

# Proteínas Canijas

Ale Zavala<sup>a</sup>, Sergio Hernández<sup>a</sup>, Pedro Miramontes<sup>a</sup>, León Martínez<sup>b</sup>

<sup>a</sup>*Facultad de Ciencias, National Autonomous University of Mexico, Mexico City 04510, Mexico*

<sup>b</sup>*Metropolitan Autonomous University*

---

## Abstract

¿Qué proponemos como abstract?

*Keywords:* Intrinsically disordered protein regions, ordered protein regions, Pfam family, protein domain, The discrete generalized beta distribution, megasecuence, familiome, self-organizing map, complexity

---

## 1. Background

Historically, the science of protein structure has privileged the study of ordered protein regions (OPRs) with regular and well-defined secondary structure, such as  $\alpha$ -helix,  $\beta$ -conformations and turns. In fact, proteins that contain mainly ordered regions are the classic text-book example of what a typical protein should look like. This is probably associated with the fact that the most commonly used technique for determining protein structure, i. e. x-ray crystallography, is biased towards resolving protein regions with regular and well-defined secondary structures. However, it has recently become increasingly clear that intrinsically disordered protein regions play an important role in the structure of proteins [1].

Intrinsically disordered protein regions (IDPRs), are segments of polypeptides, that do not natively acquire well defined, regular, or repetitive, secondary structures and instead adopt many different structural ensembles with no single, preferred lowest energy conformation and have biological activities that are dependent on this disordered state. There are some computational tools that can predict the order/disorder state of a protein region using only the complete polypeptide sequence as input [2], [3]. IDPRs have been found to be prevalent in eukaryotes and are common in bacteria and archaea too (XXX? ¿PEDRO: que significa esta nota?) [4] [5] [6] [1]. A general function

21 that has been proposed for many IDPRs is the binding to other molecules,  
22 including proteins, nucleic acids and other ligands [7] [8] [1]. However, it is  
23 possible that the importance of IDPRs goes beyond this currently accepted  
24 function, and that they may even play an important role in the folding and  
25 functioning of proteins. We would like to emphasize that in this study we  
26 have tried to jointly analyze both the OPRs and IDPRs. This approach  
27 derives from our belief that at least some of the properties of proteins are  
28 the product of the interactions between OPRs and IDPRs. Furthermore, we  
29 hypothesize that the distribution pattern of OPRs and IDPRs within a given  
30 protein domain is specific to that domain (similar to a fingerprint) and also  
31 specific to that domain’s family. We will expand on the definition of protein  
32 families and other related concepts in the following sections.

33 Here, we will only mention that studying patterns of protein order/disorder  
34 at the level protein domains (instead of single proteins) allowed us to cap-  
35 ture amino acid sequence variability that was, in principle, compatible with  
36 the tertiary structure conservation that characterizes a domain. Addition-  
37 ally, each representative of a domain acted as a new quasi-replicate thereof,  
38 allowing us to (naturally) increase the sample size of the data series to be  
39 analyzed with our proposed statistical tools (see below). Furthermore, we  
40 augmented our protein sequence data with gene ontology (GO) annotations  
41 in order to explore whether functional information of protein domains corre-  
42 lates with their order/disorder regularities.

43 Measures of complexity of patterns of OPRs and IDPRs have the poten-  
44 tial to contain information on structural and functional aspects of proteins in  
45 the sense that members of a given protein family that are close to each other  
46 in the parameter space of some measures of complexity could also be close in  
47 terms of functional properties. If IDPRs are an integral part of proteins the  
48 patterns of alternating OPRs and IDPRs of varying lengths should be rela-  
49 tively conserved among members of a protein family. Potentially, measures  
50 of the complexity of such patterns could capture enough information thereof,  
51 such that they could serve as a proxy for characterizing specific proteins.

52 In this study we will use several measures of complexity of a polypeptide  
53 pattern. At this time ”complexity” can be understood as how close or how  
54 far a sequences is from two extremes: randomness or periodicity. On one  
55 hand, the well known Kolmogorov’s complexity index (K) provides a measure  
56 of how far a sequence is from being random and, on the other, Shannon’s  
57 entropy (H) tells us how far is a sequence from the equiprobable distribution  
58 (REFSXXX). Additionally, the recently proposed Discrete Generalized Beta

59 function (DGBF) of rank-ordered distributions has also emerged as a measure  
60 of complexity. Interestingly, the DGBF can separate intermittent regimes  
61 from chaotic dynamics and may serve as an indicator of transitions between  
62 these two regimes in a wide array of phenomena [9].

63 Here, we propose the use of the above mentioned descriptors to place any  
64 given protein into a space of parameters and to discuss they way different  
65 families of proteins group in such a space. Our working model will be the Sac-  
66 charomyces cerevisiae proteome. We discuss the possibility that the DGBF of  
67 patterns of protein order/disorder may be able to detect evolutionary transi-  
68 tions in which ordered proteins acquire or expand their disordered regions or  
69 in which disordered protein start to limit their conformational repertoires.

70 • **Esta parte la va a “arreglar Pedro”.**

71 with affinities to other phenomenological descriptions such as power  
72 laws, that is able to capture the “flavor” of a given complex system and  
73 differentiate it from other similar yet distinguishable systems, while at  
74 the same time being robust to the finite size effects that often affect  
75 power law descriptions of physical phenomena.

76 Here, we propose to use a discrete version of the generalized beta function  
77 as well as Kolmogorov K and Shannon H as descriptors of the order-disorder  
78 patterns of protein domains present in the Saccharomyces cerevisiae “famil-  
79 iome” (a novel concept that we propose and that refers to the complement  
80 of Pfam families present in the baker’s yeast a proteome).

81 • **MÁS DE QUE ESPERAMOS R:.**

82 Idea 1. En este artículo utilizamos un método novedoso nunca antes  
83 aplicado en el área de la estructura de las proteínas denominado germi-  
84 beta. Nuestra hipótesis es que obtendremos una combinación de valores  
85 específica e independiente para cada familia a partir de la entropía, la  
86 complejidad algorítmica y la función beta característica de sus patrones  
87 de orden-desorden que esperamos pueda servir como un nuevo atributo  
88 para caracterizar a una familia de proteínas.

89 Si lo anterior no se cumple, idea 2: En este artículo utilizamos un  
90 método novedoso nunca antes aplicado en el área de la estructura de las  
91 proteínas denominado germibeta. Este método nos ayudará a entender  
92 (aquí no sabemos qué es exactamente lo que nos ayudará a entender)  
93 que junto con el proceso biológico de la estructura de las proteínas  
94 permiten aumentar el conocimiento de la función de las mismas en base  
95 al estudio de las RIDs/ROs. Palabras clave: describir, caracterizar,  
96 dominio, proteínas, propiedades de las familias o proteínas o dominios.

## 97 2. Materials and Methods

### 98 2.1. Construction of the database of complete sequences represented in the 99 yeast familiome

100 In this study, our aim is to characterize protein families in terms of  
101 the entropy, algorithmic complexity and characteristic beta function of their  
102 order-disorder pattern. A family of proteins is a group of proteins or pro-  
103 tein domains that share patterns of significant sequence conservation, due to  
104 common ancestry, that frequently entail functional similarity [10]. A protein  
105 domain is a substructure produced by any contiguous part of a polypeptide  
106 chain whose structural features are independent of the rest of the protein. A  
107 domain usually contains between 40 to 350 amino acids, and it is the modular  
108 unit from which many larger proteins are constructed [11]. It is important to  
109 note that any given protein does not necessarily belong to a single family, as  
110 a given protein can contain several domains with different evolutionary his-  
111 tories, and indeed, many proteins belong to several families [12]. Although  
112 most protein domains that are identified using sequence-based approaches  
113 are have well-defined and relatively stable spatial structures, some can be  
114 fully or largely disordered or can contain conserved disordered regions [13],  
115 these are known as intrinsically disordered domains (IDDs; [14]). The protein  
116 families information is provided by the Pfam database Sonnhammer et al.  
117 [15] [16] [17] [18]. Pfam is a collection of protein domains and protein fam-  
118 ilies in which each family is represented by two multiple sequence alignments  
119 and two profile-Hidden Markov Models, one of the two alignments is a high  
120 quality seed alignment [15] [16].

121 To build our yeast familiome database, we explored protein information  
122 in several biological databases. To every translated gene from *Saccharomyces*

Genome Database (v2015; <https://www.yeastgenome.org/>; [19]) we associated an UniProt identifier and its complete polypeptide sequence (“uniprot” full file, v2014; [20]). Using the UniProt information, we linked to every yeast protein its corresponding Pfam families (v28, 2015; <https://pfam.xfam.org/>; [12] [17]). Subsequently, for each Pfam family from *S. cerevisiae* famiome we downloaded its Pfam-A seed alignment file (v28, 2015; [15] [17] (XXXXXXXX ?) , which contains only the aligned segments that belong to a protein family of a variety of species. Finally, we used the UniProt identifier provided by the Pfam-A seed alignment file and the UniProt full file to obtain the complete polypeptide sequence of each protein in the Pfam-A seed alignments of the yeast famiome.

## 2.2. *Predicting intrinsically disordered residues in each polypeptide sequence of yeast famiome*

In order to assign each residue from our complete polypeptide sequence yeast famiome database to either the “ordered” or “intrinsically disordered” categories we used the open-source DisEMBL prediction software. DisEMBL is based on artificial neural networks trained to predict three different definitions of disorder and displays the disordered segments of arbitrary length within a protein sequence [21]. We used the three different algorithms of DisEMBL: loops/coils, hot loops and remark465 and the final assignment of each residue as either ordered or disordered residue was based on a majority rule decision between the three predictions. The order/disorder information was coded into the sequences as UPPERCASE/lowercase one-letter amino acid symbols, respectively. This procedure was performed for each complete polypeptide sequence of each Pfam-A seed alignment family in the yeast famiome. All the members in one family were concatenated together into one big megasequence.

## 2.3. *Transferring the ordered/disordered information to the Pfam-A families seed alignment*

In our study, we needed to associate sequences of the Pfam-A families seed alignment of the yeast famiome with the DisEMBL majority rule decision results. In order to do this we used options of the MAFFT program to maintain the Pfam-A families seed alignment unchanged, to maintain the gaps, the UPPERCASE/lowercase in the alignment and to preserve intact the order of the residues [22] [22] in the Pfam-A families seed alignment of the yeast famiome.

159 *2.4. Gene Ontology annotations in the yeast familiome*

160 To enrich our yeast familiome database, for those Pfam families where this  
161 information was available, we associated the molecular function annotation  
162 from Gene Ontology (v2018; <http://geneontology.org/>; [23] [20]). We needed  
163 a general molecular function annotation so we manually curated the specific  
164 GO molecular function annotations of the yeast familiome.

165 *2.5. Megasequence construction*

166 In order to have a sequence big enough to statistically represent the whole  
167 family in a robust way, we constructed what we called a megasequence which  
168 consisted in all the domain instances within a family glued together, so the  
169 statistical regularities will be magnified and easily observable in the so called  
170 megasequence.

171 We took all the aligned sequences for a given domain and spliced them  
172 one after another to obtain a family megasequence.

173 *2.6. The discrete generalized beta distribution (DGBD)*

174 All these megasequences of the yeast familiome were compared using a  
175 discrete generalized beta distribution (DGBD) [9] [24] which is a rank order-  
176 ing distribution that takes the form:

$$f(r) = \frac{A(N + 1 - r)b}{r^a}$$

177 Where a and b are parameters to be found, N is the number of ranks  
178 and A is a normalization constant. This rank ordering distribution has been  
179 successfully used across a wide range of different phenomena regardless of  
180 the truncated scaling behavior shown typically in most of the rank-order  
181 distributions. The approach used in our analysis is as follows. We took all  
182 the aligned sequences and merge them together one after another to make a  
183 family megasequence, then we counted the frequencies of the different words  
184 of length 2 and ordered these distributions of sizes in decreasing order. Then,  
185 through a nonlinear fit of (1) we obtained the (a,b) pair which was used to  
186 represent the distribution.

187 *2.7. Shannon Entropy ( $H(X)$ )*

188 Another attribute added to the whole yeast familiome was the calculation  
189 of Shannon's Entropy for each of the sequences of the larger. Shannon's  
190 Entropy is defined as follows:

$$H(X) = \sum \frac{p_i}{\log(p_i)} (2)$$

Where  $p_i$  is the probability of one of the  $N$  amino acids in the megasequence  $X$ . Applying  $H(X)$  to every family megasequence we have will reveal which sequences are the furthest from the normal distribution and so, which megasequence has the most structure in it [25].

## 2.8. Kolmogorov complexity

Kolmogorov complexity index when applied to a string of characters, in our case is the megasequence  $X$ , can be interpreted as the complexity of a computer program required to reproduce megasequence  $X$ . The calculation of Kolmogorov's complexity index can be approximated as follows:

$$k(seq) = \frac{\text{length}(\text{compressed}(seq))}{\text{length}(seq)} (3)$$

Where  $seq$  is the original megasequence of some family of proteins. The actual implementation of this function was done in Python computer language where zlib libraries were used to compress every sequence of the family [26]. In our experiment we used  $K$  as another attribute together with the already described, in order to understand the algorithmic complexity assumed to exist in every family. That is, if a set of instructions is behind the description of every protein in each family, we would expect that  $K$  captures this particular complexity.

## 2.9. Self-organizing map

The self-organizing map (SOM) is an unsupervised neural network used for data analysis and dimensionality reduction [27]. It has been long being applied into data analysis and biological sciences to detect similar profiles of analyzed data [28, 29]. The basic algorithm is divided in two steps. First, an initial map is formed with  $N$  neurons arranged in a lattice which will represent  $M$   $d$ -dimensional vectors. Each of the  $N$  neurons contains a single  $d$ -dimensional prototype that will be modified during the training of the map. Then, for each vector sample a winning prototype must be found. The second step consists in modifying the prototypes of all vectors within a neighborhood of this winning neuron, the magnitude of the modification is in proportion to the distance of the winning neuron. This process ends up unfolding a map where each prototype in neurons represents local averages

of data, hence nearby neurons have similar prototypes. Once the SOM is formed, locally grouped families must be assigned to a group. This step is done by a hierarchical clustering using euclidean distance. The number of groups was determined by the Davies-Bouldin index.

### 3. Results

#### 3.1. *Saccharomyces cerevisiae* familiome database

Our yeast familiome database contains 538 Pfam families (protein domains) and a total of NNNN instances of domains. The size range of the families goes from a family containing XXXX instances of a domain (Family number PF????) to a family containing YYYYY instances (Family PFLLLL). All instances of a given domain were concatenated to obtain a megasequence, thus yielding a total of 538 megasequences. Each megasequence contains information of ordered/disordered status for all its residues, as well as the values of the parameters  $\alpha$  and  $\beta$  from the expression of DGBD, Shannon's entropy, and Kolmogorov complexity. Additionally, for 260 families we have the specific GO molecular function corresponding to 18 different GO categories namely hydrolase, electron transfer activity, protein dimerization activity, isomerase, motor activity, transferase, transmembrane transporter, translation initiation factor activity, binding, translocases, antioxidant activity, structural constituent, oxidoreductase, ligase, structural molecule activity, catalytic activity, lyase, and copper chaperone activity. This information is shown in supplementary material S1-Yeast Familiome.

#### 3.2. *Discrete Generalized Beta Distribution to characterize S. cerevisiae familiome*

The Discrete Generalized Beta Distribution used in this work is a novel probability function that it has been shown to be an alternative to fit data that does not fit Zipf's law perfectly nevertheless an underlying process alike seems to be taking place [9, 30]. In some instances the  $\alpha$  exponent can be related to behaviors generating power laws, as is the case of scale invariance in turbulence in the so called inertial range where energy is transferred between different scales at the same rate, while  $\beta$  seems to be associated with chaotic, disordered fluctuations, for example the dissipative range for turbulence. In contrast with classical powerlaw like functions, the DGBD is able



254 to encompass both the scale invariance and chaotic regime in depicting the  
 255 whole process and its conflicting dynamics in the same graph. This gives a  
 256 general representation of the phenomena under study. [9].

257 There is a wide variety of a distinct phenomena studied under the DGBD  
 258 as shown in [9] and specifically there is some research in the field of genomics  
 259 as shown in [31, 30, 32] It is worth noting that the role of exponents  $\alpha$  and  $\beta$   
 260 as universality classifying parameters, as for example in critical phenomena,  
 261 remains be investigated in further detail.

262 We constructed DGBD plots for the 580 families and we are showing  
 263 the best DGBD plots according to the following criteria: first, a Pfam family  
 264 needed to have the highest square of correlation coefficient; second, an sample  
 265 size N of at least 30 different domain instances; and finally the DGBD alpha  
 266 and beta values had to be  $\geq 0$ . Every panel in each figure indicates the Pfam  
 267 family, the square of correlation coefficient, the alpha and beta values, and  
 268 its N value.

269 In the familiome database, the highest  $\alpha$  value was 2.0027 and the lowest  
 270 alpha value was 0.1003. Figure 1 shows selected cases for the scenario where  
 271  $\alpha > \beta$ .

272 In this database, we have the best 30 DGBD graphs with high alpha and  
 273 low beta values and there are 5 different general GO molecular functions be-  
 274 longing to 11 different Pfam families. Seven of these eleven have a “binding”  
 275 GO, 2 have “ligase” GO and “oxidoreductase”, “isomerase” and “lyase” on-  
 276 tologies have one each (Figure 1). Although there are seven general binding  
 277 ontologies, we cannot group them because their specific ontologies are dif-  
 278 ferent. We have two “protein binding” , one “ATP binding”, one “thiamine  
 279 binding”, one “DNA binding”, one “metal ion binding”, and one “phos-  
 280 phatidylinositol binding”. For ligase ontology, we have two different Pfam  
 281 families with specific ontologies like “aminoacyl-tRNA ligase” activity and  
 282 “aminoacyl-tRNA editing” activity. In the case of a family with lyase on-  
 283 tology, its specific ontology is “phosphatidylserine decarboxylase activity”.  
 284 For the families with oxidoreductase and isomerase, there are not specific  
 285 ontologies.

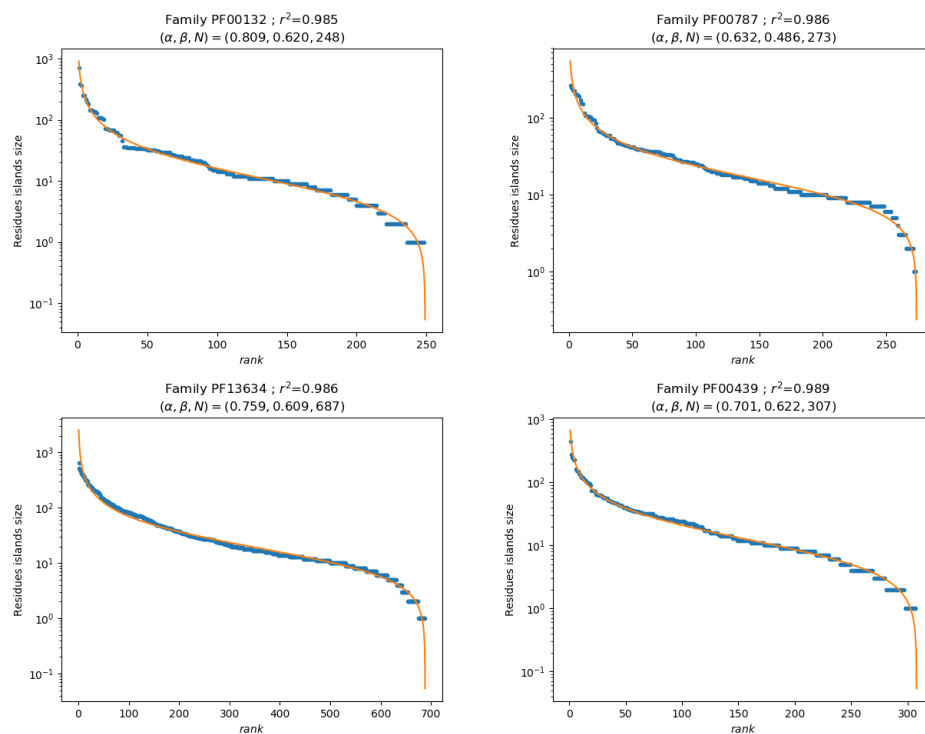


Figure 1: Semilog DGBD plots with high  $\alpha$  and low  $\beta$  values. Orange lines represent the best fitting DGBD model for the corresponding experimental data points (shown in blue). Values for the  $\alpha$  and  $\beta$  parameters of the DGBD as well as sample size are shown on top of each panel. PF00787 and PF00439 families have binding molecular function ontology, whereas PF00132 and PF13634 do not have an assigned GO. Notice that x- and y-axis scales are different among panels.

286 In the yeast familiome database, the highest beta value was 1.6529 and  
 287 the lowest beta value was 0.0236. Figure 2 show the cases in which  $\alpha < \beta$ ,  
 288 namely 4 from among the best 30 DGBD graphs with high beta and low alpha  
 289 values. There can be found 5 different general GO molecular functions associ-  
 290 ated with different families. Seven families have “binding” ontology, whereas  
 291 “transferase”, “translocase”, and “ligase” are found in only one family each  
 292 (Figure 2). Although the binding function has a clearly higher prevalence,  
 293 we have 9 different families with different molecular specific binding function.  
 294 We have 3 families with “protein binding” and one family with both “pro-  
 295 tein binding” and “ATP binding”, two families are labelled “DNA binding”  
 296 and one family is labelled “DNA binding” and “RNA polymerase activity”,  
 297 one family has “RNA binding”, one family is labelled “metal ion binding”,  
 298 and two families with other specific ontologies like “proton-transporting ATP  
 299 synthase activity” and “aminoacyl-tRNA editing activity”. In this case, the  
 300 distribution falls rapidly and with high  $\beta$  values, the rank is minor in the  
 301 distribution por lo qué..... [9].

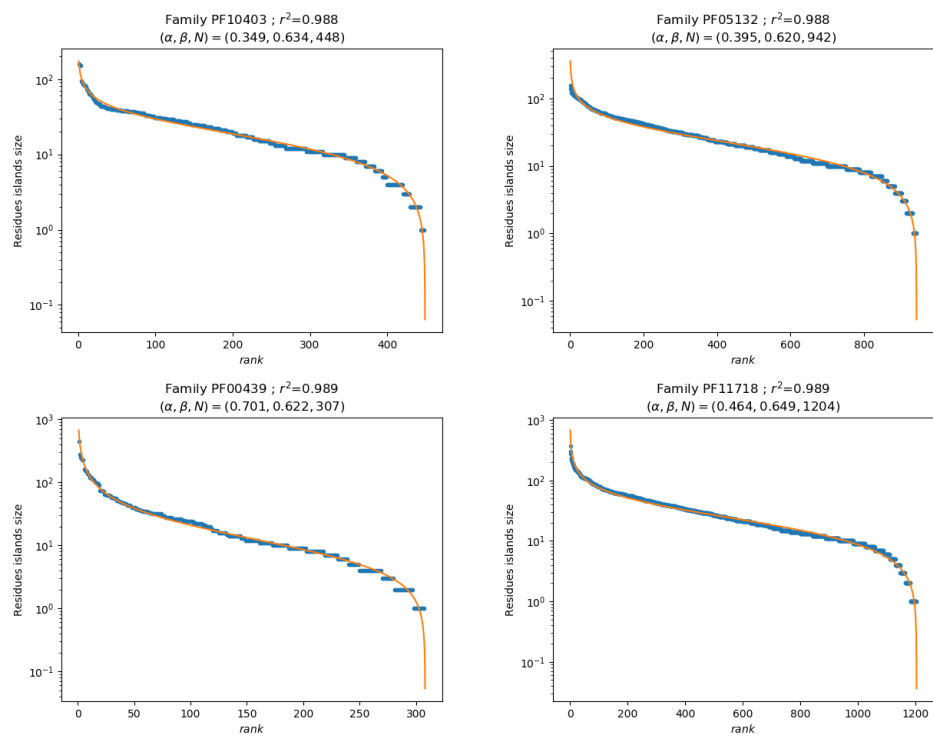


Figure 2: Semilog DGBD plots with high beta and low alpha values. PF05132 family has transferase and binding ontology. PF00439 and PF10403 share binding ontology. Notice that the scales are different.

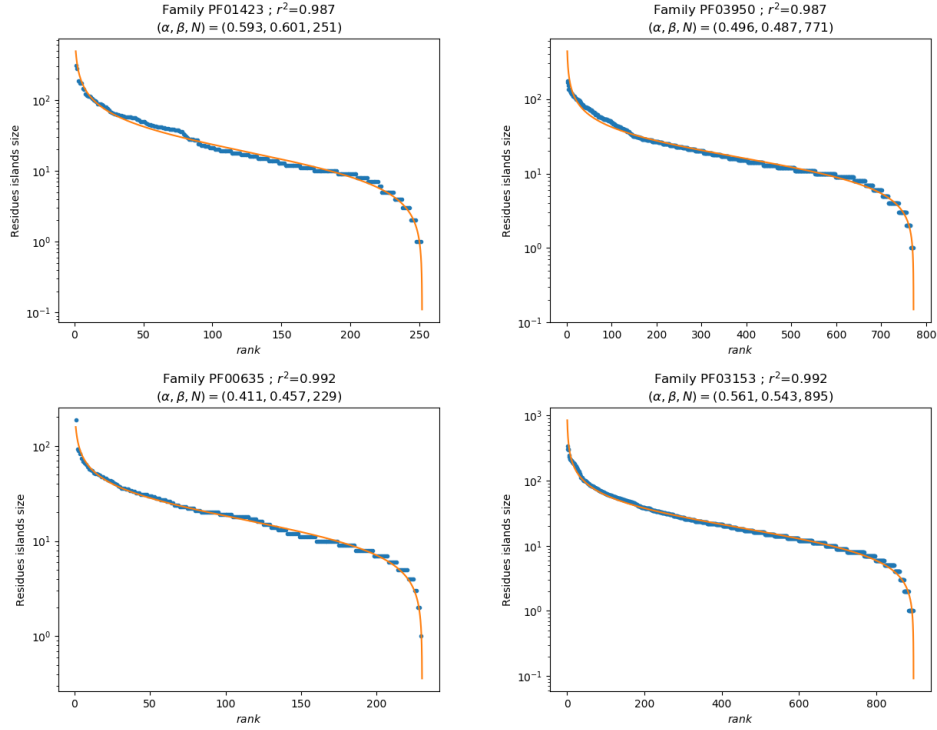


Figure 3: Semilog DGBD plots when alpha equals beta. PF03950 has binding and initiating factor ontologies. The scales are different.

302 In the yeast familiome database, 28 different families were found when  
 303 alpha equals beta with values nearest to  $1 \pm 0.10$ . There are 3 different  
 304 generals GO molecular functions belonging to different families. Six Pfam  
 305 families have a binding ontology, two have ligase ontologies, and hydrolase,  
 306 initiation factor and cooper chaperone ontologies are represented by one each  
 307 (Figure 3). We have 6 families with different binding ontologies and other  
 308 ontologies. One family have zinc ion binding and nucleic acid binding, one  
 309 family have metal ion binding and translation initiation factor activity, one  
 310 family have ATP binding, nucleotide binding and aminoacyl-tRNA ligase  
 311 activity, one family have copper ion binding and copper chaperone activity,  
 312 one family have RNA binding and the last family have ATP-dependent 3'-5'  
 313 DNA helicase activity and aminoacyl-tRNA editing activity.

314 In general, the specific ontologies are different between the parameters  
 315  $\alpha > \beta$ ,  $\alpha < \beta$ , and  $\alpha = \beta$ . There is no consensus in these ontologies, although  
 316 the three cases are represented by binding and ligase general ontologies.

317 LEÓN PATRICIO ENTIENDES LO QUÉ QUIERE DECIR EL PÁRRAFO  
 318 ANTERIOR?

### 319 3.3. *Self-organizing map.*

320 In this work we use a self-organizing map with input data from in a 7-  
 321 dimensional space. This can be seen in figure 5 and 6 where two maps were  
 322 constructed using the before mentioned groups of features respectively. As  
 323 shown in both of these figures, the lack of a more complex structure like the  
 324 one shown in figure 4 gives an idea of how the conjunction of both sets of  
 325 features is required to accomplish such rich structure. Three variables corre-  
 326 spond to gene ontology features such as binding, transferase and hydrolase,  
 327 which are coded in a binary variable whether the attribute is in the family  
 328 or not. The other four features are: parameters a and b from expression (1),  
 329 Shannon's entropy and Kolmogorv complexity.

330 The map was trained with a lattice of 14x14 units. In figure 4 we can  
 331 see the final map where we can detect 6 groups of clustered families located  
 332 in the darker blue areas and these patterns and groups cannot be formed  
 333 neither using gene ontologies nor the complexity features alone, and in figure  
 334 5 we can see in a heatmap each one of different entries of the prototypes,  
 335 this visualization allows us to see the distribution of the different attributes.  
 336 In this map all of the 7 attributes were used. From this final figure we can  
 337 infer that the gene ontology attributes are completely independent from each  
 338 other.

339 In figure 4 there are 6 groups of clustered families with different general  
 340 molecular function annotation between every group and delimited by darker  
 341 blue areas that are a group of neurons that does not have any family repre-  
 342 sentation. The cluster with blue circles have the general molecular function  
 343 annotation binding. It is the biggest cluster with 93 different families. The  
 344 molecular function annotation more specific for these families are: protein  
 345 binding, DNA binding, metal ion binding, ATP binding, RNA binding, nu-  
 346 cleic acid binding, thiamine pyrophosphate binding, GTP binding, nucleotide  
 347 binding, NAD binding, calcium ion binding, heme binding, FMN binding,  
 348 coenzyme binding, flavin adenine dinucleotide binding, phosphatidylinositol  
 349 binding, pyridoxal phosphate binding, chromatin binding, GTPase binding,

350 ubiquitin binding, iron-sulfur cluster binding, translation binding and his-  
351 tone binding. There are 16 different families from this group that had others  
352 molecular function like oxidoreductase, catalytic activity, copper chaperone  
353 activity, translation initiation factor activity and protein dimerization activ-  
354 ity. The cluster with purple "X" have the general molecular function anno-  
355 tation of transferase and have 26 different families. The molecular function  
356 annotation more specific for these families are: phosphotransferase activity,  
357 prenyltransferase activity, methyltransferase activity, and others. The red  
358 pentagon cluster have the general molecular function annotation of hidrolase  
359 and have 22 different families. The molecular function annotation more spe-  
360 cific for these families are: endonuclease activity, thiol-dependent ubiquitinyl  
361 hydrolase activity, deubiquitinase activity, and others.

362 Two clustered families have different molecular functions but join two  
363 different groups of clustered families. The yellow plus sign (+) has 8 different  
364 families and the molecular function of binding and hydrolase and the dark  
365 yellow downward pointing triangles cluster have 7 different families and the  
366 molecular function of binding and transferase. These clusters have only these  
367 two functions delimited by less dark blue areas and link two big families  
368 groups, binding and transferase, and the other families groups are binding  
369 and hydrolase.

370 Finally, the green square cluster has 49 different families with the rest of  
371 the general ontologies that do no has any particular grouping. The ontolo-  
372 gies are oxidoreductase, isomerase, ligase, lyase, translocase, transmembrane  
373 transporter, structural constituent, catalytic activity, antioxidant activity,  
374 structural molecule activity, translation initiation factor activity, motor ac-  
375 tivity, and electron transfer activity.

376 The information in each cluster is different from each one, however, the  
377 information in every neuron is very important to clustering one or more  
378 families and it has to be direct with the general ontologies and complexity  
379 features like DGBD, the Shannon entropy, and the Kolmogorov complexity.  
380 In the future, we hope we can predict the family ontology with this kind of  
381 result, however, we know that we have to obtain more clear results with this  
382 methodology.

383 FINALMENTE TRATAR DE DESCRIBIR LAS ZONAS PROHIBIDAS  
384 COMO "FACTOR DELIMITANTE". EN ESTE CASO LAS ZONAS PRO-  
385 HIBIDAS SE DESCRIBIERON COMO ZONAS AZULES OSCURAS.

386 ES EVIDENTE QUE LA INFORMACIÓN CONTENIDA EN CADA

387 UNA DE LAS ZONAS ES DIFERENTE ENTRE CADA UNA DE ELLAS.  
388 APARENTEMENTE LA ONTOLOGÍA EN LAS ZONAS ESTA DOMI-  
389 NANDO, SIN EMBARGO, ES CLARO QUE CADA UNA DE LAS NEU-  
390 RONAS PUEDE CONTENER UNA O MAS FAMILIAS NO SOLO PRE-  
391 DOMINA LA IMPORTANCIA DE LA ONTOLOGÍA SINO QUE LA GER-  
392 MIBETA, KOLMOGOROV, SHANNON TAMBIÉN JUGAN UN PAPEL  
393 DECISIVO PARA QUE HAYA DIFERENCIA ENTRE CADA UNA DE  
394 LAS NEURONAS.

395 Para Alejandra. Identificar las familias de las 30 mejores alfas mayor  
396 a beta, beta mayor a alfa y alfa igual a beta en el SOM (obviamente las  
397 que tengan ontología) y ver si visulamente están muy separadas o juntas, o  
398 forman parte de la misma neurona o de plano ni se conocen.



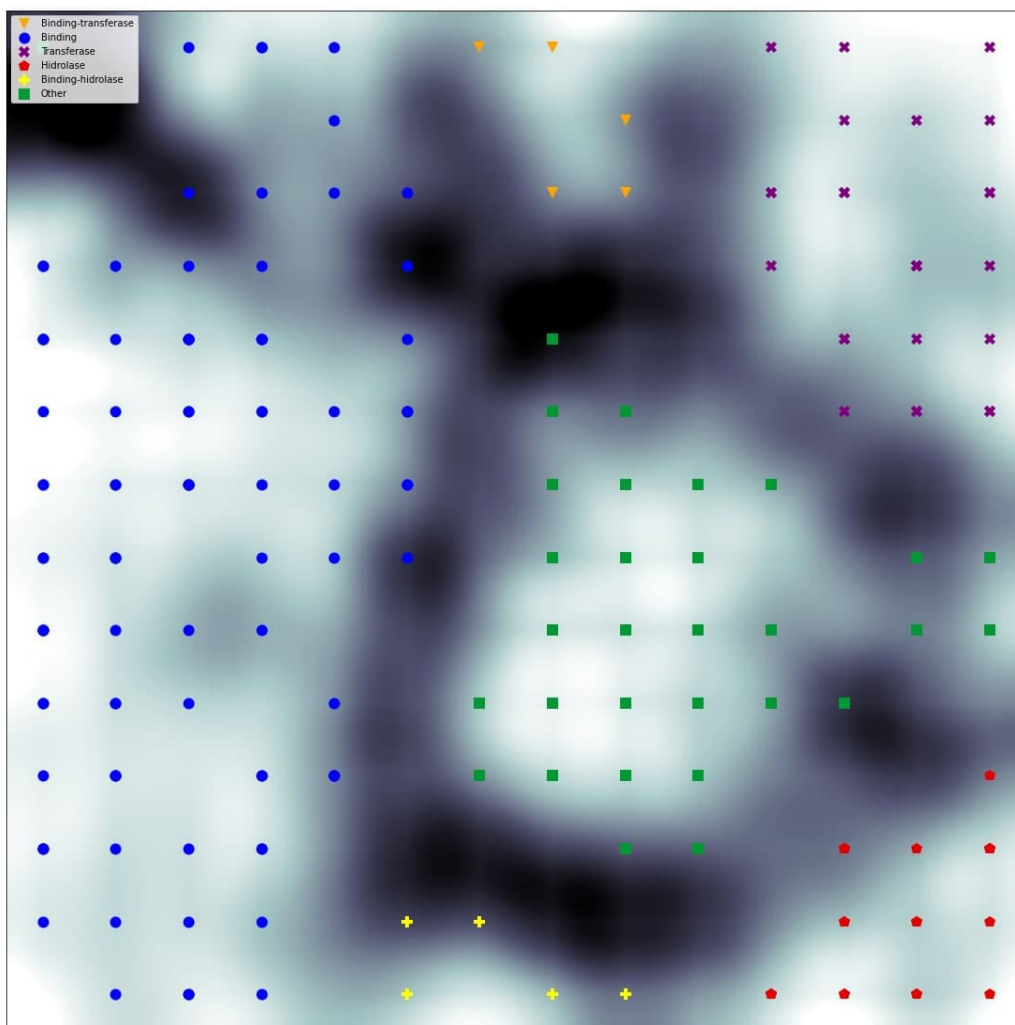


Figure 4: Self organized map of the whole familiome with all of the selected features. This map shows the clustered families projected into a 2D space preserving the topology of the data in the original dimension.

399        Para describir la figura 5.

400        Qué no se parezcan entre si significa que los componentes son independi-  
401 entes para el SOM y que se parezcan entre si, significa que estadísticamente  
402 no son independientes, es decir, que uno dependa del otro. La columna de  
403 datos, a y H, b y K dependen estadísticamente entre si.

404        Poner en algún punto que no se obtuvo la información que se deseaba, o  
405 algo así, en cuánto al uso de RIDs/ROs, o simplemente poner que no se sabe  
406 cuál es el papel de las RIDs/ROs en esta metodología.

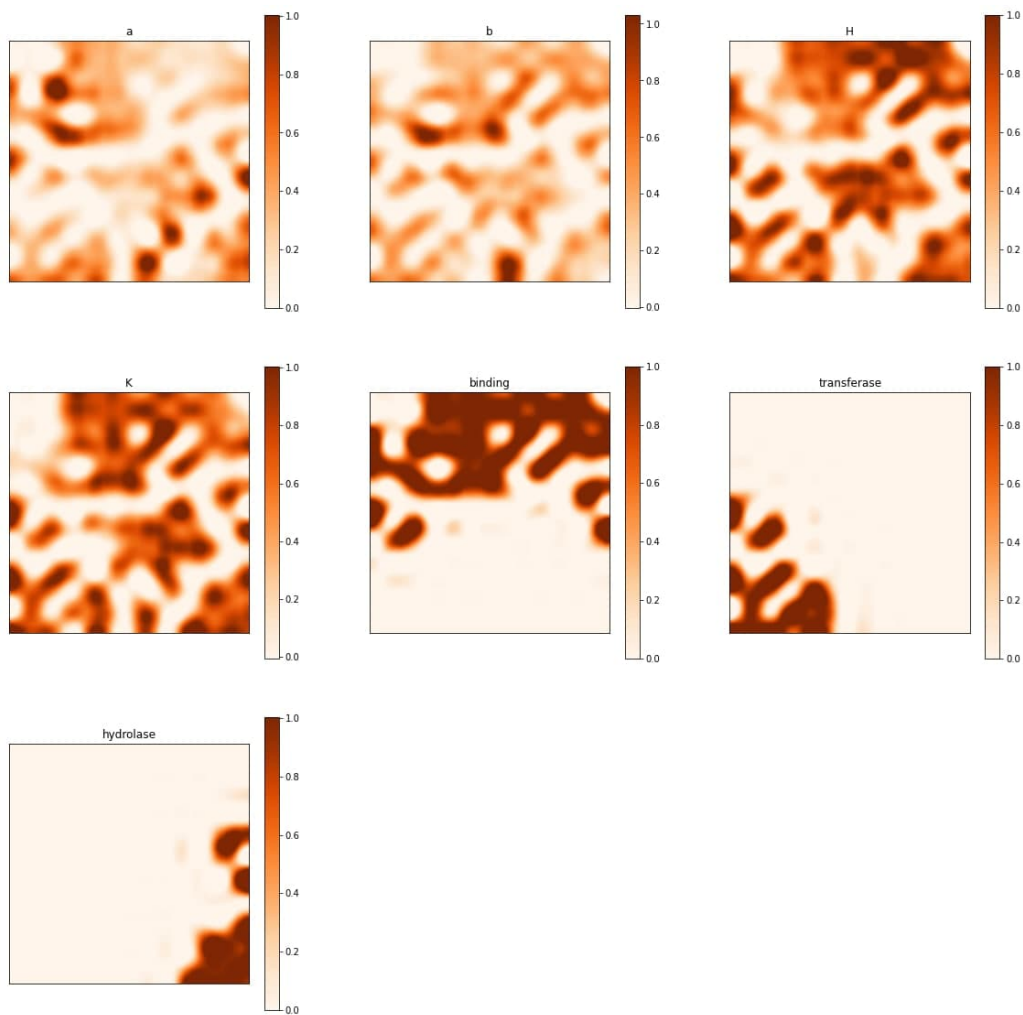


Figure 5: Components of SOM in . This figure shows the distribution of the features used to create the SOM which are a is alfa, b is beta, H is Shannon entropy and K is Kolmogorov complexity, the last three are obtained from the gene ontology functions associated with each family

#### 407 4. Conclusions.

408 NOS INTERESA SABER SI LA GERMIBETA APLICADA A LOS PA-  
409 TRONES DE ORDEN Y DESORDEN DE LAS FAMILIAS DE PROTEÍNAS  
410 NOS SIRVE PARA CARACTERIZAR A LAS FAMILIAS DE PROTEÍNAS.  
411 SE TIENEN GRÁFICAS EN DÓNDE ESTÁN POR UN LADO LAS RIDs  
412 Y LAS ROS EN LA GERMIBETA.

413 ¿qué significa caracterizar? POR EJEMPLO, SI UNA FAMILIA PERTENECE  
414 A UNA FAMILIA O A OTRA, SI EL PATRÓN DE UNA FAMILIA ES  
415 CARACTERÍSTICO O ESPECÍFICO PARA CADA UNA DE LAS FA-  
416 MILIAS O INCLUSO GRUPO DE FAMILIAS, VER SI LA SUMA DE  
417 LAS MEDICIONES QUE SE UTILIZARON AGRUPAN A LAS FAMIL-  
418 IAS. EN DATOS QUE NO SE VAN A PUBLICAR, SE HICIERON SOM's  
419 CON DISTINTOS GRUPOS DE DATOS. EN UNO SE UTILIZÓ SOLO  
420 EL PARÁMETRO DE LA ONTOLOGÍA Y NO SE OBSERVÓ NINGÚN  
421 GRUPO SEPARADO, EN OTRO SOM SOLO SE UTILIZARON LOS DATOS  
422 DE KOLMOGOROV, SHANNON Y ONTOLOGÍA Y NO SE OBSERVARON  
423 LOS GRUPOS QUE SE TIENEN EN LA FIGURA 4, EN OTRO SOM  
424 SE UTILIZARON SOLO LOS DATOS DE LA GERMIBETA Y NO SE  
425 OBSERVARON LA SEPARACIÓN DE GRUPOS. COMO CONCLUSIÓN  
426 PRELIMINAR, SE PUEDE DECIR QUE LOS CUATROS DATOS SON  
427 IMPORTANTES PARA QUE SE AGRUPEN COMO EN LA FIGURA 4.

428 EN ALGÚN PUNTO TRATA DE DISCUTIR LO DE LAS REGIONES  
429 DESORDENADAS, QUE DEBEN IR DE LA MANO CON LAS REGIONES  
430 ORDENADAS, EL TAMAÑO DE LAS REGIONES. DISCUTIR UN POCO  
431 SOBRE LA INTEGRACIÓN DE ORDEN/DESORDEN, LAS PROTEÍNAS  
432 COEXISTEN ENTRE ESAS DOS REGIONES. LAS RIDs PARA NOSOTROS  
433 SON COMO UNA ESPECIE DE HUELLA DIGITAL.

434 COMO PERSPECTIVA. QUE SE PUDIERA UTILIZAR EN LA SE-  
435 CUENCIA COMPLETA Y NO SOLO EN DOMINIOS Y EN FAMILIAS  
436 PARA LA CREACIÓN DE LA MEGASECUENCIA.

437 LA MEDICIÓN DE LA GERMIBETA, KOLMOGOROV Y SHANNON,  
438 SON UNA MÉTRICA? SI LA RESPUESTA ES POSITIVA, ENTONCES  
439 VALDRÍA LA PENA PONERLO EN LAS CONCLUSIONES.

440 EN LA DISCUSIÓN Y LAS CONCLUSIONES "DEBEMOS PONER"  
441 EL SIGNIFICADO DE LAS CURVAS, VISUALMENTE HABLANDO?

## 442 5. References.

- 443 [1] C. J. Oldfield, A. K. Dunker, Intrinsically disordered proteins and in-  
 444 trinsically disordered protein regions, *Annual review of biochemistry* 83  
 445 (2014) 553–584.
- 446 [2] S. Vucetic, Z. Obradovic, V. Vacic, P. Radivojac, K. Peng, L. M. Iak-  
 447 oucheva, M. S. Cortese, J. D. Lawson, C. J. Brown, J. G. Sikes, et al.,  
 448 Disprot: a database of protein disorder, *Bioinformatics* 21 (2004) 137–  
 449 140.
- 450 [3] D. Piovesan, F. Tabaro, I. Mičetić, M. Necci, F. Quaglia, C. J. Oldfield,  
 451 M. C. Aspromonte, N. E. Davey, R. Davidović, Z. Dosztányi, et al.,  
 452 Disprot 7.0: a major update of the database of disordered proteins,  
 453 *Nucleic acids research* 45 (2016) D219–D227.
- 454 [4] A. K. Dunker, P. Romero, Z. Obradovic, E. C. Garner, C. J. Brown,  
 455 Intrinsic protein disorder in complete genomes, *Genome Informatics* 11  
 456 (2000) 161–171.
- 457 [5] A. K. Dunker, C. J. Brown, J. D. Lawson, L. M. Iakoucheva,  
 458 Z. Obradović, Intrinsic disorder and protein function, *Biochemistry*  
 459 41 (2002) 6573–6582.
- 460 [6] A. K. Dunker, M. M. Babu, E. Barbar, M. Blackledge, S. E. Bondos,  
 461 Z. Dosztányi, H. J. Dyson, J. Forman-Kay, M. Fuxreiter, J. Gsponer,  
 462 et al., What’s in a name? why these proteins are intrinsically disordered:  
 463 Why these proteins are intrinsically disordered, *Intrinsically disordered*  
 464 *proteins* 1 (2013) e24157.
- 465 [7] A. K. Dunker, J. D. Lawson, C. J. Brown, R. M. Williams, P. Romero,  
 466 J. S. Oh, C. J. Oldfield, A. M. Campen, C. M. Ratliff, K. W. Hipps,  
 467 et al., Intrinsically disordered protein, *Journal of molecular graphics*  
 468 *and modelling* 19 (2001) 26–59.
- 469 [8] R. Van Der Lee, M. Buljan, B. Lang, R. J. Weatheritt, G. W. Daughdrill,  
 470 A. K. Dunker, M. Fuxreiter, J. Gough, J. Gsponer, D. T. Jones, et al.,  
 471 Classification of intrinsically disordered regions and proteins, *Chemical*  
 472 *reviews* 114 (2014) 6589–6631.

- 473 [9] G. Martínez-Mekler, R. A. Martínez, M. B. del Río, R. Mansilla, P. Mi-  
474 ramontes, G. Cocho, Universality of rank-ordering distributions in the  
475 arts and sciences, *PLoS One* 4 (2009) e4791.
- 476 [10] D. L. Nelson, A. L. Lehninger, M. M. Cox, *Lehninger principles of bio-*  
477 *chemistry*, 7th ed., Macmillan, 2017.
- 478 [11] B. Alberts, A. Johnson, J. Lewis, D. Morgan, M. Raff, K. Roberts,  
479 P. Walter, *Molecular Biology of the Cell in Cell*, 6th ed., Garland Sci-  
480 ence, Taylor & Francis Group, 2015.
- 481 [12] A. Bateman, E. Birney, R. Durbin, S. R. Eddy, R. D. Finn, E. L.  
482 Sonnhammer, Pfam 3.1: 1313 multiple alignments and profile hmms  
483 match the majority of proteins, *Nucleic acids research* 27 (1999) 260–  
484 262.
- 485 [13] B. Xue, R. L. Dunbrack, R. W. Williams, A. K. Dunker, V. N. Uver-  
486 sky, Ponder-fit: a meta-predictor of intrinsically disordered amino acids,  
487 *Biochimica et Biophysica Acta (BBA)-Proteins and Proteomics* 1804  
488 (2010) 996–1010.
- 489 [14] P. Tompa, M. Fuxreiter, C. J. Oldfield, I. Simon, A. K. Dunker, V. N.  
490 Uversky, Close encounters of the third kind: disordered domains and  
491 the interactions of proteins, *Bioessays* 31 (2009) 328–335.
- 492 [15] E. L. Sonnhammer, S. R. Eddy, R. Durbin, Pfam: a comprehensive  
493 database of protein domain families based on seed alignments, *Proteins:*  
494 *Structure, Function, and Bioinformatics* 28 (1997) 405–420.
- 495 [16] A. Bateman, L. Coin, R. Durbin, R. D. Finn, V. Hollich, S. Griffiths-  
496 Jones, A. Khanna, M. Marshall, S. Moxon, E. L. Sonnhammer, et al.,  
497 The pfam protein families database, *Nucleic acids research* 32 (2004)  
498 D138–D141.
- 499 [17] M. Punta, P. C. Coghill, R. Y. Eberhardt, J. Mistry, J. Tate,  
500 C. Boursnell, N. Pang, K. Forslund, G. Ceric, J. Clements, et al., The  
501 pfam protein families database, *Nucleic acids research* 40 (2011) D290–  
502 D301.
- 503 [18] R. D. Finn, P. Coghill, R. Y. Eberhardt, S. R. Eddy, J. Mistry, A. L.  
504 Mitchell, S. C. Potter, M. Punta, M. Qureshi, A. Sangrador-Vegas, et al.,

- 505 The pfam protein families database: towards a more sustainable future,  
506 Nucleic acids research 44 (2015) D279–D285.
- 507 [19] J. M. Cherry, E. L. Hong, C. Amundsen, R. Balakrishnan, G. Binkley,  
508 E. T. Chan, K. R. Christie, M. C. Costanzo, S. S. Dwight, S. R. En-  
509 gel, et al., Saccharomyces genome database: the genomics resource of  
510 budding yeast, Nucleic acids research 40 (2012) D700–D705.
- 511 [20] G. O. Consortium, Gene ontology consortium: going forward, Nucleic  
512 acids research 43 (2015) D1049–D1056.
- 513 [21] R. Linding, L. J. Jensen, F. Diella, P. Bork, T. J. Gibson, R. B. Rus-  
514 sell, Protein disorder prediction: implications for structural proteomics,  
515 Structure 11 (2003) 1453–1459.
- 516 [22] K. Katoh, M. C. Frith, Adding unaligned sequences into an existing  
517 alignment using mafft and last, Bioinformatics 28 (2012) 3144.
- 518 [23] M. Ashburner, C. A. Ball, J. A. Blake, D. Botstein, H. Butler, J. M.  
519 Cherry, A. P. Davis, K. Dolinski, S. S. Dwight, J. T. Eppig, et al., Gene  
520 ontology: tool for the unification of biology, Nature genetics 25 (2000)  
521 25–29.
- 522 [24] A. M. Petersen, H. E. Stanley, S. Succi, Statistical regularities in the  
523 rank-citation profile of scientists, Scientific reports 1 (2011) 181.
- 524 [25] B. Alberts, A. Johnson, J. Lewis, M. Raff, K. Roberts, P. B. Andersen,  
525 N. O. Finnemann, P. V. Christiansen, M. Ashburner, C. A. Ball, et al.,  
526 Adami, christoph 2004. information theory in molecular biology. physics  
527 of life reviews 1: 3–22. adams, frederick 2003. the informational turn  
528 in philosophy. minds and machi-nes 13: 471–501., Annual Review of  
529 Biochemistry 72 (1925) 291–336.
- 530 [26] M. Li, P. M. Vitanyi, Kolmogorov complexity and its applications, Cen-  
531 tre for Mathematics and Computer Science, 1989.
- 532 [27] T. Kohonen, Essentials of the self-organizing map, Neural networks 37  
533 (2013) 52–65.
- 534 [28] J. Wang, J. Delabie, H. C. Aasheim, E. Smeland, O. Myklebost, Clus-  
535 tering of the som easily reveals distinct gene expression patterns: results  
536 of a reanalysis of lymphoma study, BMC bioinformatics 3 (2002) 1–9.

- 537 [29] V. Nurminen, A. Neme, S. Seuter, C. Carlberg, The impact of the vita-  
538 min d-modulated epigenome on vdr target gene regulation, *Biochimica*  
539 *et Biophysica Acta (BBA)-Gene Regulatory Mechanisms* 1861 (2018)  
540 697–705.
- 541 [30] O. Fontanelli, P. Miramontes, Y. Yang, G. Cocho,  
542 W. Li, Beyond zipf’s law: The lavalette rank func-  
543 tion and its properties, *PLOS ONE* 11 (2016) 1–14.  
544 URL: <https://doi.org/10.1371/journal.pone.0163241>.  
545 doi:10.1371/journal.pone.0163241.
- 546 [31] K. Gupta, M. Lalit, A. Biswas, C. D. Sanada, C. Greene, K. Hukari,  
547 U. Maulik, S. Bandyopadhyay, N. Ramalingam, G. Ahuja, et al., Mod-  
548 eling expression ranks for noise-tolerant differential expression analysis  
549 of scrna-seq data, *Genome Research* 31 (2021) 689–697.
- 550 [32] W. Li, O. Fontanelli, P. Miramontes, Size distribution of function-based  
551 human gene sets and the split–merge model, *Royal Society open science*  
552 3 (2016) 160275.