

1 **Title Page**

2 **Title:** An extension of Shannon's entropy to explain taxa diversity and human
3 diseases

4 **Running title:** A mathematical interpretation of life

5

6 **Author list and full affiliations:**

7 Farzin Kamari^{*1,2}, MD, MPH; Sina Dadmand^{2,3}, PharmD.

8 ¹Neurosciences Research Centre, Tabriz University of Medical Sciences, Tabriz, Iran.

9 ²Synaptic ProteoLab, Synaptic ApS, Skt Knuds Gade 20, 5000 Odense C, Denmark.

10 ³Faculty of Pharmacy, Tabriz University of Medical Sciences, Tabriz, Iran.

11 ***Corresponding author:** Farzin Kamari

12 **Email:** kamari.farzin@gmail.com

13

14 **Total character count (with spaces):** 72,502

15 **Keywords:** origin of diseases/ protein-protein interaction/ Shannon's entropy/ taxonomic
16 classification/ tree of life

17 **Abstract**

18 In this study, with the use of the information theory, we have proposed and proved a
19 mathematical theorem by which we argue the reason for the existence of human diseases. To
20 introduce our theoretical frame of reference, first, we put forward a modification of
21 Shannon's entropy, computed for all available proteomes, as a tool to compare systems
22 complexity and distinguish between the several levels of biological organizations. We
23 establish a new approach, namely the wave of life, to differentiate several taxa and
24 corroborate our findings through the latest tree of life. Furthermore, we found that human
25 proteins with higher mutual information, derived from our theorem, are more prone to be
26 involved in human diseases. Our results illuminate the dynamics of protein network stability
27 and offer probable scenarios for the existence of human diseases and their varying occurrence
28 rates. The current study presents the fundamentals in understanding human diseases by means
29 of information theory. In practice, the theorem proposes multiple-protein approach as
30 therapeutic agents targeting protein networks as a whole, rather than approaching a single
31 receptor.

32 **Introduction**

33 The term ‘entropy’ was originally introduced by Rudolf Clausius in thermodynamics more
34 than one and a half centuries ago (Clausius, 1864). Entropy is predominantly known as a
35 measure of the disorder and uncertainty in dynamic systems (Ghahramani, 2006; Bailey,
36 2009). In information theory, entropy, also known as Shannon’s entropy, is defined as the
37 average minimum rate at which information is produced or predicted in an uncertain
38 stochastic setting (Shannon, 1948). In recent decades, information theory has been vastly
39 applied in many fields of science (Andrews *et al*, 2015). Biology is of no exception, but
40 compared to other areas, the applications of information theory in biological sciences have
41 been indeed limited (Battail, 2013). More importantly, medical sciences lack any use of
42 information theory in daily practice. The applications of information theory in molecular
43 biology have been mostly focused on genome sequence analysis (Vinga, 2013). To date, no
44 study has investigated the evolutionary nature of human diseases using information theory.

45 The backbone of evolution is random genetic mutations being selected according to the
46 natural environment. So what has been encountered in nature after some 3.5 billion years of
47 life history is a ‘selected randomness’. This is the reason why we believe information theory
48 can be a perfect language to understand life – i.e., this selected randomness. In the literature,
49 single nucleotide polymorphisms (SNPs) accounting for the main portion of this randomness
50 have been associated with inherited disease susceptibility (Bodmer & Bonilla, 2008; Wang &
51 Moul, 2001). However, such approaches have only focused on the genome investigation and,
52 in most part, neglected the human proteome and the protein-protein interactions (PPIs). PPIs
53 are the leading cause of cellular metabolic processes. They are induced-fit physical contacts
54 between macromolecules of proteins allowing the cellular function (Changeux & Edelstein,
55 2011; Keskin *et al*, 2008; Koshland Jr, 1995). In order to employ information theory in

56 medical sciences, it would be necessary to investigate diseases in detail considering their
57 molecular networks and PPIs. We believe evolutionary evidence interpreted by stochastic
58 information analysis can provide substantial help in understanding diseases of living
59 organisms.

60 In this study, to understand the nature of human diseases, we have focused on human
61 interactome and available proteomes of living organisms. To avoid confusion, by ‘human
62 diseases’ we only refer to non-communicable diseases in human with at least one reported
63 genetic basis. Also, as the term ‘proteome’ has sometimes referred to all proteins of a cell or
64 a tissue in the literature, it is to be noted that in this article, the term ‘proteome’ will refer to
65 the complete set of proteins that *can be* expressed by an *organism*. Because proteomes are
66 functional representatives of the ‘expressed genome’ of organisms, we have used them as the
67 means of our investigation. We have used Shannon’s entropy as a retrograde approach to
68 trace ~180 million proteins with more than 61 billion amino acids through the tree of life and
69 investigated the trends of complexity among organisms. We have shown that this
70 methodology agrees with the classification of phyla and may be used as a new tool in
71 taxonomy. Also, using our new mathematical theorem presented in the Materials and
72 Methods section, we have focused on *Homo sapiens*’ PPI network and discussed potential
73 clinical applications in the practice of medicine. We argue why there are only the diseases we
74 know, and not others, and discuss why some diseases are more prevalent. We also elaborate
75 on the reasonable links between our mathematical theory, Shannon’s entropy, the evolution
76 of taxa, and human diseases.

77 **Results**

78 *CAIR comparisons among taxonomic groups*

79 Calculated Average Information per Residue (CAIR) was calculated (see the Materials and
80 Methods section) for all proteomes available at the UniProt database until April 2020. Nearly
81 180 million proteins with more than 61 billion amino acids were analysed to classify ~29,000
82 organisms in 92 phyla. Table 1 shows CAIRs of the most popular proteomes and model
83 organisms (for all ~29k proteomes, see Dataset EV1). The minimum CAIR of an organism is
84 that of *Zinderia insecticola* (0.8247) and the maximum is of *Ciona savignyi* (0.9449). The
85 mean \pm standard deviation considering all organisms is 0.9210 ± 0.0129 with a median
86 (interquartile range) of 0.9250 (0.0160). Having performed a literature review of articles
87 published no later than April 2020, we have drawn the most updated tree of life for UniProt
88 taxonomic lineage data (Fig 1A). For each bifurcation point on the tree, we tested if two sides
89 of the bifurcation have developed divergent CAIRs. Fig 1B illustrates how CAIR divergence
90 is present through the different lineages of taxonomy. On all bifurcation points of Fig 1A, a
91 number is written whose respective statistical test results are demonstrated in Fig 1B with two
92 half violin plots for upper and lower sides of the bifurcation and their box-and-whisker plots.
93 Since the groups were negatively skewed, unbalanced, and heteroscedastic, their difference
94 was investigated via the two-sided Brunner-Munzel statistical test (Neuhäuser & Ruxton,
95 2009). It is noteworthy that groups with ten or fewer organisms were excluded from
96 comparisons, as the Brunner-Munzel test is statistically imprecise even with a permutation.
97 Among 56 performed tests, 48 tests demonstrated a significant difference at the point of
98 bifurcation. Interestingly, the bifurcation points of eight insignificant tests are mostly known
99 to be a matter of controversy in the scientific literature (Spang *et al*, 2017; Evans *et al*, 2019).

100 Along with the p -value significance, estimated effect sizes (ES) and 95% confidence intervals
101 (CI) are also reported. For the exact p -values of each test, please see Table 2.

102 *CAIR as a means of understanding the behaviour of natural selection*

103 The extent of natural selection's capacity to show bias in favour of selecting a spectrum of
104 organisms is open to question. Since genetic mutations are known to be generally random, the
105 CAIR density plot of such a random condition without a selection bias shall turn out to have a
106 uniform distribution. In a simulation of various random protein systems, we obtained a
107 similar distribution to the actual CAIR density plot taking into account a negative skewness
108 of -0.90 (Fig EV1). Fig EV1A shows the density plot of life and Fig EV1B depicts that of our
109 simulation. Fig EV1C, also, shows how these two distributions are similar given a tiny
110 bandwidth. To test for their similarity of distributions, we also performed two-sample
111 Kolmogorov-Smirnov tests whose mean p -value was 0.40 on 1000 iterations. It is
112 noteworthy, not unexpected though, that the natural selection is biased toward the organisms
113 with higher CAIRs. This might have stemmed from the random mutations over the course of
114 ages as discussed in the next section. Also, natural selection favours more complex and more
115 unpredictable protein systems, as they can accommodate superior functionalities. Besides, the
116 density plot of CAIRs possesses one other interesting property. Since there are significant
117 differences among taxonomic groups of the second hierarchy as shown in Fig 1B, it can
118 further be expected that all organisms are noticeable on the density plot. Fig 2 shows how
119 different taxonomic hierarchies are manifested in the density plot of organism CAIRs. On the
120 'wave of life', members of the succeeding taxonomic ranks are revealed by zooming in on the
121 preceding taxonomic group. This property of CAIR density suggests an original methodology
122 to help classifying organisms into various taxa.

123 *Human proteome analysis and the estimation of mutual information for a protein (EMIP)*

124 According to the theorem presented in the Materials and Methods section and its biological
125 inferences, EMIP has been calculated for each human protein entry using its PPI network.
126 Each Swiss-Prot, i.e. reviewed, entry has been categorized into three groups using the
127 Orphanet database of diseases. In case of a reported disease(s) related to an entry, all disease
128 point-prevalence/incidences of the entry (from Orphanet database) are summed up to obtain
129 the total occurrence of the disease, that is to say, the protein's overall malfunction. Table 3
130 shows the narrative data of the human disease categories. As seen in the table, the groups are
131 unbalanced in size and heteroscedastic which makes conventional statistical analyses
132 unfavourable. Herein, our results demonstrate how well indicators of diseases can be
133 correlated to the disease occurrence categories. In the Materials and Methods section, we
134 have explained why such independent variables were candidates of correlation and further
135 statistical analyses. Fig 3 shows the results of comparisons between disease occurrence
136 categories in four disease indicators. In each comparison, we have also included gene age
137 categories to test our hypotheses and biological inferences. The Dunnett-Tukey-Kramer
138 pairwise multiple comparison test adjusted for unequal variances and unequal sample sizes
139 (Dunnett, 1980) was performed to test overall comparisons. Results of comparisons show that
140 disease indicators correlate significantly with disease occurrence categories. Among four
141 indicators, EMIP is revealed to have the most significant differences between categories,
142 while CAIR was incongruous. The inconsistency seen in CAIR confirms that the use of
143 Shannon's entropy alone is not a good enough indicator of disease occurrence categories.
144 Additionally, since gene ages are presented as ranked data in the literature (Liebeskind *et al*,
145 2016), ranked analysis with an equal number of ranks has been performed to make all five
146 indicators comparable with one another. Fig 4 shows the Likert plots and rank comparisons.
147 As illustrated, natural Logarithm of EMIP (LEMIP) is by far better than other indicators in
148 correlating disease occurrence categories which might be stemmed from its bell-shaped

149 histogram. Unlike other indicators, the distribution of LEMIP allows it to be ranked with
150 mean and standard deviations which have been reported in the Table 4. Please refer to Table
151 4 for detailed information about the ranks in disease indicators.

152 As results suggest, EMIP seems to be a superior indicator of human diseases. Calculated
153 values of disease indicators for all reported diseases (Dataset EV2) and all human proteins
154 (Dataset EV3) are available in the Expanded View of the article. In a nutshell, we have
155 presented 16 human proteins with the highest EMIPs in Table 5. It is noteworthy that high-
156 EMIP proteins are more susceptible to have diseased networks and are clinically crucial for
157 human health. Fig EV2 shows the network topology of the same proteins (for its R code, see
158 the Data Availability section).

159 Discussion

160 *Proteome evolution during ages*

161 As shown in the Materials and Methods section, relative frequencies of residues in a protein
162 decide on its CAIR. Obviously, biochemical properties of amino acids play a central role in
163 determining their primary relative frequencies in *de novo* proteins. Such dissimilarity of
164 chemical properties would encourage unbalanced primary relative frequencies, thus lesser
165 CAIRs, as shown in our results for younger proteins in Fig 3D. This finding follows the
166 theory of *de novo* gene birth from non-coding DNA (Neme *et al*, 2017; Wilson *et al*, 2017).
167 Nevertheless, during the course of evolution, residues are subjected to random mutation
168 which equalizes their relative frequencies. This, in turn, increases the CAIR of proteins as
169 they age which also agrees with the trend of CAIRs in Fig 3D. This can be a corroborating
170 rationale for the study carried out with a different methodology in which it is shown that
171 intrinsic disorder of proteins negatively correlates with gene age (Banerjee & Chakraborty,
172 2017). Random mutations aside, natural selection's bias in favour of more complex proteins
173 may have also contributed to increasing the CAIR in older proteins. This is also noticeable
174 from the results seen in Fig EV1 verifying the identical behaviour of natural selection toward
175 all living organizations. Accordingly, it is not surprising that the human proteome includes a
176 negatively-skewed distribution whose lesser CAIRs are mostly associated with proteins
177 expressed by younger genes. Unfortunately, the literature lacks any thorough investigation on
178 linguistic complexity of proteins and gene ages.

179 Moreover, it is evident from Fig 3C that the younger eukaryotic proteins are shorter than their
180 older counterparts. A significant decline is noticed, however, in the length, interactions, and
181 mutual information of proteins during the old ages. This refers to the different evolutionary
182 rates of prokaryotic and eukaryotic genes and is compatible with the findings of previous

183 studies (Alba & Castresana, 2005; Wolf *et al*, 2009). Even for proteins of prokaryotic ages,
184 the trend of increase in protein length and interactions is observed; nonetheless, the trend line
185 is discrete from that of eukaryotic proteins. Interestingly, this decline is not seen in Fig 3D, as
186 the CAIR, in both eukaryotic and prokaryotic settings, is identically affected by directional
187 selection. The dome shape increase of disease indicators from younger to older proteins
188 refutes the justification of the study by Elhaik *et al* claiming that the slower rate of evolution
189 in older genes is ‘an artefact of increased genetic distance (Elhaik *et al*, 2006)’.

190 Furthermore, it is demonstrated in Fig 3B that interactions increase as genes age which agrees
191 with the literature (Saeed & Deane, 2006). However, the rate of increase seems to be slower
192 in a prokaryotic setting. That means the rate of interaction turnover in eukaryotic proteomes
193 may be comparatively higher. This might be due to the denser networks of eukaryotes and the
194 gene duplication (Wagner, 2003). It is also noteworthy that the rate of interaction turnover
195 seems to be non-decreasing as eukaryotic genes get older. Lastly, mutual information might
196 be considered as an assembly of protein information and interactions. So as seen in Fig 3A,
197 the trend of EMIP is also increasing by time for both eukaryotic and prokaryotic genes.

198 *Dynamics of protein networks*

199 A protein network is considered ‘stable’ when the odds of network malfunction is tiny.
200 Among various protein networks, particular ones malfunction often, and they generally are
201 responsible for the networks of non-communicable diseases. According to the theorem
202 presented in the Materials and Methods section, the number of interactions is negatively
203 correlated to network stability. This deduction is contrary to what is seen in Fig 3B, 4b, and
204 the literature (Jonsson & Bates, 2006; Oti *et al*, 2006; Xu & Li, 2006). According to the
205 literature, the PPI networks of disease genes are different in topology containing more
206 interactions, as it has been similarly shown in the mentioned figure panels. Previously, it was

207 clarified that the interactions increase with a stationary rate as a gene ages. As a matter of
208 fact, there are also many genes encoding for proteins of the cellular organisms that have not
209 established any interactions. Hence, a confounding factor that might be overlooked would be
210 the primary stability of the protein itself. In other words, proteins that are prone to
211 malfunction need substantial interactions to compensate for their malfunctions within their
212 network. That is the reason for increasing interactions along with the disease occurrences. In
213 the literature, the finding of the increased mutation rate of disease genes may reason their
214 increasing interactions (Smith & Eyre-Walker, 2003). Natural selection might be the other
215 reason which would favour interactions merely for faulty networks rejecting in others.

216 As mentioned previously, an inaccurate estimate of a protein's stability would be its CAIR.
217 Fig 3D illustrates that CAIRs of disease proteins are significantly higher than those of non-
218 disease proteins. However, the incongruous decrease of CAIRs in rare disease comparing
219 with extremely rare diseases had not expected in the inferences of our theorem. This might
220 be, in part, the influence of abundant younger proteins responsible for rare diseases that have
221 not developed interactions yet. Proteins with less CAIR are more stable, yet they are
222 negatively selected as they do not contain sufficient information. Complexity allows a protein
223 to have the potential capacity to carry more intricate functions (Babu, 2016). Thus, in order to
224 obtain higher functionalities, proteins grow both in their sizes and CAIRs. Consequently, this
225 necessitates new interactions to arrive. However, new protein interactors may take millions of
226 years to appear and reform the instability of the network (Fraser *et al*, 2002). These cycles
227 begin all over again as the new proteins come into existence. All these aside, the essentiality
228 of a protein's function is an argument that should not be dismissed. Generally, old proteins
229 are more crucial to forming life than younger ones (Chen *et al*, 2012). The less crucial roles
230 of younger proteins render their diseases to be less severe. As human ages, numerous
231 interactors in various networks are cancelled out causing these networks to malfunction. The

232 coincidental rush of diseases at the late ages of human life may be caused by the
233 accumulative effect of interactors removal and the loss of proteostasis (Kaushik & Cuervo,
234 2015; Labbadia & Morimoto, 2015). So, even substantial interactions cannot fully guarantee
235 networks of complex unstable proteins.

236 *Gene age, hubs, and human diseases*

237 Our results presented in Fig 4C illustrates the critical role of proteins encoded by the older
238 genes to be responsible for a broad spectrum of diseases. This finding is in total agreement
239 with a previous study in the literature highlighting the importance of ancient genes in human
240 genetic disorders (Domazet-Lošo & Tautz, 2008). We also discussed that the older genes
241 might take the leading role in creating the necessary fundamentals of life. A series of papers
242 by Barabási et al dedicatedly demonstrate the distinction between disease genes and the
243 essential genes. According to their work, disease genes are in most cases non-essentials being
244 located at the periphery of the network, rather than being a central hub (Barabási *et al*, 2011;
245 Domazet-Lošo & Tautz, 2008; Goh *et al*, 2007). In Fig 4D, although the trend of the increase
246 in CAIR is parallel to more odds of disease to occur, the proportion of rare to extremely rare
247 diseases suggests that more prevalent diseases are not associated with the proteins with the
248 highest CAIRs. This is also evident from the overall comparison of CAIRs in disease
249 occurrence categories in Fig 3D. This finding puts forward an argument that the most
250 complex proteins located at the centre of networks are encoded by the essential old genes that
251 in case of their malfunction, the condition would be fatal or cause an extremely rare and
252 severe disease. However, the more prevalent diseases are caused by comparably less complex
253 proteins that are indeed younger than central hubs. So, the proportion of rare to extremely
254 rare diseases in Fig 4 is of great importance and should be noticed as they agree with the
255 scenarios presented by Barabási et al.

256 *EMIP and human diseases*

257 A disease may occur when the process of network compensation works inaccurately. For any
258 network, the removal of the nodes would weaken its stability. This fact is evidently derived
259 from our theorem. As we discussed, it is the reason why many of the unstable networks have
260 urged to have many interactions. Integrating the effect of interactions with the protein
261 information is what mutual information tries to demonstrate. In a sense, the mutual
262 information of a protein is the information that is identical between the protein and its
263 network. Thus, in the case of a malfunction, it would not result in a disease, as the network
264 already carries the identical information. Estimation of mutual information, mathematically,
265 is equivalent to the difference between the information of the network as a whole and the
266 scalar summation of information that interactors carry when they are not interacting within a
267 network. This, of course, will show how much capacity the network carries to compensate for
268 its malfunction. In Fig 3A and Fig 4A the superior relation between EMIP and disease
269 occurrence categories stems from this fact.

270 *On the existence of diseases*

271 Based on what is discussed above, the scenario of a gene to cause a disease or not is being
272 summarized as the following (Fig 5). When the evolution was in its initial periods, the
273 involved proteins for life network was mainly the crucial ones. The metabolic pillars of life
274 owe the fundamental and critical components to the first formation of these networks. There
275 are two types of proteins at this stage, i.e. either 'robust' or 'weak'. To define, robust proteins
276 are those which are structurally not very susceptible to malfunction. These proteins have
277 evolved without many interactions comparing to others. Hence, during the path of evolution,
278 and still, we cannot detect as many interactions for them. On the other hand, weak proteins
279 would have caused vital errors that result in severe diseases. It would be reasonable to assume
280 that the incidence of such diseases would have been higher at ancient ages. So we may expect

281 to observe many interactions of them till now. According to our theorem, it can be reasoned
282 that these interactions would bypass the hub protein in case it malfunctions. So the incidence
283 of such diseases would have been drastically lessened by now. As a general point, for an
284 unstable weak protein that initiates a network, there are two prospects to be considered, i.e.
285 being able to develop a mature network at the time of the investigation, or still getting
286 involved in an immature one. By definition, mature networks would be those which have
287 totally cancelled out the adverse outcomes of the malfunctioning hub protein. Based on the
288 theorem, maturity is an ideal that no network can satisfy it in a biology setting.

289 Having mentioned all these about the archaic proteins, it should also be noticed that the
290 proteins which have come into existence in more chronologically proximal periods would
291 have by far lower chances of being a vital hub. The functioning network of life, in one piece,
292 has less to do with a mammalian protein than an archaic cellular respiration protein. So, the
293 other category of proteins is that of the contemporary period on which the logical assumption
294 would be that they are either non-hubs or if hubs do not function in places crucial to life. The
295 contemporary category, i.e. the young proteins, again, is subdivided into robust and weak
296 proteins. Robust young proteins are those without interactions which are not very much
297 susceptible to cause diseases, because if otherwise, evolution would bring interactions for
298 them. It is noteworthy that very young proteins are mainly very stable, simple, and proteins
299 with minor functional capacities. Stable young proteins may by time change to unstable
300 complex proteins because of the random mutations and the natural selection's bias toward
301 more complex proteins. In the way of transformation, most of the weak young proteins arise
302 which are primarily responsible for the prevalent metabolic diseases of the current
303 evolutionary era. Since the average evolutionary rate interaction turnover has not satisfied the
304 optimum number of network nodes for them, they are apprentices susceptible to malfunction.
305 It would be easy to infer that these proteins cause less severe diseases comparing to archaic

306 proteins, as their functions are still not vital. Nevertheless, they cause diseases with higher
307 incidences.

308 *Applicability of the method*

309 The methodology we have presented in the next section is not sensitive to the level of
310 taxonomy, i.e. whether the calculation is for a species, a genus, an order, a kingdom, or the
311 complete tree of life. The reason for this fact is that calculating the amino acid frequencies is
312 the same considering a species proteome or any other more extensive proteomic combination
313 of taxonomic levels. Also, we can calculate the Shannon entropy for a single protein, or a
314 peptide. This insensitivity to the size of the network in calculation enables a homogenous
315 analysis through the whole tree of life.

316 *A perspective of future studies*

317 Future studies may focus on each of the non-communicable diseases to elaborate more on the
318 speculations made in this study. Communicable diseases may also be the focus of further
319 studies to investigate the CAIRs of organisms and probable relations to their pathogenicity.
320 Treatments that are targeting networks with high-CAIR interacting protein crowds would be
321 an option to be explored. Moreover, better estimations for mutual information would be of
322 great interest. Besides, the notation of the wave of life and CAIR comparisons may bring
323 further arguments to taxonomists. Lastly, the theorem may be used in various fields of
324 science which are shaped by networks.

325 **Materials and Methods**

326 *Introducing the Calculated Average Information per Residue (CAIR) and the Protein* 327 *Information (PI)*

328 Proteins, from a mathematical point of view, are once randomly-occurred sequences of
329 residues that have gone through a process of selection in nature. Considering this fact, a
330 protein structure can be defined using a random sequence that carries mathematical
331 information. The information that a protein carries may be defined as to be equivalent to the
332 amount of uncertainty in predicting its residues. It is to be highlighted that regardless of the
333 protein conformation, the information of a protein is determined merely by its primary
334 structure. The average information carried by a residue in a protein is calculated by
335 Shannon's entropy (H) equation as below:

$$H = - \sum_{i=1}^s p_i \log_2 p_i \quad (1)$$

336 where p_i is the probability of state i , and s is the total number of possible states. In the
337 current context, the CAIR notion is introduced to be the same as Shannon's entropy except
338 for the logarithm base which is 22 in the former and 2 in the latter. In other words, CAIR is
339 the 22-ary of Shannon's entropy and is formulated as:

$$\text{CAIR} = - \sum_{r=1}^t p_r \log_{22} p_r \quad (2)$$

340 in which r is a numeral given to each residue, t is the total number of residues, p_r is the
341 relative frequency of r^{th} residue in the protein. More simply, the CAIR could be written as:

$$\text{CAIR} = kH \quad (3)$$

342 in which H is Shannon's entropy in equation (1), and k is a constant equivalent to:

$$k = \log_{22} 2 \cong 0.224243$$

343 As it is evident from equation (3), the patterns in the results obtained in the current article
344 were independent of the base of the logarithm; however, the scale of entropy would be much
345 more tangible considering the base of 22, as the ideal proteinogenic alphabet contains 22
346 letters. Deriving from CAIR, protein information (PI) is the amount of information carried by
347 all the residues of a protein as of the following equation:

$$\begin{aligned} \text{PI} \\ &= -l \sum_{r=1}^t p_r \log_{22} p_r \end{aligned} \quad (4)$$

348 in which the variables are the same as those in equation (2), and l is the length of the protein.
349 The notations of PI and CAIR are proposed, instead of the conventional H , for the fact that
350 they are more expressive and pertinent for the field of proteomics. It would also be humbly
351 proposed – with an analogy to Shannon’s bit – to use ‘pit’, i.e. protein unit, for the unit of
352 CAIR, PI, and EMIP in order to be readily comprehensible. As an example, one kilo-pit
353 would be equivalent to the PI of a 1000-lengthed protein whose residues have equal
354 frequencies.

355

356 *Mathematical Theorem*

357 Suppose $\{X\}, \{Y_1\}, \{Y_2\} \dots, \{Y_n\}$ are sets of sequences from which $\{Y_1\}, \{Y_2\} \dots, \{Y_n\}$ are all
358 dependent to $\{X\}$, but are pairwise independent from each other, not necessarily identically
359 distributed random variables having characteristic functions of $\varphi_1, \varphi_2, \dots, \varphi_n$, distribution
360 functions of f_1, f_2, \dots, f_n , and entropies of H_1, H_2, \dots, H_n . Let Φ be the characteristic function
361 of $\{X\}$, F_x be its distribution function, and H_x be its entropy. Then:

$$\lim_{n \rightarrow \infty} H(\Phi | \varphi_1, \varphi_2, \dots, \varphi_n) = 0$$

362

363 **Proof.** According to the definition of mutual information, we can write the following

364 relations (Cover & Thomas, 2012):

$$\begin{aligned} I(X; Y_i) &= I(Y_i; X) \\ I(Y_i; X) &= H(X) - H(X|Y_i) \\ \sum_{i=1}^{\infty} H(X|Y_i) &= H(X) - \sum_{i=1}^{\infty} I(X; Y_i) \end{aligned} \quad (5)$$

365 Corollary. Non-negativity of mutual information (Cover & Thomas, 2012):

$$I(X; Y) \geq 0$$

366 with equality iff X and Y are independent.

367 Based on the corollary of non-negativity of mutual information, and because Y_i are all

368 dependent on X , the mutual information of X with respect to all Y_i is always positive:

$$I(X; Y_i) > 0 \quad (6)$$

369 from which it can be inferred that infinite sum of positive values yields not to infinite, but to

370 the maximum mutual information possible, i.e., the entropy of X :

$$\sum_{i=1}^{\infty} I(X; Y_i) = H(X) \quad (7)$$

371 So, substituting equation (7) in equation (5):

$$\sum_{i=1}^{\infty} H(X|Y_i) = H(X) - H(X)$$

$$\sum_{i=1}^{\infty} H(X|Y_i) = 0 \quad (8)$$

372 Also, since $\sum_{i=1}^{\infty} H(X|Y_i)$ is the same as $\lim_{n \rightarrow \infty} H(\Phi|\varphi_1, \varphi_2, \dots, \varphi_n)$, thus:

$$\lim_{n \rightarrow \infty} H(\Phi|\varphi_1, \varphi_2, \dots, \varphi_n) = 0 \quad (9)$$

373 *Biological inferences and hypotheses*

374 In the above theorem, supposing Φ be a protein with the amino acid sequence of $\{X\}$,
375 interacting with n number of other proteins, namely $\varphi_1, \varphi_2, \dots, \varphi_n$, with sequences of
376 $\{Y_1\}, \{Y_2\} \dots, \{Y_n\}$:

- 377 1) Said interactions are mathematically interpreted as the dependency of F_x on
378 f_1, f_2, \dots, f_n . To elucidate, the probability distribution functions of proteins
379 correspond to their Boltzmann distributions. Because of the induced-fit nature of
380 biochemical interactions, it might be plausible to consider the distribution
381 functions to be dependent on one another. For that reason, each interaction is
382 deduced as the dependency of two distributions in the theorem.
- 383 2) As the theorem suggests, H_x indicates Shannon's entropy of the protein Φ with an
384 amino acid sequence of $\{X\}$. This might be interpreted as the extent of probable
385 variations that could lay in the primary structure of a protein affecting its function.
386 This measure would be an inaccurate estimate of a protein's malfunction as no
387 biochemical conditions have been taken into account. Despite its inaccuracy, we
388 have included the CAIR as an indicator of diseases in our analysis. The intuitive
389 hypothesis would maintain that the CAIR is significantly more in disease proteins
390 than non-disease ones for their surplus odds of having potential disadvantageous
391 variations in the primary structure leading to malfunction and disease.

- 392 3) According to the assumptions of the theorem, the notation of ‘information’ could
393 also apply to any system containing proteins, i.e. a particular metabolic pathway,
394 an interactome, diseasome, or the entire living organisms. Measuring the
395 information is independent of the size of the network. This property allows us to
396 calculate the information among the different taxonomic hierarchies and to utilize
397 the results in practice. Also, in a single organism like *Homo sapiens*, it would
398 allow us to compare different disease pathways, track hub proteins, and discern
399 potential disease proteins from non-disease ones.
- 400 4) The length of a protein sequence determines the PI, as shown in equation (4). We
401 have not included the PI itself as an indicator in the study, but have included both
402 the CAIR and the protein length separately. Proteins with longer sequences carry
403 more information and are more prone to malfunction as the overall odds of a
404 faulty residue in their sequence is higher compared to a protein with a shorter
405 sequence.
- 406 5) Considering equation (9), it is understandable that the conditional entropy of Φ
407 with respect to the knowledge of n number of interactors, i.e. $\varphi_1, \varphi_2, \dots, \varphi_n$, would
408 equal zero when n approaches infinity. The conditional entropy designates the
409 new information carried by the Φ protein when functioning in its network with
410 other proteins of φ_1 to φ_n . So, as a preliminary and naive inference, it could be
411 easily inferred from the theorem that networks with more interactions are more
412 stable. Therefore, we have included the number of interactions as an indicator of
413 human diseases to test our hypothesis.
- 414 6) Although equation (9) is a relatively straightforward approach to calculate the
415 stability of a network, it can be shown that its exact quantitative calculation is not
416 possible in proteomic analysis. An equivalent measure of network stability would

417 be to use equation (6) in which mutual information is shown to be always positive.
 418 Unlike the conditional entropy which is negatively correlated to the network
 419 stability, the mutual information is positively correlated. To quantitate mutual
 420 information in the current context, we propose the following estimation,
 421 henceforth referred to as EMIP (μ):

$$\begin{aligned} \mu_{\Phi} = & -\left[\sum_{\varphi=1}^n l_{\varphi} + l_{\Phi}\right] \sum_{\varphi=1}^n \sum_{r=1}^t \frac{p_{(r,\varphi)} l_{\varphi} + p_{(r,\Phi)} l_{\Phi}}{l_{\varphi} + l_{\Phi}} \log_{22} \frac{p_{(r,\varphi)} l_{\varphi} + p_{(r,\Phi)} l_{\Phi}}{l_{\varphi} + l_{\Phi}} \\ & + \sum_{\varphi=1}^n \sum_{r=1}^t l_{\varphi} p_{(r,\varphi)} \log_{22} p_{(r,\varphi)} \end{aligned} \quad (10)$$

422 in which μ_{Φ} is the mutual information for the Φ protein, n is the number of
 423 interactions, l_{φ} is the length of the φ^{th} interactor, l_{Φ} is the length of protein Φ , r is
 424 a numeral given to each residue in proteins, t is the total number of residues,
 425 $p_{(r,\varphi)}$ is the relative frequency of r^{th} residue in φ^{th} interactor, and $p_{(r,\Phi)}$ is the
 426 relative frequency of r^{th} residue in Φ protein. EMIP has also been included as an
 427 indicator of diseases in our analysis.

428 7) An essential element in the course of life evolution and human disease analysis is
 429 to consider the gene age of proteins. It is conceivable to hypothesize that the
 430 chronological data of genes can be very much associated with the indicators of
 431 human diseases. One reason is the additive effect of gene ages to introduce new
 432 interactions. The second reason is based on the assumption that natural selection
 433 can show bias towards proteins with a specific range of CAIRs. Also, the third
 434 reason is the possibility that the older proteins can grow to have longer sequences
 435 during evolution. Therefore, we have also included the gene ages as a covariate of
 436 human diseases in our analysis.

437 8) According to the theorem, it is also inferred that in an ideal condition with an
438 infinite number of interactors for a protein, the functionality of such a hub protein
439 reaches an infallible state, i.e. no disease would ever happen. This could also be
440 applied to other fields of science to reason why there is an order in complex
441 systems with innumerable components. The rate of decline in error in our theorem
442 as shown in equation (9) under random conditions is generally consistent with the
443 simple statistical rule of \sqrt{n} as proposed in physicist Schrödinger's 'what is life'
444 (Schrödinger, 1944).

445

446 *Protein database*

447 For taxonomic comparisons, all complete proteomes were extracted from the UniProt
448 (Consortium, 2019) FTP server, freely available at
449 ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/complete/. Both
450 SwissProt and TrEMBL files were downloaded in FASTA format. We calculated protein
451 entropies separately for each entry in both files including a total of more 61 billion amino
452 acids and merged them with every individual organism. Then, we used the calculations for
453 further statistical analyses.

454 Besides, for human proteins analyses, the complete list of *Homo sapiens* proteins was
455 downloaded directly from the UniProt website in a tab-separated (.tab) format containing the
456 following columns: 'Entry', 'Length', 'Sequence', 'Orphanet', and 'Involvement in disease'.
457 Data were updated on 22nd April 2020 with UniProt release 2020_02.

458

459 *Protein-protein interactions database*

460 Protein-protein interactions in human proteome were obtained from Protein InteraCtion
461 KnowLedgebasE (Gioutlakis *et al*, 2017) (PICKLE) meta-database, the release of 2.5.
462 PICKLE is a cross-checked integration of all available human PPIs included in BioGRID,
463 IntAct, HPRD, MINT, DIP databases. The default filter mode was selected to download
464 191,113 binary interactions among 16,418 UniProtKB/SwissProt entries. All interactions
465 were included in our study for further analysis.

466

467 *Taxonomy database*

468 Taxonomy data of organisms were also extracted from the UniProt FTP, the release of
469 2020_02. All organisms were included and matched according to their ‘OX’, i.e. organism
470 number used by UniProt and other databases. The evolutionary tree was plotted based on a
471 landmark study (Hug *et al*, 2016) by Hug *et al*, published in 2016, with a review of updates
472 (Cavalier-Smith *et al*, 2014; Eloë-Fadrosh *et al*, 2016; Hahnke *et al*, 2016; Kirkegaard *et al*,
473 2016; Munoz *et al*, 2016; Hamilton *et al*, 2016; Eme *et al*, 2017; Momper *et al*, 2017;
474 Jungbluth *et al*, 2017; Jay *et al*, 2018; Momper *et al*, 2018; Pavan *et al*, 2018; Cavalier-Smith
475 *et al*, 2018; Carr *et al*, 2019; Dombrowski *et al*, 2019; Ward *et al*, 2019; Carnevali *et al*,
476 2019; Martinez *et al*, 2019; Youssef *et al*, 2019; Wang *et al*, 2019; Zhou *et al*, 2020; Kevbrin
477 *et al*, 2020) since then until April 2020. All updates were added to the tree and were matched
478 with UniProt taxonomy data.

479

480 *Diseases database*

481 According to cross-references between UniProt and Orphanet, related epidemiological data
482 were downloaded and extracted from the Orphanet database (Weinreich *et al*, 2008) in XML
483 format. The file was then used in Python code for further analysis. Only diseases with at least
484 one reported worldwide occurrence were included. Normally, in the Orphanet database,

485 occurrences are reported from one or more of the following categories: annual incidence,
486 cases/families, lifetime prevalence, point prevalence, and prevalence at birth. In the case of
487 more than one reported occurrence, the priority of selection was for incidence, prevalence at
488 birth, point prevalence, respectively. Accumulative occurrence data, i.e. cases/families and
489 lifetime prevalence, were excluded from the analysis.

490

491 *Analysis of taxonomic hierarchies*

492 The following steps were carried out to implement the theorem over the taxonomic
493 hierarchies (Fig 6A):

- 494 1) TrEMBL and SwissProt FASTA files of complete proteomes were downloaded from
495 UniProt KnowledgeBase. A total number of ~180 million protein entries were
496 included.
- 497 2) The frequencies of all amino acid residues were calculated with regards to all protein
498 entries.
- 499 3) All proteins were grouped according to their organisms using organism IDs.
- 500 4) Organisms that are not proteomes were excluded. Duplicates are also removed.
- 501 5) Viruses are excluded from the study.
- 502 6) Shannon's entropy was calculated according to the residue frequencies for all
503 included non-virus proteomes.
- 504 7) Taxonomic data were downloaded directly from the UniProt website.
- 505 8) The most updated tree of life was drawn after a thorough review of the literature until
506 April 2020.
- 507 9) Organisms were grouped with respect to taxonomic hierarchies.
- 508 10) Organisms with unknown taxonomic lineage were excluded from the analysis.

509 11) Brunner-Munzel test was performed for every bifurcation node through the tree of
510 life. The level of significance was 0.05. Preparation of data and protein information
511 calculations were all executed in Python 3.8.0 (Van Rossum & Drake, 2009) using
512 NumPy (Oliphant, 2006), Pandas (McKinney & others, 2010), and Biopython (Cock
513 *et al*, 2009) libraries. Statistical analysis and violin plots were carried out using R-
514 3.6.0 (R Core Team, 2019) with brunnermunzel (Ara, 2020) and Plotly (Sievert,
515 2018) packages.

516

517 *Analysis of human disease proteins*

518 The following steps were carried out to implement EMIP and analyse disease occurrence
519 categories (Fig 6B):

- 520 1) The human proteome was downloaded from the UniProt database with the organism
521 ID of 9606. 20,350 of reviewed and 54,473 of unreviewed proteins were included in
522 the study.
- 523 2) CAIR was calculated for all human entries using the sequences of residues.
- 524 3) Protein-protein interactions were extracted from the PICKLE database as the default
525 UniProt normalized file with a total number of ~190k interactions.
- 526 4) Interactions were altered in order to match the UniProt 'interactions with' column.
527 This was done to keep the homogeneity of the data.
- 528 5) EMIP was then calculated for all entries with the help of the PPIs.
- 529 6) The unreviewed proteins were excluded after the calculations of disease indicators
530 because they have not been reported to cause any diseases.
- 531 7) Ordinal age categories of genes were merged to the file using the consensus data
532 article (Liebeskind *et al*, 2016).

- 533 8) Disease categories are ranked using R into three groups of no diseases, extremely rare
534 diseases, and rare diseases.
- 535 9) Disease indicators were also categorized into eight groups to draw Likert plots.
- 536 10) Statistical differences between groups were tested among variables with the DTK test.
537 Significance levels of the DTK test were added to the comparison graphs with stars.
538 The highest significance level was set to 0.05. Graphs of comparisons were plotted in
539 R with DTK (Lau, 2013) and ggpubr (Kassambara, 2019) packages. Likert plots were
540 plotted with HH (Heiberger, 2019) package. Networks were plotted using igraph
541 (Csardi & Nepusz, 2006) package.

542 **Data availability**

543 Data used for analysis is available as the following files. FASTA files of all Swiss-Prot and
544 TrEMBL entries are publicly available from UniProt's FTP server at
545 https://ftp.expasy.org/databases/uniprot/current_release/knowledgebase/complete/. Also, all
546 non-redundant proteomes could be downloaded from UniProt website:
547 [https://www.uniprot.org/proteomes/?query=redundant:no&format=tab&force=true&columns](https://www.uniprot.org/proteomes/?query=redundant:no&format=tab&force=true&columns=id,name,organism-id,lineage&compress=yes)
548 [=id,name,organism-id,lineage&compress=yes](https://www.uniprot.org/proteomes/?query=redundant:no&format=tab&force=true&columns=id,name,organism-id,lineage&compress=yes). Tab-separated format of human proteome data
549 used in our analysis is achievable from
550 [https://www.uniprot.org/uniprot/?query=proteome:UP000005640&format=tab&force=true&](https://www.uniprot.org/uniprot/?query=proteome:UP000005640&format=tab&force=true&columns=id,reviewed,genes(PREFERRED),protein%20names,sequence,database(Orphanet),comment(INVOLVEMENT%20IN%20DISEASE),interactor&compress=yes)
551 [columns=id,reviewed,genes\(PREFERRED\),protein%20names,sequence,database\(Orphanet\),](https://www.uniprot.org/uniprot/?query=proteome:UP000005640&format=tab&force=true&columns=id,reviewed,genes(PREFERRED),protein%20names,sequence,database(Orphanet),comment(INVOLVEMENT%20IN%20DISEASE),interactor&compress=yes)
552 [comment\(INVOLVEMENT%20IN%20DISEASE\),interactor&compress=yes](https://www.uniprot.org/uniprot/?query=proteome:UP000005640&format=tab&force=true&columns=id,reviewed,genes(PREFERRED),protein%20names,sequence,database(Orphanet),comment(INVOLVEMENT%20IN%20DISEASE),interactor&compress=yes). Additionally,
553 PICKLE interactions are freely available from
554 http://www.pickle.gr/Data/2.5/PICKLE2_5_UniProtNormalizedTabular-default.zip. Orphanet
555 data is also freely available from http://www.orphadata.org/data/xml/en_product9_prev.xml.
556 Data of gene ages are adopted from the Gene-Ages GitHub repository at
557 https://github.com/marcottelab/Gene-Ages/raw/master/Main/main_HUMAN.csv.
558 Supplementary information is available in the online version.

559 All Python and R codes necessary to reproduce all parts of the analysis and for the illustration
560 of the figures are available under the MIT license on our GitHub repository at
561 https://github.com/synaptic-proteolab/CAIR_EMIP, or the Zenodo link at
562 <https://zenodo.org/record/3970210>. For executing codes online on cloud servers, Google
563 Colab links are also available on the GitHub page. Python
564 (<https://www.python.org/downloads/>), Jupyter Notebook (<https://jupyter.org/install>), R

565 (<https://cran.r-project.org/>), and RStudio (<https://rstudio.com/products/rstudio/download/>)

566 are all freely available for the public.

567 **Expanded View** for this article is available online.

568

569 **Acknowledgments**

570 Authors would like to thank Dr Saeed Sadigh-Eteghad, PhD in neurosciences, Tabriz
571 University of Medical Sciences, and Dr Pedram Dindari, a PhD candidate in computer
572 sciences, University of Tabriz, for their assistance during the preparation of the manuscript.

573 No funding to declare.

574

575 **Author contributions**

576 FK proposed the mathematical theorem and its proof. SD gathered and handled the large data
577 and prepared the data for analysis using Python. FK reviewed the literature to illustrate the
578 tree of life. FK analysed the data using R. SD and FK prepared the codes for open-source
579 publishing. FK wrote the manuscript, and SD agreed with all sections. FK designed the
580 figures and SD prepared the tables of the manuscript.

581

582 **Conflict of interest**

583 The authors have no conflict of interest to disclose.

584 **References**

- 585 Alba MM & Castresana J (2005) Inverse relationship between evolutionary rate and age of
586 mammalian genes. *Mol. Biol. Evol.* **22**: 598–606
- 587 Andrews JG, Dimakis A, Dolecek L, Effros M, Medard M, Milenkovic O, Montanari A,
588 Vishwanath S, Yeh E & Berry R (2015) A perspective on future research directions in
589 information theory. *arXiv Prepr. arXiv1507.05941*
- 590 Ara T (2020) brunnermunzel: (Permuted) Brunner-Munzel Test. Available at: [https://cran.r-](https://cran.r-project.org/package=brunnermunzel)
591 [project.org/package=brunnermunzel](https://cran.r-project.org/package=brunnermunzel)
- 592 Babu MM (2016) The contribution of intrinsically disordered regions to protein function,
593 cellular complexity, and human disease. *Biochem. Soc. Trans.* **44**: 1185–1200
- 594 Bailey KD (2009) Entropy systems theory. *Syst. Sci. Cybern. Eolss Publ. Oxford, UK*: 152–
595 169
- 596 Banerjee S & Chakraborty S (2017) Protein intrinsic disorder negatively associates with gene
597 age in different eukaryotic lineages. *Mol. Biosyst.* **13**: 2044–2055
- 598 Barabási A-L, Gulbahce N & Loscalzo J (2011) Network medicine: a network-based
599 approach to human disease. *Nat. Rev. Genet.* **12**: 56–68
- 600 Battail G (2013) Biology needs information theory. *Biosemiotics* **6**: 77–103
- 601 Bodmer W & Bonilla C (2008) Common and rare variants in multifactorial susceptibility to
602 common diseases. *Nat. Genet.* **40**: 695
- 603 Carnevali PBM, Schulz F, Castelle CJ, Kantor RS, Shih PM, Sharon I, Santini JM, Olm MR,
604 Amano Y, Thomas BC & others (2019) Hydrogen-based metabolism as an ancestral trait
605 in lineages sibling to the Cyanobacteria. *Nat. Commun.* **10**: 1–15

- 606 Carr SA, Jungbluth SP, Eloë-Fadrosh EA, Stepanauskas R, Woyke T, Rappé MS & Orcutt
607 BN (2019) Carboxydrotrophy potential of uncultivated Hydrothermarchaeota from the
608 subseafloor crustal biosphere. *ISME J.* **13**: 1457–1468
- 609 Cavalier-Smith T, Chao EE & Lewis R (2018) Multigene phylogeny and cell evolution of
610 chromist infrakingdom Rhizaria: contrasting cell organisation of sister phyla Cercozoa
611 and Retaria. *Protoplasma* **255**: 1517–1574
- 612 Cavalier-Smith T, Chao EE, Snell EA, Berney C, Fiore-Donno AM & Lewis R (2014)
613 Multigene eukaryote phylogeny reveals the likely protozoan ancestors of opisthokonts
614 (animals, fungi, choanozoans) and Amoebozoa. *Mol. Phylogenet. Evol.* **81**: 71–85
- 615 Changeux J-P & Edelstein S (2011) Conformational selection or induced fit? 50 years of
616 debate resolved. *F1000 Biol. Rep.* **3**:
- 617 Chen W-H, Trachana K, Lercher MJ & Bork P (2012) Younger genes are less likely to be
618 essential than older genes, and duplicates are less likely to be essential than singletons of
619 the same age. *Mol. Biol. Evol.* **29**: 1703–1706
- 620 Clausius R (1864) *Abhandlungen über die mechanische Wärmetheorie* F. Vieweg
- 621 Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T,
622 Kauff F, Wilczynski B & others (2009) Biopython: freely available Python tools for
623 computational molecular biology and bioinformatics. *Bioinformatics* **25**: 1422–1423
- 624 Consortium U (2019) UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res.*
625 **47**: D506–D515
- 626 Cover TM & Thomas JA (2012) *Elements of information theory* John Wiley & Sons
- 627 Csardi G & Nepusz T (2006) The igraph software package for complex network research.

- 628 *InterJournal Complex Sy*: 1695 Available at: <http://igraph.org>
- 629 Domazet-Lošo T & Tautz D (2008) An ancient evolutionary origin of genes associated with
630 human genetic diseases. *Mol. Biol. Evol.* **25**: 2699–2707
- 631 Dombrowski N, Lee J-H, Williams TA, Offre P & Spang A (2019) Genomic diversity,
632 lifestyles and evolutionary origins of DPANN archaea. *FEMS Microbiol. Lett.* **366**:
633 fnz008
- 634 Dunnett CW (1980) Pairwise multiple comparisons in the unequal variance case. *J. Am. Stat.*
635 *Assoc.* **75**: 796–800
- 636 Elhaik E, Sabath N & Graur D (2006) The ‘inverse relationship between evolutionary rate
637 and age of mammalian genes’ is an artifact of increased genetic distance with rate of
638 evolution and time of divergence. *Mol. Biol. Evol.* **23**: 1–3
- 639 Eloë-Fadrosh EA, Paez-Espino D, Jarett J, Dunfield PF, Hedlund BP, Dekas AE, Grasby SE,
640 Brady AL, Dong H & Briggs BR (2016) Global metagenomic survey reveals a new
641 bacterial candidate phylum in geothermal springs. *Nat. Commun.* **7**: 10476
- 642 Eme L, Spang A, Lombard J, Stairs CW & Ettema TJG (2017) Archaea and the origin of
643 eukaryotes. *Nat. Rev. Microbiol.* **15**: 711
- 644 Evans PN, Boyd JA, Leu AO, Woodcroft BJ, Parks DH, Hugenholtz P & Tyson GW (2019)
645 An evolving view of methane metabolism in the Archaea. *Nat. Rev. Microbiol* **17**: 219–
646 232
- 647 Fraser HB, Hirsh AE, Steinmetz LM, Scharfe C & Feldman MW (2002) Evolutionary rate in
648 the protein interaction network. *Science (80-.)*. **296**: 750–752
- 649 Ghahramani Z (2006) Information theory. *Encycl. Cogn. Sci.*

- 650 Gioutlakis A, Klapa MI & Moschonas NK (2017) PICKLE 2.0: A human protein-protein
651 interaction meta-database employing data integration via genetic information ontology.
652 *PLoS One* **12**:
- 653 Goh K-I, Cusick ME, Valle D, Childs B, Vidal M & Barabási A-L (2007) The human disease
654 network. *Proc. Natl. Acad. Sci.* **104**: 8685–8690
- 655 Hahnke RL, Meier-Kolthoff JP, García-López M, Mukherjee S, Huntemann M, Ivanova NN,
656 Woyke T, Kyrpides NC, Klenk H-P & Göker M (2016) Genome-based taxonomic
657 classification of Bacteroidetes. *Front. Microbiol.* **7**: 2003
- 658 Hamilton TL, Bovee RJ, Sattin SR, Mohr W, Gilhooly III WP, Lyons TW, Pearson A &
659 Macalady JL (2016) Carbon and sulfur cycling below the chemocline in a meromictic
660 lake and the identification of a novel taxonomic lineage in the FCB superphylum,
661 *Candidatus Aegiribacteria*. *Front. Microbiol.* **7**: 598
- 662 Heiberger RM (2019) HH: Statistical Analysis and Data Display: Heiberger and Holland.
663 Available at: <https://cran.r-project.org/package=HH>
- 664 Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, Butterfield CN,
665 Hernsdorf AW, Amano Y & Ise K (2016) A new view of the tree of life. *Nat. Microbiol.*
666 **1**: 16048
- 667 Jay ZJ, Beam JP, Dlakić M, Rusch DB, Kozubal MA & Inskeep WP (2018) Marsarchaeota
668 are an aerobic archaeal lineage abundant in geothermal iron oxide microbial mats. *Nat.*
669 *Microbiol.* **3**: 732
- 670 Jonsson PF & Bates PA (2006) Global topological features of cancer proteins in the human
671 interactome. *Bioinformatics* **22**: 2291–2297
- 672 Jungbluth SP, Amend JP & Rappé MS (2017) Metagenome sequencing and 98 microbial

- 673 genomes from Juan de Fuca Ridge flank subsurface fluids. *Sci. data* **4**: 1–11
- 674 Kassambara A (2019) ggpubr: ‘ggplot2’ Based Publication Ready Plots. Available at:
675 <https://cran.r-project.org/package=ggpubr>
- 676 Kaushik S & Cuervo AM (2015) Proteostasis and aging. *Nat. Med.* **21**: 1406
- 677 Keskin O, Guroy A, Ma B & Nussinov R (2008) Principles of protein– protein interactions:
678 What are the preferred ways for proteins to interact? *Chem. Rev.* **108**: 1225–1244
- 679 Kevbrin V, Boltyanskaya Y, Grouzdev D, Koziyeva V, Park M & Cho J-C (2020)
680 *Natronospirillum operosum* gen. nov., sp. nov., a haloalkaliphilic satellite isolated from
681 decaying biomass of a laboratory culture of cyanobacterium *Geitlerinema* sp. and
682 proposal of *Natronospirillaceae* fam. nov., *Saccharospirillaceae* fam. nov. and Gynuell.
683 *Int. J. Syst. Evol. Microbiol.* **70**: 511–521
- 684 Kirkegaard RH, Dueholm MS, McIlroy SJ, Nierychlo M, Karst SM, Albertsen M & Nielsen
685 PH (2016) Genomic insights into members of the candidate phylum Hyd24-12 common
686 in mesophilic anaerobic digesters. *ISME J.* **10**: 2352
- 687 Koshland Jr DE (1995) The key–lock theory and the induced fit theory. *Angew. Chemie Int.*
688 *Ed. English* **33**: 2375–2378
- 689 Labbadia J & Morimoto RI (2015) The biology of proteostasis in aging and disease. *Annu.*
690 *Rev. Biochem.* **84**: 435–464
- 691 Lau MK (2013) DTK: Dunnett-Tukey-Kramer Pairwise Multiple Comparison Test Adjusted
692 for Unequal Variances and Unequal Sample Sizes. Available at: [https://cran.r-](https://cran.r-project.org/package=DTK)
693 [project.org/package=DTK](https://cran.r-project.org/package=DTK)
- 694 Liebeskind BJ, McWhite CD & Marcotte EM (2016) Towards consensus gene ages. *Genome*

- 695 *Biol. Evol.* **8**: 1812–1823
- 696 Martinez MA, Woodcroft BJ, Espinoza JCI, Zayed AA, Singleton CM, Boyd JA, Li Y-F,
697 Purvine S, Maughan H, Hodgkins SB & others (2019) Discovery and ecogenomic
698 context of a global *Caldiserica*-related phylum active in thawing permafrost, *Candidatus*
699 *Cryosericotia* phylum nov., Ca. *Cryosericia* class nov., Ca. *Cryosericales* ord. nov., Ca.
700 *Cryoseriaceae* fam. nov., comprising the four species *C. Syst. Appl. Microbiol.* **42**: 54–
701 66
- 702 McKinney W & others (2010) Data structures for statistical computing in python. In
703 *Proceedings of the 9th Python in Science Conference* pp 51–56.
- 704 Momper L, Jungbluth SP, Lee MD & Amend JP (2017) Energy and carbon metabolisms in a
705 deep terrestrial subsurface fluid microbial community. *ISME J.* **11**: 2319–2333
- 706 Momper LM, Aronson H & Amend JP (2018) Genomic description of ‘*Candidatus*
707 *Abyssobacteria*,’ a novel subsurface lineage within the candidate phylum
708 *Hydrogenedentes*. *Front. Microbiol.* **9**: 1993
- 709 Munoz R, Rosselló-Móra R & Amann R (2016) Revised phylogeny of Bacteroidetes and
710 proposal of sixteen new taxa and two new combinations including *Rhodothermaeota*
711 phyl. nov. *Syst. Appl. Microbiol.* **39**: 281–296
- 712 Neme R, Amador C, Yildirim B, McConnell E & Tautz D (2017) Random sequences are an
713 abundant source of bioactive RNAs or peptides. *Nat. Ecol. Evol.* **1**: 1–7
- 714 Neuhauser M & Ruxton GD (2009) Distribution-free two-sample comparisons in the case of
715 heterogeneous variances. *Behav. Ecol. Sociobiol.* **63**: 617–623
- 716 Oliphant TE (2006) *A guide to NumPy* Trelgol Publishing USA

- 717 Oti M, Snel B, Huynen MA & Brunner HG (2006) Predicting disease genes using protein–
718 protein interactions. *J. Med. Genet.* **43**: 691–698
- 719 Pavan ME, Pavan EE, Glaeser SP, Etchebehere C, Kämpfer P, Pettinari MJ & López NI
720 (2018) Proposal for a new classification of a deep branching bacterial phylogenetic
721 lineage: transfer of *Coprothermobacter proteolyticus* and *Coprothermobacter platensis* to
722 *Coprothermobacteraceae* fam. nov., within *Coprothermobacterales* ord. nov.,
723 *Coprothermobacte*. *Int. J. Syst. Evol. Microbiol.* **68**: 1627–1632
- 724 R Core Team (2019) R: A Language and Environment for Statistical Computing. Available
725 at: <https://www.r-project.org>
- 726 Van Rossum G & Drake FL (2009) Python 3 Reference Manual Scotts Valley, CA:
727 CreateSpace
- 728 Saeed R & Deane CM (2006) Protein protein interactions, evolutionary rate, abundance and
729 age. *BMC Bioinformatics* **7**: 128
- 730 Schrödinger E (1944) What is life? The physical aspect of the living cell and mind
731 Cambridge University Press Cambridge
- 732 Shannon CE (1948) A mathematical theory of communication. *Bell Syst. Tech. J.* **27**: 379–
733 423
- 734 Sievert C (2018) plotly for R. Available at: <https://plotly-r.com>
- 735 Smith NGC & Eyre-Walker A (2003) Human disease genes: patterns and predictions. *Gene*
736 **318**: 169–175
- 737 Spang A, Caceres EF & Ettema TJG (2017) Genomic exploration of the diversity, ecology,
738 and evolution of the archaeal domain of life. *Science (80-.).* **357**:

- 739 Vinga S (2013) Information theory applications for biological sequence analysis. *Brief.*
740 *Bioinform.* **15**: 376–389
- 741 Wagner A (2003) How the global structure of protein interaction networks evolves. *Proc. R.*
742 *Soc. London. Ser. B Biol. Sci.* **270**: 457–466
- 743 Wang Y, Wegener G, Hou J, Wang F & Xiao X (2019) Expanding anaerobic alkane
744 metabolism in the domain of Archaea. *Nat. Microbiol.* **4**: 595–602
- 745 Wang Z & Moulton J (2001) SNPs, protein structure, and disease. *Hum. Mutat.* **17**: 263–270
- 746 Ward LM, Cardona T & Holland-Moritz H (2019) Evolutionary Implications of Anoxygenic
747 Phototrophy in the Bacterial Phylum Candidatus Palusbacterota (WPS-2). *BioRxiv*:
748 534180
- 749 Weinreich SS, Mangon R, Sikkens JJ, Teeuw ME & Cornel MC (2008) Orphanet: a
750 European database for rare diseases. *Ned. Tijdschr. Geneeskd.* **152**: 518–519
- 751 Wilson BA, Foy SG, Neme R & Masel J (2017) Young genes are highly disordered as
752 predicted by the preadaptation hypothesis of de novo gene birth. *Nat. Ecol. Evol.* **1**: 1–6
- 753 Wolf YI, Novichkov PS, Karev GP, Koonin E V & Lipman DJ (2009) The universal
754 distribution of evolutionary rates of genes and distinct characteristics of eukaryotic
755 genes of different apparent ages. *Proc. Natl. Acad. Sci.* **106**: 7273–7280
- 756 Xu J & Li Y (2006) Discovering disease-genes by topological features in human protein–
757 protein interaction network. *Bioinformatics* **22**: 2800–2805
- 758 Youssef NH, Farag IF, Hahn CR, Jarett J, Becraft E, Eloë-Fadrosh E, Lightfoot J, Bourgeois
759 A, Cole T, Ferrante S & others (2019) Genomic characterization of candidate division
760 LCP-89 reveals an atypical cell wall structure, microcompartment production, and dual

761 respiratory and fermentative capacities. *Appl. Environ. Microbiol.* **85**: e00110--19

762 Zhou Z, Liu Y, Xu W, Pan J, Luo Z-H & Li M (2020) Genome-and Community-Level

763 Interaction Insights into Carbon Utilization and Element Cycling Functions of

764 Hydrothermarchaeota in Hydrothermal Sediment. *mSystems* **5**:

765

766 **Figure legends**

767 **Figure 1. CAIR comparisons through the tree of life.**

768 **A** The most updated tree of life comprising all second hierarchies stemming from
769 cellular organisms. For bacteria superkingdom, the second hierarchy includes all 56 bacterial
770 phyla; Candidate Phyla Radiation (CPR) is included as one separate phylum. For Archaea
771 and Eukaryota superkingdoms, the second hierarchy encompasses 20 archaeal phyla and 16
772 eukaryote supergroups and divisions. On each bifurcation point of the tree, the arbitrary
773 number corresponds to the test number and the associated plots (B). The blue and light green
774 colours indicate the superior and inferior arm of the bifurcation point, respectively, which
775 correspond to the left and right sides of the violin- and box-and-whisker plots. The red colour
776 indicates a bifurcation point with both arms from which at least one arm contains less than
777 ten organisms and thus the Brunner-Munzel test would not be reliable. The numbers in
778 parentheses designate the total number of available complete proteomes in each group.

779 **B** Violin plots of each bifurcation point in the tree of life, except for those in red. The
780 vertical axes refer to the CAIR in all plots. The box-and-whiskers are overlaid within each
781 violin plot, and the white dashed line in each box indicates the CAIR mean in the
782 corresponding group of organisms. None of the outliers were excluded from the analysis.
783 Asterisks after each test number indicate the significance level of tests. ES; estimated effect
784 size. CI; 95% confidence interval.

785 **Figure 2. From the whole life to *E. coli* as an exemplary organism illustrating the**
786 **CAIR density of proteomes in several ranks of taxonomy.**

787 **A** The ‘wave of life’ denotes the CAIR density of proteomes through the tree of life. As
788 it is argued, the wave of life is proposed to be a summation of skewed distributions. As a
789 result, the wave of life holds an interesting property; namely, zooming in on the wave of life
790 by narrowing the range on the horizontal axis reveals the members in the next rank of
791 taxonomy. The peaks on the plots (A-G) have been named according to the most abundant
792 phyla with the closest median to the peak.

793 **B** CAIR density of the proteomes in the Proteobacteria phylum.

794 **C** Zooming in further on (B) and narrowing the range of horizontal axis to the
795 distribution of organisms under Gammaproteobacteria class, i.e. CAIRs of 0.900 – 0.936,
796 reveals the taxonomic orders.

797 **D-G** Proceeding further to zoom in on the peaks of the previous plot shows a perfect
798 agreement with all taxa lineage. Lastly, several strains of *E. Coli* form a bell-shaped
799 distribution (G).

800 **Figure 3. Comparisons of disease indicators in three occurrence groups across gene**
801 **age categories.**

802 **A** EMIP increases significantly as the occurrence of diseases increases. Generally, EMIP
803 has the highest level of significance among disease indicators. Not surprisingly, the trend of
804 EMIP is also increasing as the genes age. The decline of EMIP in primordial gene age eras
805 are due to the eukaryotic branching and the nucleic genetic material.

806 **B** The number of interactions is the second-best indicator of the disease occurrence
807 category. Similarly, as expected, the trend of interactions is increasing as genes grow older.
808 Evolution brings new interactions and adds new nodes to the network. Similarly, there is a
809 decline in the number of interactions in the last two gene age eras.

810 **C** The bigger the size of a protein, the more likely it is to be involved in a disease. Also,
811 it is noticeable that the gene ages correlates positively with the protein size in both eukaryotic
812 and prokaryotic settings. However, the disparity between these two settings is easily
813 discernible.

814 **D** Unlike what is expected, extremely rare diseases account for the proteins with the
815 highest CAIRs. This observation has been further elucidated in the Discussion section.
816 Nonetheless, the trend of gene ages agrees with the expectation as the complexity of proteins
817 increases with age. Error bars illustrate mean \pm 95% confidence intervals, and the
818 significance test is the Dunnett-Tukey-Kramer pairwise multiple comparison test adjusted for
819 unequal variances and unequal sample sizes.

820 **Figure 4. Likert plots of disease indicators (all classified into 8 ranks) with regard**
821 **to occurrence categories.**

822 **A** Ranked data of EMIP show a robust relationship with the disease categories. Note that
823 the occurrence of diseases increases as EMIP increases. Because the log-transformed of
824 EMIP (LEMIP) forms a bell-shaped curve, the ranking has been done by the mean and three
825 standard deviations of LEMIP.

826 **B** Ranked data of interactions increase with the occurrences, maintaining the order of
827 ranks. However, the number of interactions is not as well correlated to occurrences as EMIP.
828 It is noteworthy in (B), (D), and (E) panels, since indicator distributions are inconsistent with
829 a Gaussian distribution, rankings have been accomplished by median and equal percentile
830 intervals.

831 **C** The link between gene age ranks and disease categories is satisfactory considering the
832 first and last ranks; however, the overall order of ranks does not match with the order of
833 occurrence categories. Gene age ranks are the same as the eight age groups presented
834 previously in the literature (Liebeskind *et al*, 2016).

835 **D** Among Likert plots, CAIR ranks have shown the least correlation with the disease
836 categories which is in line with Fig 3D.

837 **E** The ranks of protein length are associated with disease occurrence categories
838 maintaining the order of ranks, except for the sixth rank.

839 **Figure 5. The cone of life and the evolution of diseases.**

840 The cone of life summarizes the probable scenarios for proteins put forward in the discussion
841 section. It is to be noted that the shape of the cone is schematic and an exemplary instance of
842 every possible occasion has been schematically illustrated. Also, for the clearness of the
843 drawing, the dark blue cylinder has been cut in half in order not to block other elements. The
844 diameter of the cone in any cross section shows the amount of existing protein material in
845 that given time period. The proteins are born from the circumference of any cone base. As
846 shown in the figure, a general rule would be that diseases are a result of weak proteins with
847 immature networks. Details have been depicted in the figure for every protein scenario. It is
848 to be highlighted that the eukaryotic and prokaryotic proteomes have not been discerned in
849 the figure. s, size; w/, with.

850 **Figure 6. Detailed steps needed to carry out the presented methodology.**

851 **A** Flowchart shows how the proteins were included, grouped according to their
852 respective organisms, and the exclusion criteria.

853 **B** A similar flowchart explaining the steps used to integrate human proteome data,
854 PICKLE interactions, and Orphanet diseases. Green rectangles, steps; purple rectangles,
855 executions; red flags, exclusions.

856 **Tables and their legends**

857 **Table 1. Proteome CAIRs of organisms mainly used in biological models and**
 858 **studies.**

Organism	CAIR	# of proteins	Organism	CAIR	# of proteins
<i>Arabidopsis thaliana</i>	0.9366	39,364	<i>Mus musculus</i> (Mouse)	0.9376	55,398
<i>Caenorhabditis elegans</i>	0.9419	26,850	<i>Neurospora crassa</i> [74A]	0.9323	10,257
<i>Chlamydomonas reinhardtii</i>	0.8887	18,829	<i>Rattus norvegicus</i> (Rat)	0.9378	29,951
<i>Ciona savignyi</i>	0.9444	20,004	<i>Saccharomyces cerevisiae</i> [S288c]	0.9336	6,049
<i>Danio rerio</i> (Zebrafish)	0.9398	46,848	<i>Schizosaccharomyces pombe</i>	0.9347	5,141
<i>Drosophila melanogaster</i> (Fruit fly)	0.9390	21,973	<i>Tetrahymena thermophila</i>	0.9119	26,976
<i>Escherichia coli</i> [K12]	0.9328	4,391	<i>Xenopus tropicalis</i> (Western clawed frog)	0.9410	55,258
<i>Homo sapiens</i> (Human)	0.9392	74,823	<i>Zea mays</i> (Maize)	0.9341	99,254
<i>Medicago truncatula</i> (Barrel medic)	0.9392	57,065	<i>Zinderia insecticola</i>	0.8247	206

859 The table has been sorted in alphabetical order. Where there were different proteomes of a
 860 single organism, the number of proteins refers to that of the most popular proteome used in
 861 the literature. Terms enclosed in parentheses are the common names used colloquially, and
 862 those enclosed in brackets are the strain names of organisms that have different indexed
 863 proteomes for their strains in the UniProt database. Please note that the prevalent model
 864 organisms of humans, like rats, mice, fruit flies, and zebrafish are very similar to us in terms
 865 of CAIR. #; number.

866 **Table 2. Exact p -values of the two-sided Brunner-Munzel tests.**

Test no.	p-value	Test no.	p-value	Test no.	p-value	Test no.	p-value
T1	5×10^{-137}	T18	2.72×10^{-10}	T42	0.000965	T66	4.29×10^{-08}
T2	3.42×10^{-17}	T19	0.14649	T44	0.35037	T67	0
T3	6.91×10^{-31}	T20	6.98×10^{-07}	T46	0	T68	0
T4	7.42×10^{-131}	T21	0.000417	T47	2.88×10^{-28}	T69	5.64×10^{-32}
T5	2.92×10^{-76}	T26	2.74×10^{-42}	T49	2.51×10^{-07}	T70	0.008022
T6	0	T28	1.66×10^{-05}	T51	7.23×10^{-205}	T73	0.004432
T7	6.41×10^{-47}	T29	0.99445	T54	0.000363	T75	2.71×10^{-263}
T8	0.32303	T32	8.86×10^{-17}	T55	5.04×10^{-07}	T76	7.62×10^{-08}
T9	0.019045	T33	0.14933	T56	5.13×10^{-05}	T77	9.92×10^{-05}
T10	6.98×10^{-08}	T36	3.83×10^{-09}	T57	8.48×10^{-56}	T79	1.30×10^{-08}
T14	7.70×10^{-08}	T38	0.012529	T59	0.87909	T80	7.16×10^{-05}
T15	0.001135	T39	0.030999	T60	0.46556	T81	0.000134
T16	0.010167	T40	5.65×10^{-16}	T62	0.55498	T83	6.78×10^{-89}
T17	2.45×10^{-17}	T41	2.73×10^{-51}	T64	0.002841	T88	1.35×10^{-34}

867 Test no. refer to the test labels illustrated in Fig 1. Because there was limited space in the
868 illustration, the exact p -values are reported herein. Please note that no test was performed
869 when either one or both groups contained less than 10 organisms. The numbers of available
870 organisms in each phylum are enclosed by parentheses in Fig. 1A. no.; number.

871 **Table 3. Narrative data of protein groups in three human disease categories.**

Disease category	Occurrence	Group size	# of interactions median (IQR)	Protein length median (IQR)
Rare diseases	> 1:1,000,000	819	16.0 (36.0)	599.0 (668.0)
Extremely rare diseases	≤1:1,000,000	1,356	11.0 (23.0)	527.0 (487.5)
No diseases	-	16,296	5.0 (15.0)	386.0 (393.0)

872 Swiss-Prot entries have been matched with their Orphanet cross-references to obtain the total
873 occurrences classifying diseases in three categories. Please note that extremely rare diseases
874 are a vast number of diseases whose names are not even familiar to general practitioners as
875 they mainly consist of case reports from around the world. Rare diseases are a group of
876 diseases that have been mainly the focus of attention in scientific literature and medical
877 books. IQR; interquartile range. #; number.

Disease indicator	Rank 1	Rank 2	Rank 3	Rank 4	Rank 5	Rank 6	Rank 7	Rank 8
LEMIP ^a	[-9, 36.45]	(36.44, 44.94]	(44.94, 53.44]	(68.38, 61.94]	(61.94, 70.44]	(70.44, 78.93]	(78.93, 87.43]	(87.43, 104.00]
Number of interactions ^b	[0]	[1]	(1, 3]	(3, 6]	(6, 11]	(11, 19]	(19, 39]	(39, 2136]
Protein length ^b	[2, 161]	(161, 250]	(250, 328]	(328, 415]	(415, 518]	(518, 671]	(671, 968]	(968, 34350]
CAIR ^b	[0.0217, 0.8803]	(0.8803, 0.8982]	(0.8982, 0.9088]	(0.9088, 0.9164]	(0.9164, 0.9230]	(0.9230, 0.9290]	(0.9290, 0.9351]	(0.9351, 0.9567]
Gene age ^c	Mammalia	Vertebrata	Eumetazoa	Opisthokonta	Eukaryota	Euk_Archaea	Euk+Bac	Cellular_organisms

878 **Table 4. Eight ranks of disease indicators and their quantitative intervals.**

879 a. Because of the bell-shaped curve distribution, mean and standard deviations have been used to classify ranks. b. Because of the non-normal
880 distribution, median and percentiles have been used to classify ranks. c. Ranks are based on scientific literature. LEMIP; Log_e of Estimation of
881 Mutual Information of Proteins. CAIR; Calculated Average Information per Residue.

Entry	Gene	Protein name	Length	CAIR	# of ints	EMIP	Disease
Q8WZ42	TTN	Titin (EC 2.7.11.1) (Connectin)	34350	0.9190	119	33839	Yes
P05067	APP	Amyloid-beta precursor protein (APP)	770	0.9309	2136	23441	Yes
P0CG48	UBC	Polyubiquitin-C [Cleaved into: Ubiquitin]	685	0.8785	916	14231	No
Q8WXI7	MUC16	Mucin-16 (MUC-16) (Ovarian cancer-related tumor marker CA125)	14507	0.8487	3	12440	No
Q9NRI5	DISC1	Disrupted in schizophrenia 1 protein	854	0.9002	650	11312	Yes
P04637	TP53	Cellular tumor antigen p53 (Tumor suppressor p53)	393	0.9250	787	10222	Yes
Q09472	EP300	Histone acetyltransferase p300 (p300 HAT)	2414	0.9145	541	9730	Yes
P00533	EGFR	Epidermal growth factor receptor (EC 2.7.10.1)	1210	0.9439	679	9697	Yes
P62993	GRB2	Growth factor receptor-bound protein 2 (Adapter protein GRB2)	217	0.9387	615	8704	No
Q8NF91	SYNE1	Nesprin-1 (Enaptin)	8797	0.9058	31	8583	Yes
P63104	YWHAZ	14-3-3 protein zeta/delta (Protein kinase C inhibitor protein 1)	245	0.9018	517	8282	No
P78362	SRPK2	SRSF protein kinase 2 (EC 2.7.11.1) (SFRS protein kinase 2)	688	0.9278	441	8161	No
Q03001	DST	Dystonin (230 kDa bullous pemphigoid antigen)	7570	0.9152	58	7806	Yes
Q5VST9	OBSCN	Obscurin (EC 2.7.11.1)	7968	0.9118	14	7769	Yes
P38398	BRCA1	Breast cancer type 1 susceptibility protein (EC 2.3.2.27)	1863	0.9160	456	7655	Yes
Q15149	PLEC	Plectin (PCN) (PLTN)	4684	0.8815	138	7611	Yes

882 **Table 5. Proteins with the highest EMIP values.**

883 Entry is the UniProt entry of the protein; Gene is the preferred gene names used in the literature; Length denotes to the length of the protein
884 sequence; # of ints is the number of interactions adapted from PICKLE database; Disease is represented as a dichotomous variable adapted from
885 UniProt's 'Involvement in disease' column.

886 **Expanded View Figure legends**

887 **Figure EV1. A simulation of natural selection shows the selection's bias towards**
888 **higher CAIRs.**

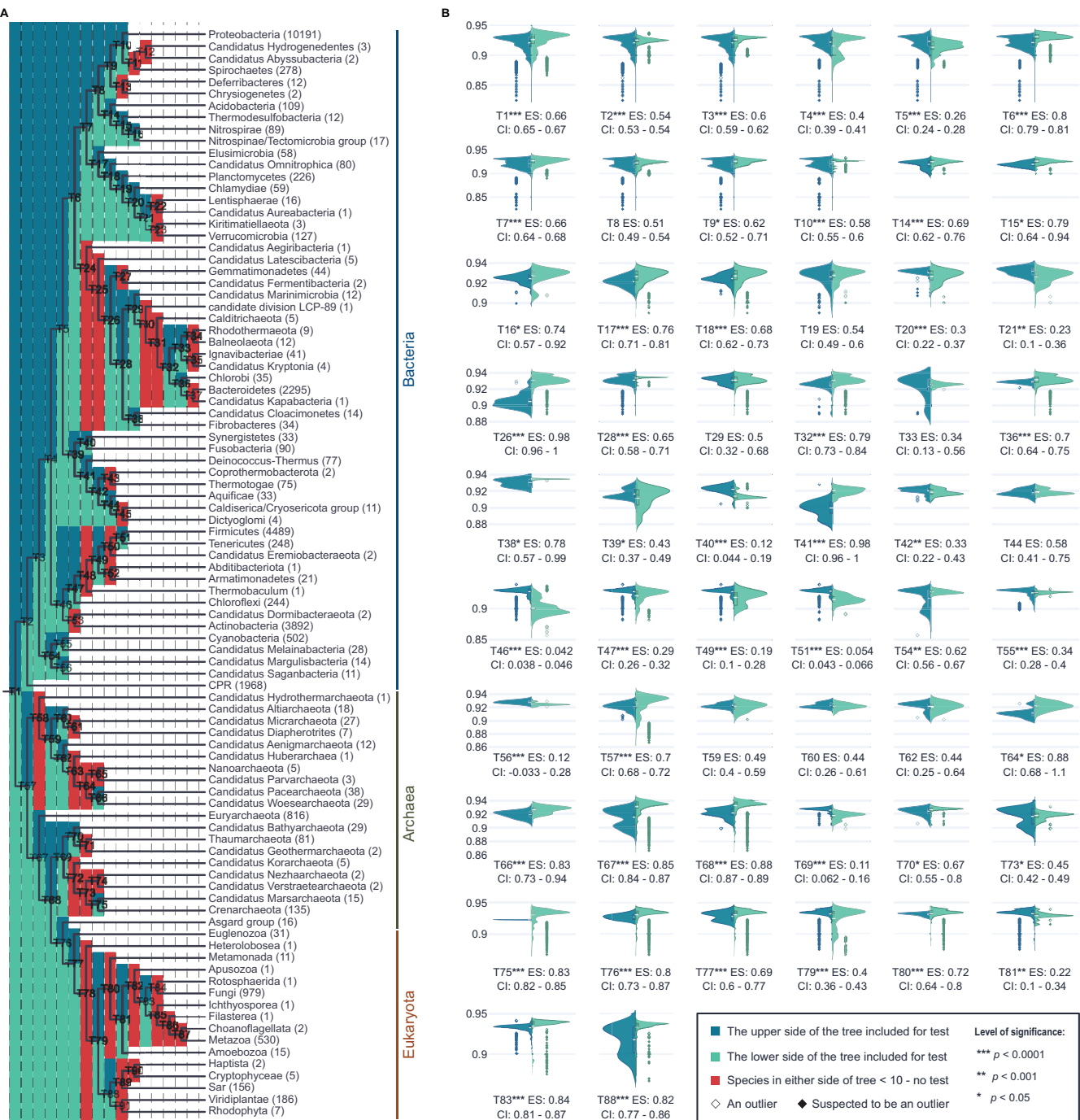
889 **A** CAIR density plot of the organisms considering the interquartile range (IQR) of phyla
890 sizes through the tree of life after removing the tiny phyla with less than 10 organisms (Q_1 - Q_3
891 $16.5 - 171$). From the 27 phyla within the IQR, random sampling was performed with a size
892 of $Q_1 \sim 17$. IQR inclusion and the subsequent sampling was done to remove the size effects of
893 populated phyla on the density plot. The mean of phyla skewness is -0.82 in Q_0 - Q_4 and -0.74
894 in Q_1 - Q_3 .

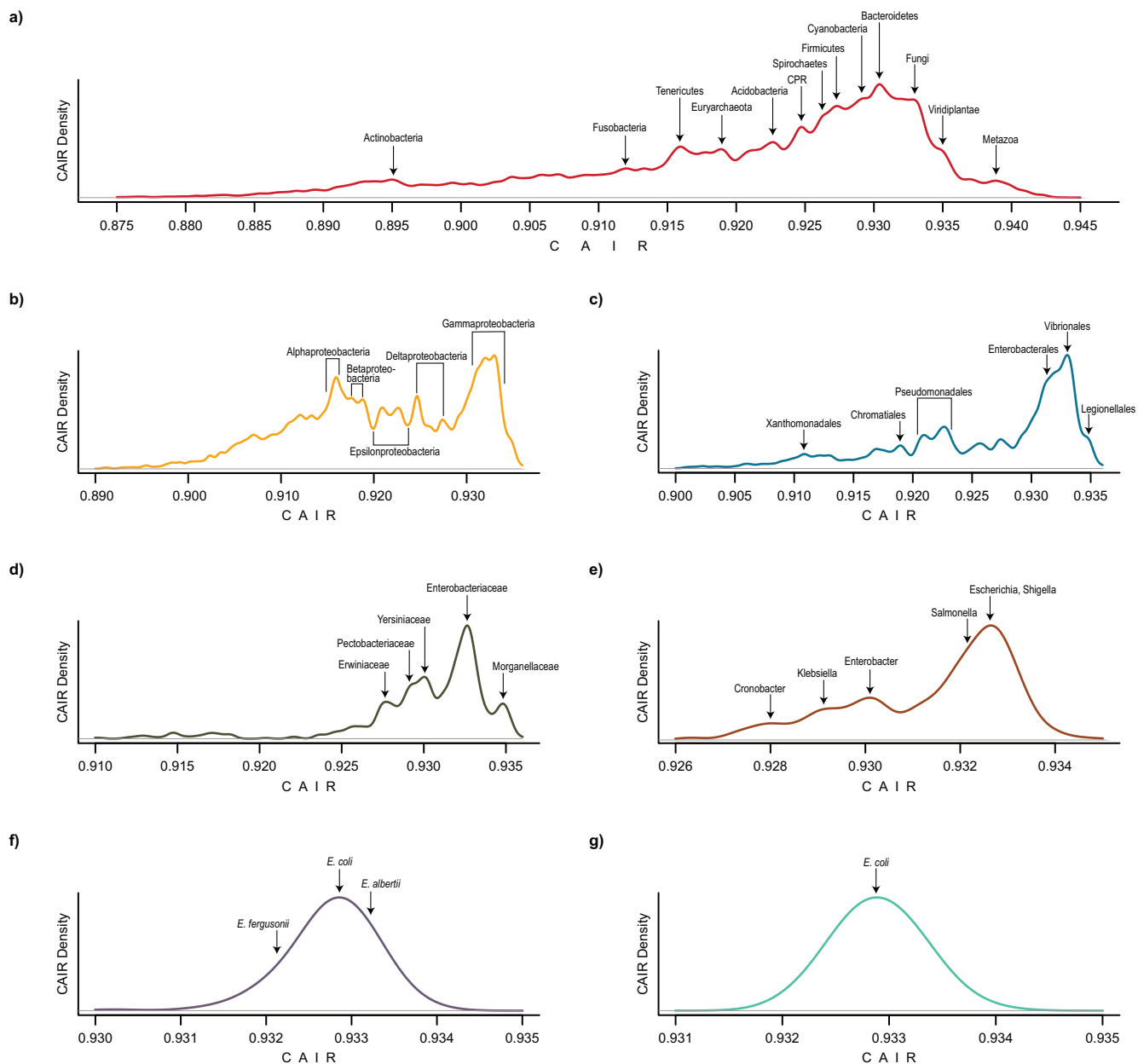
895 **B** CAIR simulation of the tree of life with 27 negatively skewed normal distributions.
896 The means of these simulated random distributions equal to the respective medians of the 27
897 phyla (marginal rugs) explained in (A). The function was iterated 1000 times to find the
898 skewness in which the Kolmogorov-Smirnov (KS) test has the maximum p -value. The
899 simulation revealed a skewness of -0.90 .

900 **C** Both (A) and (B) are overlaid with a lesser bandwidth to show the details of the
901 distributions. KS test reveals a p -value of 0.40 not rejecting the null hypothesis that
902 distributions are identical. Negative skewness of the red wave (depicting phyla data) suggests
903 that the natural selection is biased in favour of organisms with higher CAIRs. Q; quartile.

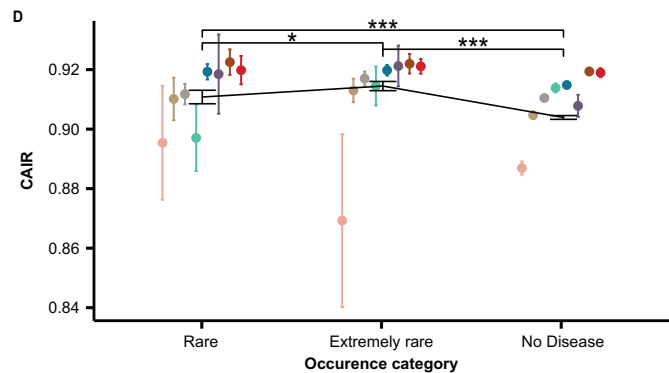
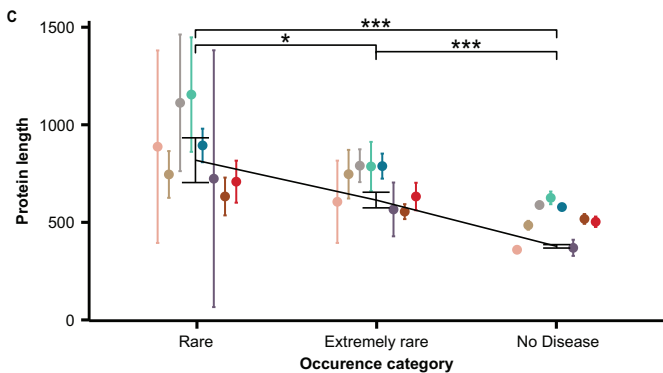
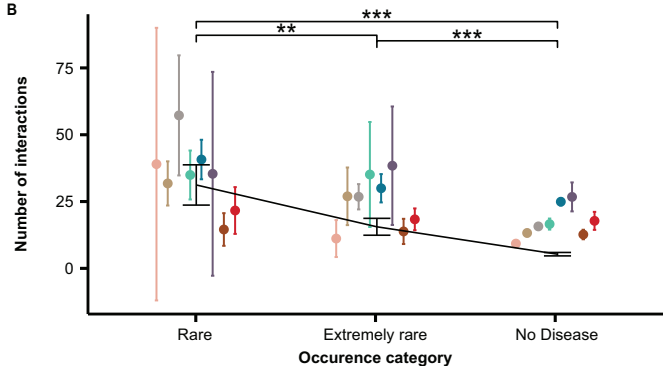
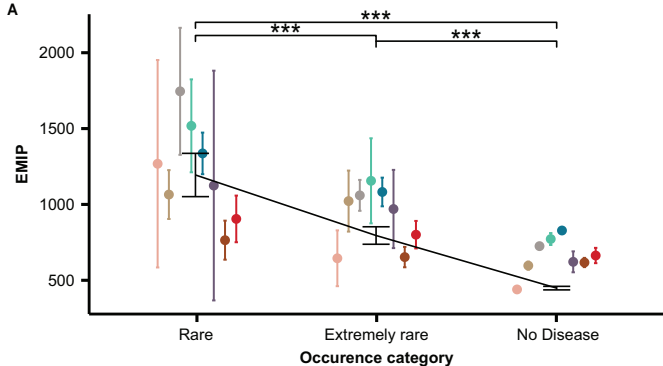
904 **Figure EV2. Protein networks of the proteins with the highest EMIPs.**

905 **A-P** The illustrated networks represent the proteins in Table 4 as the 16 proteins with the
906 highest EMIPs in the human proteome. The size of each circle represents the EMIP of the
907 protein. According to the UniProt database, (C), (D), (I), (K), and (L) are non-disease
908 networks and the rest are networks involved in at least one disease. The main proteins are
909 illustrated in colours other than sky blue, and all other interactors are coloured in sky blue.
910 The figure has been drawn with the interactions data available at the UniProt website up to
911 the second-degree interactions. It is noteworthy that the illustrated proteins with the highest
912 EMIP values are markedly present in various disease networks.





— All organisms (Wave of Life)
 — Proteobacteria phylum
 — Gammaproteobacteria class
 — Enterobacterales order
— Enterobacteriaceae family
 — Escherichia genus
— Escherichia coli species



Gene age categories:

- Mammalia
- Eumetazoa
- Eukaryota
- Euk+Bac
- Vertebrata
- Opisthokonta
- Euk_Archaea
- Cellular_organisms

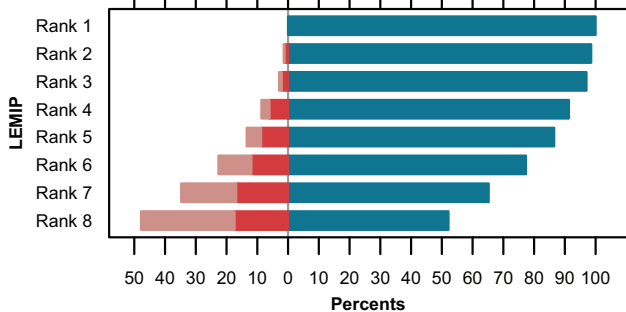
Overall comparisons:



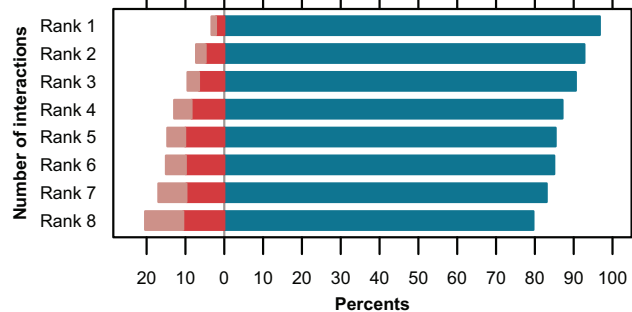
Significance levels:

- *** $p < 0.0001$
- ** $p < 0.001$
- * $p < 0.05$

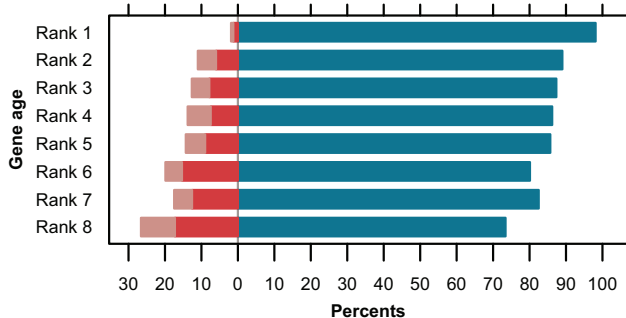
A



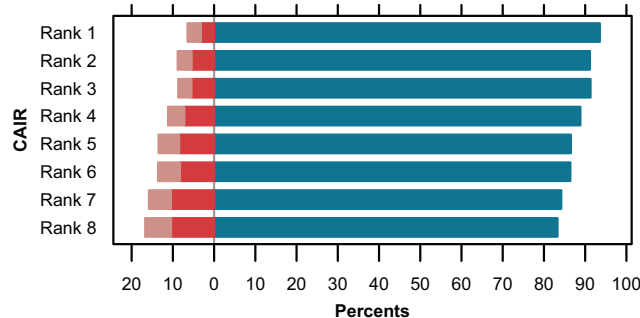
B



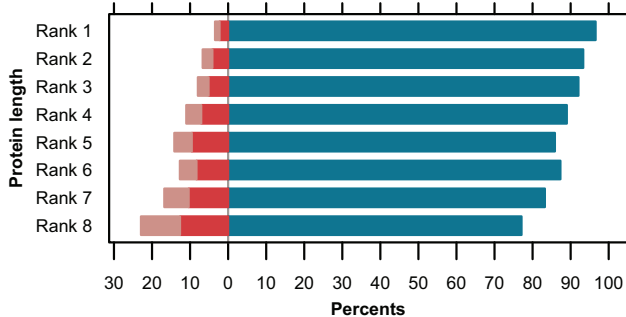
C



D



E

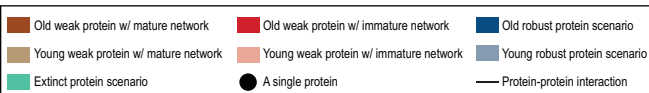
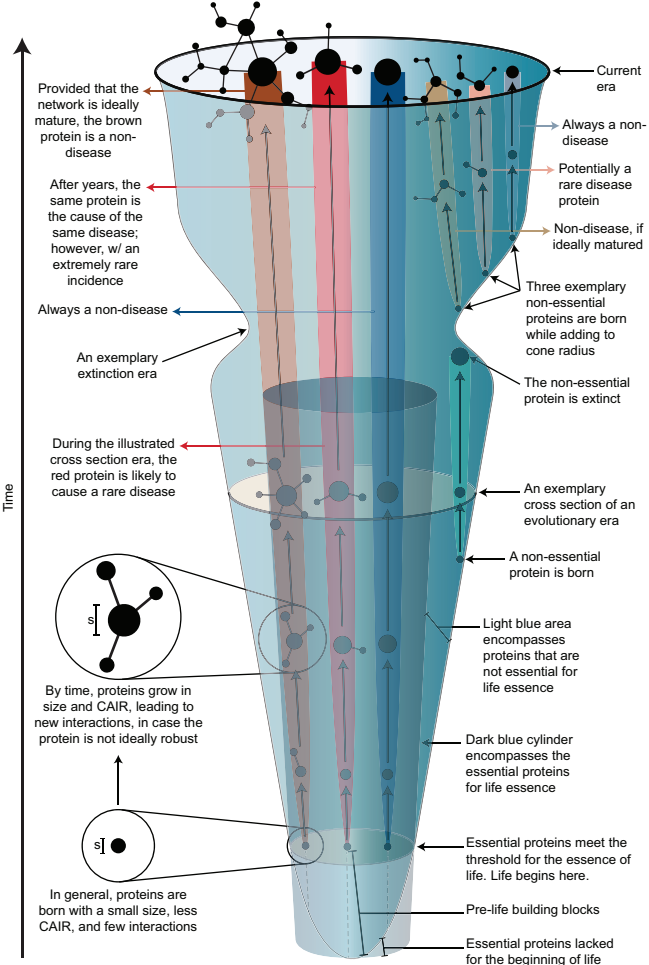


Disease categories:

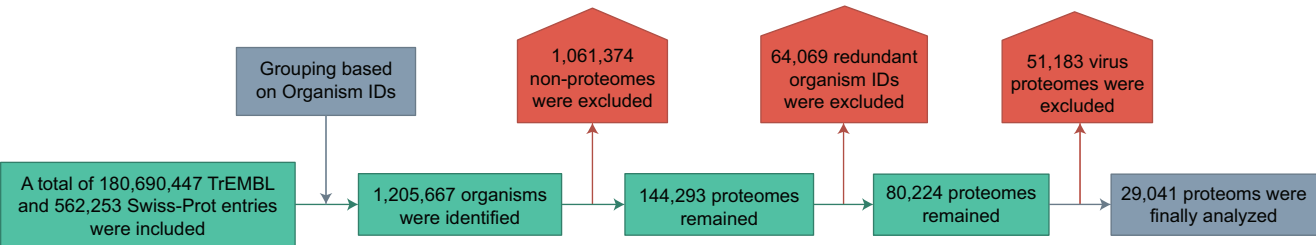
Rare

Extremely rare

No Disease



A



B

