random forest algorithm  
152 (Figure 3D), which classifies spacers based on their amino acid properties into two  
153 classes, those with  $\Delta > 0.1$ , and those with  $\Delta \leq 0.1$  nm (see Methods for details). We  
154 trained this classifier on a dataset of 256 naturally occurring disordered sequences whose  
155  $\Delta$  values had been previously calculated by molecular dynamics simulations<sup>24</sup>. This  
156 classifier achieved an accuracy of  $\sim 0.88$  in 100 independent runs with the data split into  
157 a training set (80% of the data) and a testing set (20%) (Figure S4 see Methods for

158 details). Using this classifier, we found that ~94.8% of all spacers in Dcp2 have a  
159 predicted value of  $\Delta$  below 0.1 nm. We repeated this analysis for the spacers in the prion-  
160 like domain of the FET proteins and found that in these proteins too, most spacers  
161 (~99.3%) have predicted  $\Delta \leq 0.1$  nm. Altogether, these results indicate that both fungal  
162 Dcp2 sequences and vertebrate FET proteins have spacers that can self-associate and  
163 promote phase-separation in these proteins.

164 In summary, our work reveals that evolution has maintained a characteristic length  
165 scale of multivalent sticker sequences in two classes of multivalent proteins during ~600  
166 million years. Our results extend the previous observation by Martin *et al.*<sup>22</sup> that a uniform  
167 patterning of tyrosine residues in few members of FET proteins promotes phase-  
168 separation and inhibits the aggregation of these proteins. This scaling not only promotes  
169 phase-separation, but may also increase the robustness of proteins to DNA mutations  
170 such as indels and truncations. Dcp2 plays an important role in the assembly of RNA  
171 processing bodies, and FET proteins play such a role in the assembly of stress granules.  
172 Biomolecular condensates like these are sensitive to environmental stressors such as  
173 heat shock and energy depletion<sup>16,25-27</sup>. Our results thus also raise the intriguing  
174 possibility that evolution may have modulated the multivalency of proteins in membrane-  
175 less organelles to help organisms cope with new environments. To provide experimental  
176 support for this possibility is an exciting question for future work.

177

178

179



180 **Methods:**

181 **Data compilation and the generation of neutrally-evolved sequences**

182 In this study, we used 48 orthologous coding sequences of fungal Dcp2, and ~200-300  
183 orthologs of six members of the FET family of proteins (FUS, EWS, HNRNPA1,  
184 HNRNPA3, HNRNPR, and TAF15; overall 1480 sequenced). We downloaded these  
185 sequences from the NCBI<sup>28</sup>, ENSEMBL<sup>29</sup>, and KEGG<sup>30</sup> databases. Throughout, we  
186 worked with the amino acid sequences of these proteins, except for the simulation of  
187 neutral evolution, where we represented protein sequences on the level of DNA.

188

189 **Simulation of neutrally evolved sequences**

190 We simulated protein evolution using the Evolver package within the PAML suite<sup>31</sup>. In  
191 brief, Evolver uses Monte Carlo simulations to generate codon sequences using a  
192 specified phylogenetic tree with given branch lengths, nucleotide frequencies,  
193 transition/transversion bias ( $\kappa$ ), and the ratio of the rate of nonsynonymous to  
194 synonymous substitutions ( $dN/dS$ )<sup>31</sup>. To simulate neutral sequence evolution, we used  
195 the standard genetic code with codon frequencies from our study proteins sequences.  
196 Specifically, we used codon frequencies from the set of 48 Dcp2 sequences to model  
197 neutral evolution in fungal Dcp2, and codon frequencies from FUS orthologs for neutral  
198 evolution in vertebrate proteins. We used a consensus phylogenetic tree for the fungal  
199 species from the yeast genome browser<sup>32</sup> and for the vertebrate species from the  
200 TimeTree database<sup>33</sup>. To model neutral evolution, we set  $dN/dS$  to 1 and used a  
201 transition/transversion rate ratio of 2.3, and 2.9 for fungal and mammalian sequences.

202 We estimated these values by fitting the codon model M1 to the phylogenetic tree and  
203 the sequences of these proteins. This model assumes that all branches of the  
204 phylogenetic tree have the same rate of evolution. We evaluated the number of neutrally-  
205 evolved HLMs for various values of  $dN/dS$  and the transition/transversion rate ratio to  
206 ensure that our results do not depend on the choice of these parameters (Table S1).  
207 Overall, we generated  $10^4$  evolved sequences using this sequence evolution model.

208

### 209 **Detection of HLM motifs and their distinct flanking regions**

210 We used regular expression matching to search for HLM motifs that matched the LL-x $\phi$ -  
211 L pattern, where  $\phi$  is a hydrophobic residue (one of the amino acids L, I, V, A, P, and F),  
212 and x represents any amino acid. To distinguish HLMs from HLM-like patterns we used  
213 the classification approaches of logistic regression and random forests implemented in  
214 the Python package scikit-learn. In these classifications, positive and the negative sets  
215 correspond to the flanking regions of HLMs and HLM-like motifs, respectively. The size  
216 of the training and the test set was 80%, and 20% of the whole dataset.

217

218

### 219 **Random forest classification of spacers**

220 We used the random forest algorithm to develop a classifier of spacer conformation from  
221 the protein sequence. To this end, we used the average deviation of inter-residue  
222 distances of a spacer sequence from the same distances in a Flory Random Coil as the  
223 measure for the prediction of spacer types<sup>24</sup>. This deviation, known as the  $\Delta$  parameter,  
224 can take positive and negative values. Disordered sequences with  $\Delta \leq 0.1$  have the

225 propensity to form compact conformations. We used a binary classification and classified  
226 proteins into a positive set ( $\Delta \leq 0.1$ ) and a negative set ( $\Delta > 0.1$ ). To train our classifier  
227 we used a dataset of 256 disordered sequences for which we had calculated  $\Delta$  by all-  
228 atom molecular dynamic simulations.

229 To build features for the classification, we calculated the average value of 500  
230 physicochemical properties for each sequence in the positive and the negative sets. This  
231 yielded two feature matrices, one for sequences with  $\Delta \leq 0.1$ , and another for sequences  
232 with  $\Delta > 0.1$ . To apply random forest classification, we used the randomForest package  
233 of R<sup>34</sup>, and evaluated the best number of trees (*nTree*) and the number of variables  
234 randomly sampled at each split (*mtry*) in the random forest algorithm. To do so, we  
235 systematically varied *nTree* and *mtry*, and calculated the accuracy of classification with  
236 10-fold cross-validation in 3 replicates. We defined accuracy as the percentage of  
237 correctly identified classes of spacers ( $\Delta \leq 0.1$  and  $\Delta > 0.1$ ) out of all spacers. The  
238 combination of *nTree*=5000 trees and *mtry*=10 variables achieved the highest accuracy  
239 of ~88%. Here, we define accuracy as the ratio of the number of true positives to the sum  
240 of true positives and false negatives. We then used these parameters to perform 100  
241 random forest clusterings, in which we randomly assigned proteins to the training and the  
242 testing datasets. To quantify the accuracy of classification we counted the number of true  
243 positive and false positive predictions and calculated the area under the curve (AUC). We  
244 represented these values by receiver operating characteristic curves (ROC) in Figure S2.  
245 We performed all statistical analyses using R. Scripts and input files for classification, as

246 well as for evolutionary simulations are available at:  
247 [https://github.com/dasmeh/multivalency\\_evolution](https://github.com/dasmeh/multivalency_evolution)

248

## 249 **Acknowledgement:**

250 The authors would like to thank Simon Alberti for careful reading of the manuscript and  
251 for helpful discussions on the evolution of liquid-liquid phase separation.

252

## 253 **References**

- 254 1. Banani, S.F., Lee, H.O., Hyman, A.A. & Rosen, M.K. Biomolecular condensates:  
255 organizers of cellular biochemistry. *Nature reviews Molecular cell biology* **18**, 285-298  
256 (2017).
- 257 2. Li, P. *et al.* Phase transitions in the assembly of multivalent signalling proteins. *Nature*  
258 **483**, 336-340 (2012).
- 259 3. Brangwynne, C.P., Tompa, P. & Pappu, R.V. Polymer physics of intracellular phase  
260 transitions. *Nature Physics* **11**, 899-904 (2015).
- 261 4. Posey, A.E., Holehouse, A.S. & Pappu, R.V. Phase separation of intrinsically disordered  
262 proteins. in *Methods in Enzymology*, Vol. 611 1-30 (Elsevier, 2018).
- 263 5. Xing, W., Muhlrads, D., Parker, R. & Rosen, M.K. A quantitative inventory of yeast P body  
264 proteins reveals principles of composition and specificity. *Elife* **9**, e56525 (2020).
- 265 6. Parker, R. & Sheth, U. P bodies and the control of mRNA translation and degradation.  
266 *Molecular cell* **25**, 635-646 (2007).
- 267 7. Rao, B.S. & Parker, R. Numerous interactions act redundantly to assemble a tunable size  
268 of P bodies in *Saccharomyces cerevisiae*. *Proceedings of the National Academy of Sciences*  
269 **114**, E9569-E9578 (2017).
- 270 8. Kroschwald, S. *et al.* Promiscuous interactions and protein disaggregases determine the  
271 material state of stress-inducible RNP granules. *elife* **4**, e06807 (2015).
- 272 9. Anderson, P. & Kedersha, N. RNA granules: post-transcriptional and epigenetic  
273 modulators of gene expression. *Nature reviews Molecular cell biology* **10**, 430-436 (2009).
- 274 10. Aizer, A. *et al.* Quantifying mRNA targeting to P-bodies in living human cells reveals their  
275 dual role in mRNA decay and storage. *Journal of cell science* **127**, 4443-4456 (2014).
- 276 11. Jonas, S. & Izaurralde, E. The role of disordered protein regions in the assembly of  
277 decapping complexes and RNP granules. *Genes & development* **27**, 2628-2641 (2013).
- 278 12. Hoell, J.I. *et al.* RNA targets of wild-type and mutant FET family proteins. *Nature*  
279 *structural & molecular biology* **18**, 1428-1431 (2011).
- 280 13. Schwartz, J.C., Cech, T.R. & Parker, R.R. Biochemical properties and biological functions  
281 of FET proteins. *Annual review of biochemistry* **84**, 355-379 (2015).
- 282 14. Svetoni, F., Frisone, P. & Paronetto, M.P. Role of FET proteins in neurodegenerative  
283 disorders. *RNA biology* **13**, 1089-1102 (2016).

- 284 15. Aguzzi, A. & Calella, A.M. Prions: protein aggregation and infectious diseases.  
285 *Physiological reviews* **89**, 1105-1152 (2009).
- 286 16. Franzmann, T.M. & Alberti, S. Prion-like low-complexity sequences: Key regulators of  
287 protein solubility and phase behavior. *Journal of Biological Chemistry* **294**, 7128-7136  
288 (2019).
- 289 17. Burke, K.A., Janke, A.M., Rhine, C.L. & Fawzi, N.L. Residue-by-residue view of in vitro  
290 FUS granules that bind the C-terminal domain of RNA polymerase II. *Molecular cell* **60**,  
291 231-241 (2015).
- 292 18. Hofweber, M. *et al.* Phase separation of FUS is suppressed by its nuclear import receptor  
293 and arginine methylation. *Cell* **173**, 706-719. e13 (2018).
- 294 19. Patel, A. *et al.* A liquid-to-solid phase transition of the ALS protein FUS accelerated by  
295 disease mutation. *Cell* **162**, 1066-1077 (2015).
- 296 20. Wang, J. *et al.* A molecular grammar governing the driving forces for phase separation of  
297 prion-like RNA binding proteins. *Cell* **174**, 688-699. e16 (2018).
- 298 21. Choi, J.-M., Holehouse, A.S. & Pappu, R.V. Physical principles underlying the complex  
299 biology of intracellular phase transitions. *Annual Review of Biophysics* **49**, 107-133 (2020).
- 300 22. Martin, E.W. *et al.* Valence and patterning of aromatic residues determine the phase  
301 behavior of prion-like domains. *Science* **367**, 694-699 (2020).
- 302 23. Bremer, A. *et al.* Deciphering how naturally occurring sequence features impact the phase  
303 behaviors of disordered prion-like domains. *bioRxiv* (2021).
- 304 24. Harmon, T.S., Holehouse, A.S., Rosen, M.K. & Pappu, R.V. Intrinsically disordered  
305 linkers determine the interplay between phase separation and gelation in multivalent  
306 proteins. *elife* **6**, e30294 (2017).
- 307 25. Boeynaems, S. *et al.* Protein phase separation: a new phase in cell biology. *Trends in cell*  
308 *biology* **28**, 420-435 (2018).
- 309 26. Zarin, T. *et al.* Proteome-wide signatures of function in highly diverged intrinsically  
310 disordered regions. *Elife* **8**(2019).
- 311 27. O'Connell, J.D. *et al.* A proteomic survey of widespread protein aggregation in yeast.  
312 *Molecular BioSystems* **10**, 851-861 (2014).
- 313 28. Pruitt, K.D., Tatusova, T. & Maglott, D.R. NCBI reference sequences (RefSeq): a curated  
314 non-redundant sequence database of genomes, transcripts and proteins. *Nucleic acids*  
315 *research* **35**, D61-D65 (2006).
- 316 29. Hubbard, T. *et al.* The Ensembl genome database project. *Nucleic acids research* **30**, 38-  
317 41 (2002).
- 318 30. Kanehisa, M. The KEGG database. in *Novartis Foundation Symposium* 91-100 (Wiley  
319 Online Library, 2002).
- 320 31. Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular biology and*  
321 *evolution* **24**, 1586-1591 (2007).
- 322 32. Byrne, K.P. & Wolfe, K.H. The Yeast Gene Order Browser: combining curated homology  
323 and syntenic context reveals gene fate in polyploid species. *Genome research* **15**, 1456-  
324 1461 (2005).
- 325 33. Hedges, S.B., Dudley, J. & Kumar, S. TimeTree: a public knowledge-base of divergence  
326 times among organisms. *Bioinformatics* **22**, 2971-2972 (2006).
- 327 34. Liaw, A. & Wiener, M. Classification and regression by randomForest. *R news* **2**, 18-22  
328 (2002).
- 329

330

331

332

333

334 **Figure legends:**

335

336 **Figure 1. The length scale of multivalent interactions is evolutionary conserved in**

337 **fungi species Dcp2.** A) Architecture of Dcp2 in *S. cerevisiae* with the regulatory domain

338 (in red), the NUDIX catalytic domain (in orange), and the disordered C-terminal domain

339 (in white). Within the disordered C-terminal domain helical leucine-rich short linear motifs

340 are responsible for the multivalency of Dcp2. B) Multiple sequence alignment of the C-

341 terminal domain of Dcp2 in 48 fungal species within the phylum of *Ascomycota* spanning

342 ~ 600 million years of evolution. HLMs, shown as blue columns, are highly conserved

343 within the C-terminal domain of Dcp2. C) The number of HLMs positively correlates with

344 the length of the C-terminal domain of Dcp2 in fungi (Spearman correlation;  $R=0.44$ ,

345  $p=0.0017$ ). D) The incidence of HMLs in biological sequences (shown in red) is ~ 35 times

346 higher than that of neutrally evolved sequences (shown in blue). E) The distribution of

347 spacer lengths (sequences that separate HLMs) in real Dcp2 sequences (shown in red),

348 and in neutrally evolved sequences (shown in blue). We compared the two distributions

349 and calculated the  $p$ -value for rejecting the null hypothesis that these distributions are

350 indistinguishable by Kolmogorov-Smirnov (KS) test.

351

352

353

354

355

356 **Figure 2. The length scale of multivalent interactions is evolutionary conserved in**  
357 **the FET family of vertebrate proteins.** A) The number of arginine (R) and tyrosine (Y)  
358 residues of six different FET family members and their orthologs in vertebrate species  
359 (1180 proteins overall) versus their sequence length. B) The coefficient of determination  
360 ( $R^2$ ) between the number of different amino acids and the length of FET proteins and their  
361 orthologs. For a robust estimation of  $R^2$ , we used a linear regression model with 5-fold  
362 cross-validation that we repeated 10 times. C) The distribution of distances between  
363 tyrosine residues in FET proteins (shown in red), and in neutrally-evolved sequences  
364 (shown in blue). We compared the two distributions and calculated the  $p$ -value for  
365 rejecting the null hypothesis that these distributions are indistinguishable by Kolmogorov-  
366 Smirnov (KS) test.

367  
368 **Figure 3. The disordered spacers in fungal Dcp2 and vertebrates FET proteins**  
369 **adopt conformations that promote phase-separation.** A) The fraction of charged  
370 residues (FCR) can distinguish the conformation of spacer segments in multivalent  
371 proteins. Proteins with  $FCR > 0.5$  preferentially adopt extended conformations like  
372 idealized self-avoiding random coils. Proteins with  $FCR < 0.3$  can form compact globules.  
373 Sequences with intermediate values of FCR form conformations similar to Flory random  
374 coils where the net attractive and repulsive forces between residues and solvent  
375 molecules are in balance. The fraction of charged residues for B) fungal Dcp2 sequences,  
376 and C) vertebrate HNRNPA1a, a member of the FET family. D) Schematic for machine-  
377 learning random-forest classification to classify spacer types from their amino acid



378 sequence. In brief, we used the sequences of naturally occurring disordered sequences  
379 that connect different domains, calculated the average of 500 amino acid properties for  
380 each sequence, and used this dataset to classify these sequences into the two categories  
381 of self-avoiding random coils, and Flory-random coils and compact globules. E) The  
382 fraction of spacers that adopt compact conformations (Flory random coils, and compact  
383 globules) and those that adopt extended conformations (self-avoiding random coils) in  
384 fungal Dcp2 and vertebrate FET proteins.  
385

Figure 1

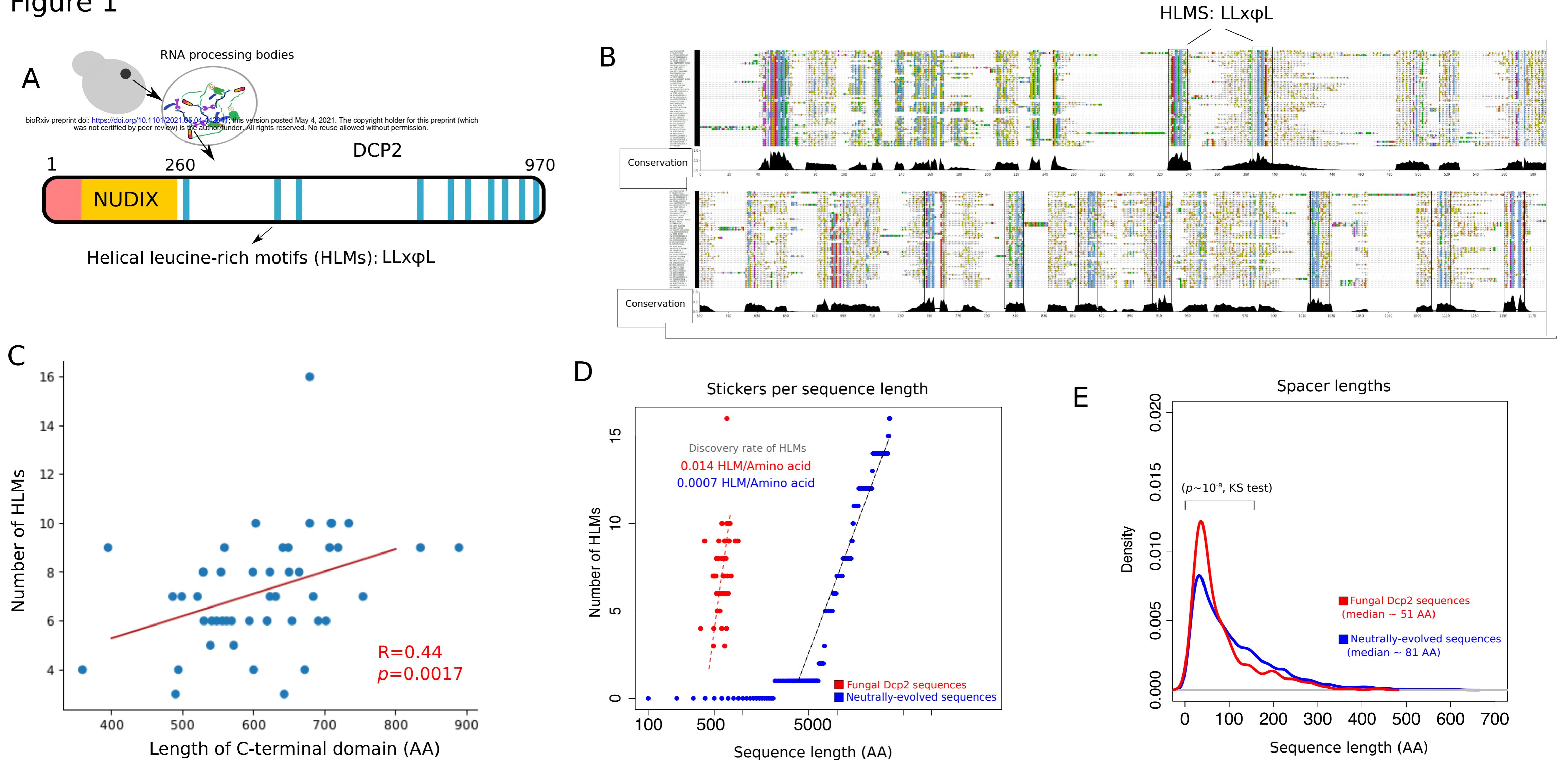


Figure 2

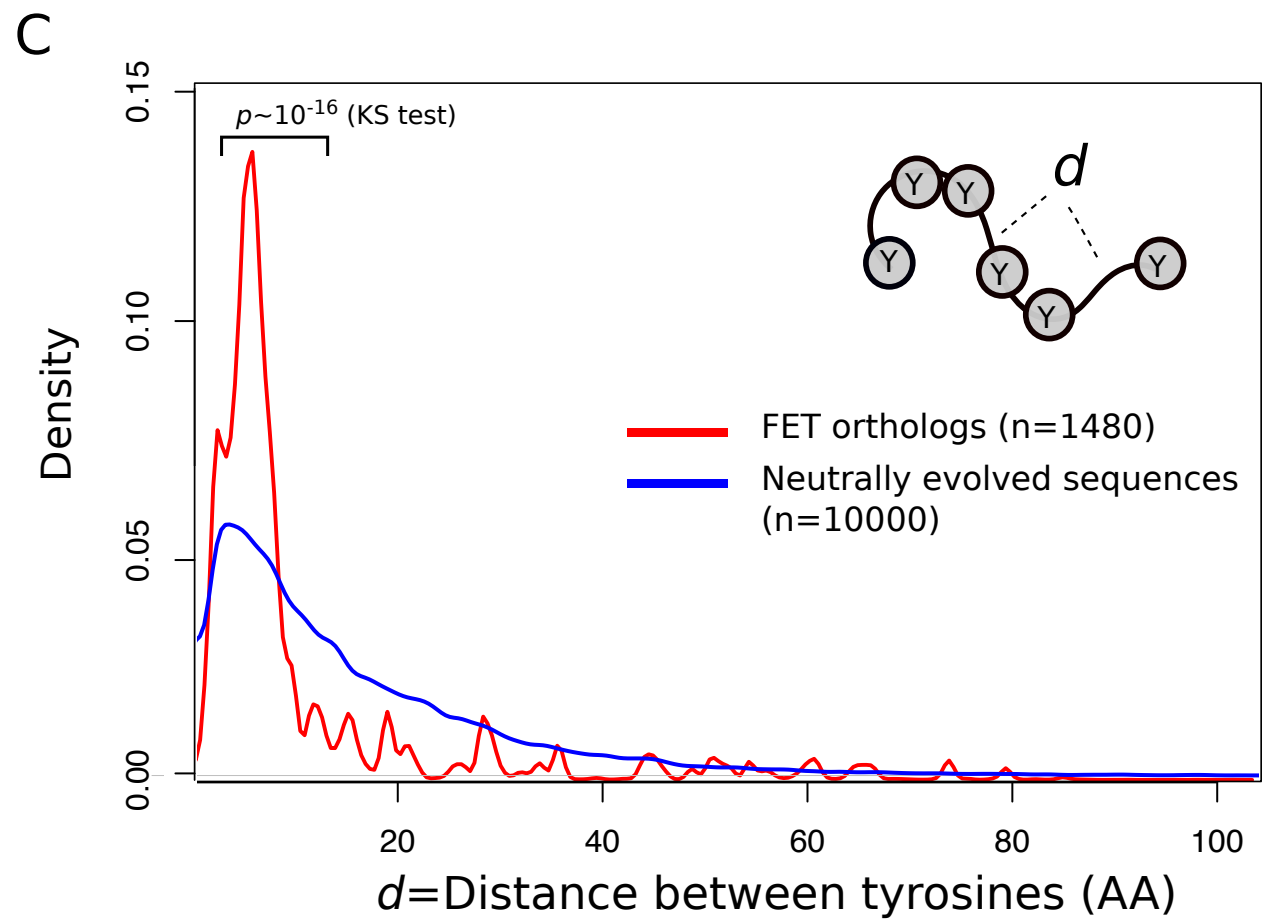
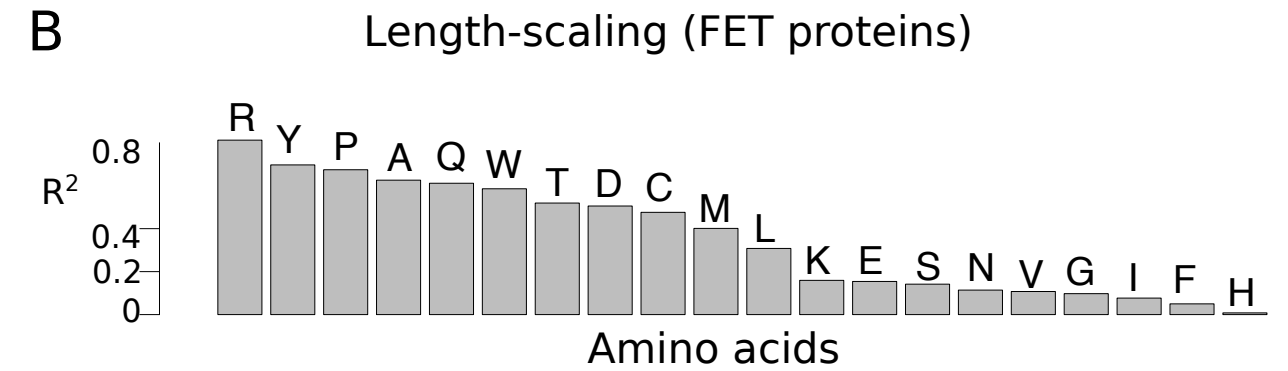
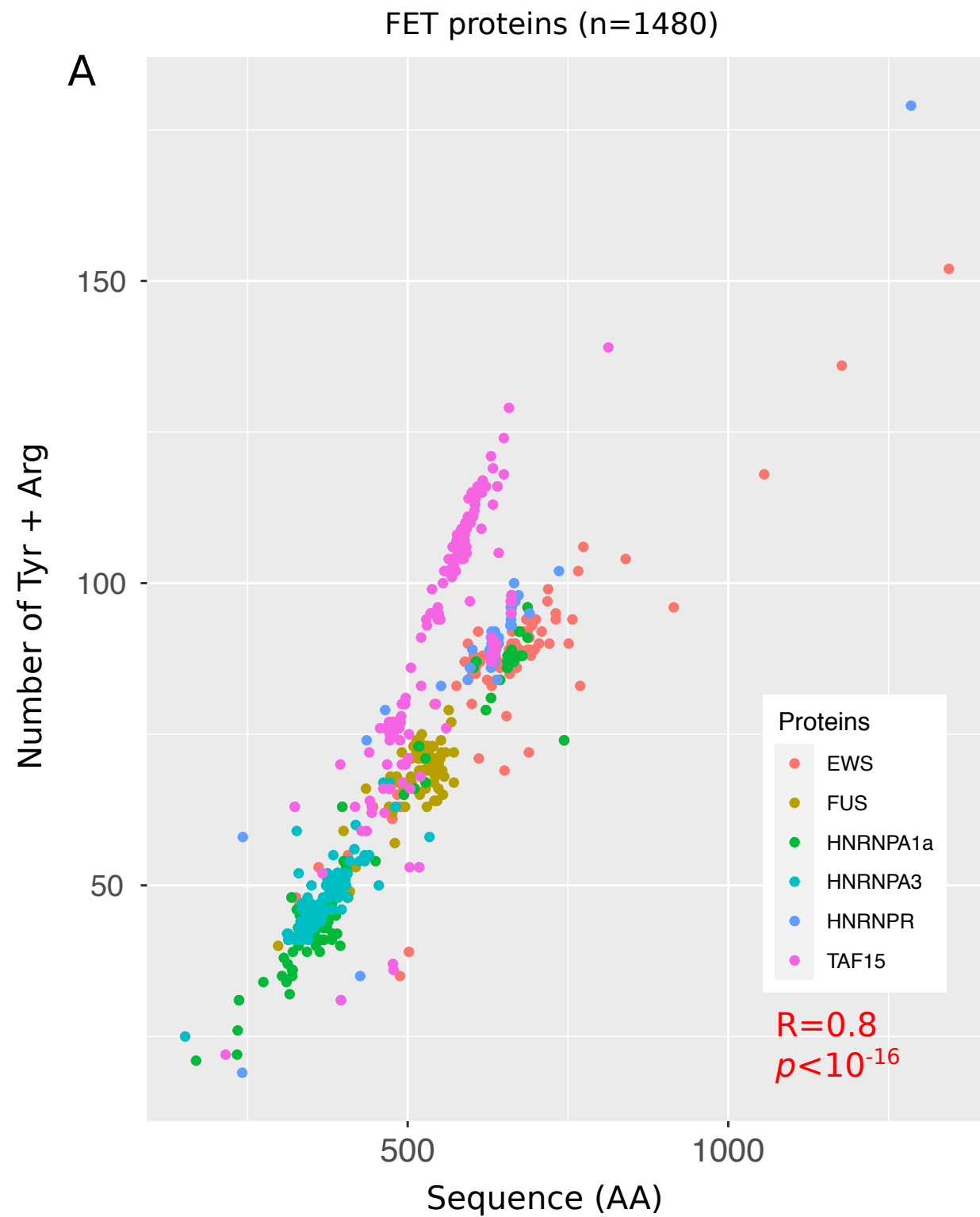


Figure 3

