

# Evolutionary Models of Amino Acid Substitutions Based on the Tertiary Structure of their Neighborhoods.

Elias Primetis<sup>1,2</sup>, Spyridon Chavlis<sup>2</sup>, and Pavlos Pavlidis<sup>3</sup>

<sup>1</sup>University of Crete, Department of Biology

<sup>2</sup>Foundation for Research and Technology, Hellas, Institute of Molecular Biology and Biotechnology, 100 Nikolaou Plastira str., Vassilika Vouton, Heraklion, Crete GR 700 13, Greece

<sup>3</sup>Foundation for Research and Technology, Hellas, Institute of Computer Science, 100 Nikolaou Plastira str., Vassilika Vouton, Heraklion, Crete GR 700 13, Greece

## Abstract

Intra-protein residual vicinities depend on the involved amino acids. Energetically favorable vicinities (or interactions) have been preserved, while unfavorable vicinities have been eliminated during evolution. We describe, statistically, the interactions between amino acids using resolved protein structures. Based on the frequency of amino acid interactions, we have devised an amino acid substitution model that implements the following idea: amino acids that have similar neighbors in the protein tertiary structure can replace each other, while substitution is more difficult between amino acids that prefer different spatial neighbors. Using known tertiary structures for  $\alpha$ -helical membrane (HM) proteins, we build evolutionary substitution matrices. We constructed maximum likelihood phylogenies using our amino acid substitution matrices and compared them to widely-used methods. Our results suggest that amino acid substitutions are associated with the spatial neighborhoods of amino acid residuals.

## 1 Introduction

2 Structure and functionality of a protein are closely related (Worth et al., 2009). Its amino acid sequence  
3 as well as cofactors, ligands and other parts of the same or other proteins form a complex network  
4 of interactions that is the basis of the unique physicochemical properties of each protein related to its  
5 function (Worth et al., 2009). During the last decade, a multitude of tertiary protein structures have  
6 been determined by using techniques such as crystallography, NMR, electron microscopy and hybrid  
7 methods (Egli, 2010). Consequently, the number of the available tertiary structures in the RCSB-PDB  
8 database (Bernstein et al., 1977) amino acid is rapidly increasing.

9 Computational methods facilitate structural, functional and evolutionary characterization of the pro-  
10 teins (Nath Jha et al., 2011). Protein function is tightly linked to its tertiary structure, and consequently  
11 to the vicinities, or interactions, amino acid residues have formed. For example, globular and membrane

12 proteins are characterized by totally different physicochemical environments. Membrane proteins show  
13 a limited interaction with water molecules, and, in contrast, they are able to interact with the lipid  
14 bilayer. Thus, their trans-membrane region adopts a single type of secondary structure – either a helix  
15 or a beta-sheet. Largely, the secondary structure is defined by non-neighboring (on the sequence level)  
16 amino acid residue interactions (Nath Jha et al., 2011).

17 Recently, based on the ideas of (Nath Jha et al., 2011), we developed PrInS (freely available from <http://pop-gen.eu/wordpress/software/prins-protein-residues-interaction-statistics> (Protein In-  
18 teraction Statistics; Pavlidis et al., unpublished) an open-source software to score proteins based on the  
19 frequency of their residue interactions. PrInS uses protein structures stored in Protein Data Bank (PDB)  
20 to construct a statistical model of intra-protein amino acid residue interactions for a certain class of  
21 proteins (e.g., membrane proteins). PrInS scores every amino acid  $a$  proportionally to the number of  
22 ‘unexpected residues’ that interact with it. The term ‘unexpected residues’ means residues that they are  
23 rarely found to be in the vicinity of  $a$  if we consider all protein structures of a given dataset. Therefore, it  
24 is able to pinpoint residues characterized by a large number of ‘less frequent’ interactions, and therefore  
25 they may represent functional areas of the proteins or targets of natural selection based on the following  
26 assumption: even though a large number of unlikely interactions characterizes these amino acids, nature  
27 has preserved them. Here, we use PrInS to describe statistically the intra-protein amino acid interactions  
28 and consequently to construct an amino acid substitution matrix endowed with the principle that amino  
29 acids with similar (residual) neighborhoods can substitute each other during evolution.  
30

31 Protein evolution comprises two major principles. The first principle suggests that protein structure  
32 is more conserved than the sequence (Siltberg-Liberles et al., 2011). The second principle suggests that  
33 the physicochemical properties of amino acids constrain the structure, the function and the evolution of  
34 proteins (Hatton and Warr, 2015). Protein evolutionary rate is strongly correlated with fractional residue  
35 burial(Siltberg-Liberles et al., 2011). This is due to the fact that the core of a protein is mostly formed  
36 by buried residues, which often play a crucial role in the stability of the folded structure (Franzosa and  
37 Xia, 2008). The three-dimensional structure of the protein determines its evolutionary rate, since most  
38 mutations in the core of a protein tend to destabilize the protein.

39 Within protein families the backbone changes are infrequent, thus, preserving the folding properties  
40 over relatively long evolutionary distances, while substitutions are found often at the side chains. For  
41 example, for proteins with binding function, the binding interface is under functional constraint and  
42 may evolve the slowest, with differences in rate between affinity-determining and specificity-determining  
43 residues (Siltberg-Liberles et al., 2011). In addition, the secondary structural elements of a protein  
44 evolve at different rates. Beta sheets evolve more slowly than helical regions and random coils evolve the  
45 fastest (Siltberg-Liberles et al., 2011). Secondary structure changes may eventually occur due to varying  
46 helix/sheet propensity. Some of these changes in secondary structural composition may be evolutionary

47 neutral, whereas some structural transition may involve negative or positive selection. In the latter case,  
48 a new mutationally accessible fold may enable the development of a new favorable function that was not  
49 possible within the previous fold (Siltberg-Liberles et al., 2011).

50 In the context of folding, the thermodynamic stability of the proteins with a stable unique tertiary  
51 structure is important. Thermodynamic stability is maintained throughout evolution despite the desta-  
52 bilizing effect of the non-synonymous mutations, which are often removed from the populations as a  
53 result of negative selection. The protein structure is important because it acts as a scaffold for properly  
54 orientating functional residues, such as a binding interface and a catalytic residue. As a result, the selec-  
55 tive pressure for particular sequences (and not structures) over longer evolutionary periods is decreased,  
56 generating a neutral network of sequences interconnected via mutational changes (Siltberg-Liberles et al.,  
57 2011).

58 Choi and Kim (Choi and Kim, 2006), based on the most common structural ancestor (CSA), showed  
59 that not all present-day proteins evolved from one single set of proteins in the last common ancestral  
60 organism, but new common ancestral protein were born at different evolutionary times. These proteins  
61 are not traceable to one or two ancestral proteins, but they follow the rules of the "multiple birth model"  
62 for the evolution of protein sequence families (Choi and Kim, 2006).

63 In this study we used the scoring matrices that are obtained from PrInS algorithm to examine the evo-  
64 lution of proteins and we focused on the evolution of  $\alpha$ -helical membrane proteins. The main hypothesis  
65 is that protein evolution is related to the three dimensional neighborhoods of amino acids. Specifically,  
66 amino acids that have the same amino acid neighbourhoods can substitute each other during evolution.

## 67 **Methods**

### 68 **Dataset retrieval and name conversion**

69 In eukaryotes,  $\alpha$ -helical proteins exist mostly in the plasma membrane or sometimes in the outer cell  
70 membrane. In prokaryotes, they are present in their inner membranes. We used 82  $\alpha$ -helical membrane  
71 proteins, which were also scrutinized previously by Nath Jha et al. (2011) to describe the statistical  
72 properties of amino acid interactions within  $\alpha$ -helical membrane proteins. First, for each of the protein  
73 in the dataset, we downloaded multiple sequence alignments of homologous protein sequences from the  
74 UCSC Genome Browser (Kent et al., 2002) for 16 primate and 3 non-primate mammalian species. Second,  
75 three-dimensional structures were downloaded from the Protein Data Bank (PDB) database.

### 76 **The PrInS software**

77 We applied PrInS on the three dimensional structures downloaded from PDB. PrInS constructs a scoring  
78 matrix  $M$  as follows: If the tertiary structure of a protein is represented by  $P$  and the total number of

79 amino acid residues is  $l$ , then residues  $1 \leq k, m \leq l$  interact if and only if:

$$A_{km} = \begin{cases} 1, & \text{if } d(C_\alpha^k - C_\alpha^m) \leq 6.5\text{\AA} \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

80 where  $d(C_\alpha^k - C_\alpha^m)$  denotes the distance between the  $C_\alpha$  atoms of the  $k$  and  $m$  amino acids. In other  
81 words, the  $k^{\text{th}}$  and the  $m^{\text{th}}$  residues of a protein interact if and only if the Euclidean distance between  
82 their  $C_\alpha$  atoms is less than  $6.5\text{\AA}$  and they are not located on adjacent positions on the amino acid chain.

### 83 Defining the environment of each amino acid

84 According to Nath Jha et al. (2011), amino acids of a HM protein can be classified in three environ-  
85 ments: The first environment comprises amino acids at different distances from the lipid bilayer. The  
86 second environment classifies the helices (and thus, their amino acids) on the basis of their inter-helical  
87 interactions inside the membrane. Finally, the third environment classifies amino acid residues based  
88 on the number of interactions they make. In our study we have used the third type of environments  
89 (residue-contact-based interaction) based on the results from Nath Jha et al. (2011), who demonstrated  
90 that the residue-contact-based environment description of residue interaction is more accurate for the  
91  $\alpha$ -helical proteins they studied. Thus, every amino acid residue  $P_i$  in protein  $P$  can be assigned to the  
92 pair  $(A, K)$ , where  $A$  indexes the amino acid type of the  $k^{\text{th}}$  residue (e.g. Alanine) and  $K$  represent the  
93 environment of the  $k^{\text{th}}$  residue. Based on the results of Nath Jha et al. (2011), we used the number of  
94 non-covalent contacts each amino acid makes to define its environment. Thus, environment I comprises  
95 all amino acids with 1-5 contacts, environment II, amino acids with exactly 6 contacts and environment  
96 III amino acids with more than 6 contacts.

97 Initially, PrInS was used to construct a matrix  $M$ , a  $60 \times 60$  matrix (or equivalently nine  $20 \times 20$   
98 scoring matrices), that scores amino acid interactions in all environment pairs. For example,  $M_{i,j}$ ,  $i, j \leq 20$   
99 describe score interactions between amino acids of the environment I. In  $M$  the pairs of amino acids with  
100 the lower scores are those that interact frequently. In general, the interaction between amino acid A  
101 ( $1 \leq A \leq 20$ ) and amino acid B ( $1 \leq B \leq 20$ ) that belong to the environments  $K$  ( $1 \leq K \leq 3$ ) and  $Q$   
102 ( $1 \leq Q \leq 3$ ), respectively, is given by:

$$M_{ij} = -\ln\left(\frac{n_{A,K-B,Q}}{g \times (S_A/S) \times (S_B/S) \times E_{K,Q}}\right) \quad (2)$$

103 .

104 The coordinates  $i$  and  $j$  are given by  $i = 20(K - 1) + A$  and  $j = 20(Q - 1) + B$ , respectively.  
105 The parameter  $g = 2$  if the contacting pair comprises different amino acids in the same environment  
106 ( $A \neq B, K = Q$ ), and  $g = 1$ , otherwise.  $S_A$  and  $S_B$  are the total number of amino acids  $A$  and  $B$  in

107 the dataset, respectively.  $S$  is the total number of amino acids in the dataset and  $E_{K,Q}$  equals to the  
108 number of interactions between environment  $K$  and  $Q$ . Finally,  $n_{A,K-B,Q}$  is the number of interactions  
109 between amino acids  $A$  in environment  $K$  and amino acid  $B$  in environment  $Q$ . It is evident that  $M_{i,j}$   
110 assumes large (positive) values when the observed number of interactions  $n_{A,K-B,Q}$ , between  $A$  and  $B$ ,  
111 is much lower than the interactions expected based on the total frequency of the amino acids  $A$  and  $B$   
112 in the dataset and the number of interactions between the  $K$  and  $Q$  environments. In other words the  
113 value of  $M_{i,j}$  is large for rare interactions.

#### 114 **Distance matrix**

115 The  $M$  scoring matrix was split in nine  $20 \times 20$  matrices and distance metrics (Euclidean, Manhattan and  
116 Pearson Squared) between the elements of each matrix (amino acids) were calculated. Let  $G_d$  represent  
117 the average  $20 \times 20$  distance matrix for the distance metrics  $d$  ( $d \in \{\text{Euclidean, Manhattan, Pearson}$   
118  $\text{Squared}\}$ ). We name these matrices **neighborhood Euclidean-distance based** ( $nEd$ ),  $nMd$  and  $nPd$ ,  
119 respectively. For comparison purposes, we created an additional  $20 \times 20$  distance matrix based on the  
120 Hamming distance between a pair of codons showing the minimum number of nucleotide changes required  
121 to transform one amino acid to another. Smaller distance values indicate higher similarity between the  
122 respective amino acids. Distance matrices have been visualized as heatmaps (**Supplementary Files**  
123 **S1a, S1c and S1d**). Two amino acids that have similar neighborhoods, *i.e.* similar amino acid neighbors  
124 in the tertiary structure will, thus, have a small distance (high similarity) between them. Consequently,  
125 according to our hypothesis, they will substitute each other during evolution at a higher rate.

#### 126 **Substitution rate matrix**

127 Each of the four distance matrices were transformed into rate matrices  $D_r$  by using a similar procedure  
128 as in Dayhoff et al. (1978). In particular, the following procedure was followed:

- 129 1. Find the maximum value of the distance matrix  $D_g$ ,  $D_{max}$ .
- 130 2. Subtract each matrix element from the  $D_{max}$  and divide by  $D_{max}$ . This will transform  $D_g$  to  
131 an amino acid "similarity" symmetric matrix, where 1 denotes maximum 'similarity' between two  
132 amino acids.
- 133 3. Similar to the substitution matrices BLOSUM62 and PAM120, one element of the matrix is defined  
134 as a reference. All other elements are scaled proportionally to the reference.
- 135 4. The diagonal of the matrix is defined as the negative sum of all other elements of the respective  
136 row.

137 If an amino acid pair is characterized by a large (relative) substitution rate, then they can substitute each

138 other during evolution at a high rate. The substitution rate matrix that is derived with the aforementioned  
139 algorithm is symmetric.

#### 140 **Evaluation of PrInS ability to predict amino acid substitutions**

141 For every multiple alignment, a substitution matrix  $S$  was created from the alignment, by counting  
142 the occurrences of the two most common amino acids for every alignment site, and assuming that a  
143 substitution between these two amino acids has occurred during evolution. To evaluate the relation of our  
144 neighborhood-based substitution rate matrices  $D_r$  and the multiple sequence alignments, we calculated  
145 the Pearson correlation coefficients between  $S$  and  $D_r$  (for each distance metric). High Pearson correlation  
146 coefficient values suggest that the amino acid substitutions in the alignments are concordant with the  
147 neighborhood-based substitution rates distances.

148 Moreover, we manually analyzed the amino acid pairs with the lowest distance values in the Euclidean,  
149 Manhattan and Pearson distance matrices with the tool ‘Common Substitution Tool’ in the ‘Amino Acid  
150 Explorer’ ([https://www.ncbi.nlm.nih.gov/Class/Structure/aa/aa\\_explorer.cgi](https://www.ncbi.nlm.nih.gov/Class/Structure/aa/aa_explorer.cgi)) of NCBI (Bulka  
151 et al., 2006). Common Substitution Tool sorts an amino acid list from the most common to the least com-  
152 mon substitutions, given an amino acid as input, based on the BLOSUM62 substitution matrix (Henikoff  
153 and Henikoff, 1992).

#### 154 **Overall and Site Likelihoods**

155 Likelihood is a function of the parameters of a statistical model for given data. Since substitution rate  
156 matrices are part of the model, we compared the per site and the overall likelihood values obtained by our  
157 substitution matrices and those obtained by using BLOSUM62 and PAM120. Both the phylogeny and  
158 the equilibrium frequencies of the amino acids were considered known in this analysis. Thus, by keeping  
159 the remaining parameters the same for both models, we aimed to assess which rate matrix can better  
160 explain the observed multiple alignments. As overall likelihood, we defined the total likelihood of the  
161 alignment, while a site likelihood is the likelihood of a single position in the amino acid alignment. The  
162 phylogenetic tree was retrieved from the UCSC Genome Browser (Kent et al., 2002), and the equilibrium  
163 frequencies of amino acids were obtained from the equilibrium frequencies in the BLOSUM62 model. We  
164 used the same phylogenetic tree and equilibrium frequencies for all comparisons.

## 165 **Results**

### 166 **Evaluation of PrInS Ability to Predict Amino Acid Substitutions**

167 The substitution matrix counts the amino acid substitutions that occurred in the 224 downloaded multiple  
168 alignments using the two most frequent amino acid residues at each alignment site. For each amino acid in

169 this matrix, we evaluated its correlation (Pearson correlation coefficient) with the same amino acid in the  
170 distance matrices. Thus, if, for a given amino acid, both matrices suggest similar substitution preferences,  
171 then neighborhood-based substitution models are concordant with the observed substitutions. Since we  
172 used distance matrix in this analysis, negative correlation coefficient indicate concordance. Results suggest  
173 that there is concordance between the substitution matrix and our neighborhood-based distance matrices  
174 for most of the amino acids. Figure 1 shows the Pearson correlation coefficients for amino acids using  
175 the substitution matrix and the *nEd*. The correlation plots between the substitution matrix and other  
176 distances are illustrated in the Supplementary material (Figures S2a-S2c).

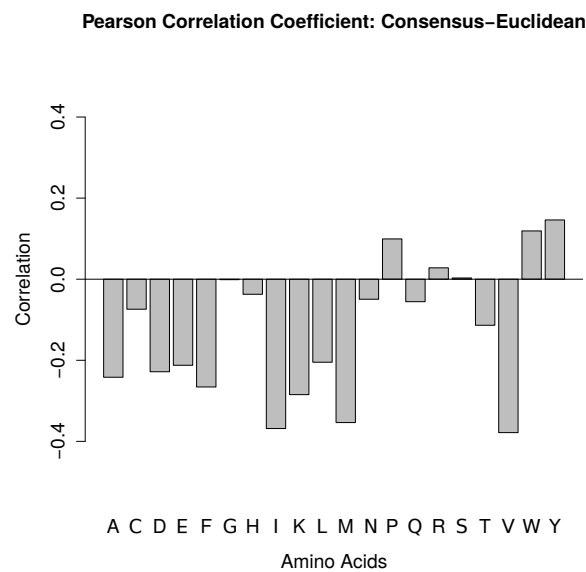


Figure 1: Pearson Correlation Coefficient between the substitution and *nEd*. Concordance is indicated with negative values, while discordance with positive values.

177 As it is shown in Figure 1, most of the amino acids are concordant when we compare the Euclidean  
178 distance matrix and the substitution matrix. For example, both the uncorrelated and discordant amino  
179 acids were glycine, proline, arginine, serine, tryptophan and tyrosine.

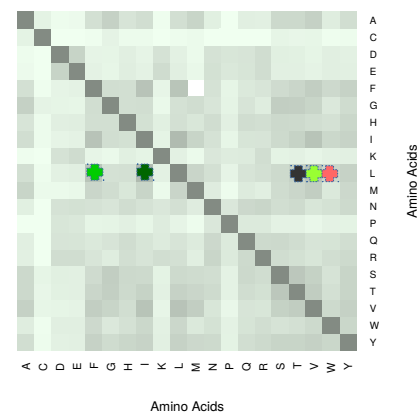
180 In the Supplementary Figure S2b (substitution vs Manhattan distance), the discordant amino acids  
181 were the five out of six amino acids of Figure 1. Similar results are shown in Supplementary Figure S2c,  
182 where the substitution matrix is compared against the squared Pearson distance matrix. Finally, there  
183 are no discordant amino acids in the correlation between the Consensus matrix and the Genetic Code  
184 matrix (Supplementary Figure S2a).

185 Furthermore, the amino acid pairs with the lowest values in the distance matrices are at the top of  
186 the amino acid list that the Common Substitution Tool of NCBI Amino Acid Explorer reports. This  
187 means that amino acids with similar three dimensional neighborhoods tend to substitute each other. For  
188 example, in Euclidean distance matrix, the Leukine-Isoleukine pair has the lowest distance value and  
189 according to the substitution list of Amino Acid Explorer, Isoleukine is the most common substitution

190 of Leucine. Figure 2 shows the first five most common amino acid substitutions for Leucine according  
 191 to ‘Common Substitution Tool’ (Figure 2a) and according to the Euclidean distance matrix (Figure 2b).  
 192 Evidently, the Euclidean distance matrix can predict the most common substitutions of Leucine, but in  
 193 slightly different order. For example, Common Substitution Tool reports the following order of amino  
 194 acids: Isoleucine, Methionine, Valine, Phenylalanine and Alanine. According to the Euclidean distance  
 195 matrix the order is: Isoleucine, Phenylalanine, Valine, Tryptophan and Tyrosine. Methionine is the  
 196 sixth most common substitution of Leucine, while Alanine is the tenth most common substitution of  
 197 Leucine according to the Euclidean distance matrix. On the other hand, Tryptophan and Tyrosine are  
 198 the seventh and the eighth most common substitutions of Leucine according to the Common Substitution  
 199 Tool. Neglecting the order of the first five substitutions, drawing five amino acids and observing three  
 200 common is marginally significant ( $p$ -value = 0.07, right tail of hypergeometric distribution). Differences  
 201 between two reports are possibly due to the fact that Euclidean distance matrix is solely based on the  
 202 structural information of a protein family, while the Common Substitution Tool is based on BLOSUM62  
 203 substitution matrix that was created using alignments of multiple protein families.

1-letter code	3-letter code	Chemistry	Potential H-bonds	Molecular Weight	Isoelectric Point	Hydrophobicity
L	Leu	CH <sub>2</sub> -C-C-	0	113	6.0	0.918
I	Ile	CH <sub>2</sub> -C-C-	0	113	6.0	1.000
M	Met	C-S-C-C-	0	131	5.7	0.811
V	Val	CH <sub>2</sub> -C-C-	0	99	6.0	0.923
E	Phe	CH <sub>2</sub> -C-C-	0	147	5.5	0.951
A	Ala	CH <sub>2</sub> -C-C-	0	71	6.0	0.805

(a)



(b)

Figure 2: The first five most common substitutions of Leucine according to (a) Common Substitution Tool and (b) Euclidean Distance Matrix. The same color code is used in both figures. Dark green shows the first most common substitution, while the red shows the fifth most common substitution. Methionine and alanine are not found among the first five most common substitutions of Leucine in Euclidean distance matrix, while Tryptophan and Tyrosine are not found among the first most common substitutions in Common Substitution Tool.

## 204 Pairwise comparison of proteins in different species

205 Using the Euclidean-based distance matrix, we scored the differences of each protein sequence to its human  
 206 homologue. Results are presented as a heatmap and are clustered hierarchically based on their distance



207 from the human homologue (Figure 3). Darker gray tone in Figure 3 denotes high similarity between  
208 human and other species homologues, whereas lighter tones suggest lower similarity. As expected, proteins  
209 from species that are evolutionarily distant from humans are more dissimilar to human homologues. This  
210 is especially true for a group of proteins clustered together using the Euclidean distance-based substitution  
211 matrix (Figure 3 proteins labeled with three dashes).

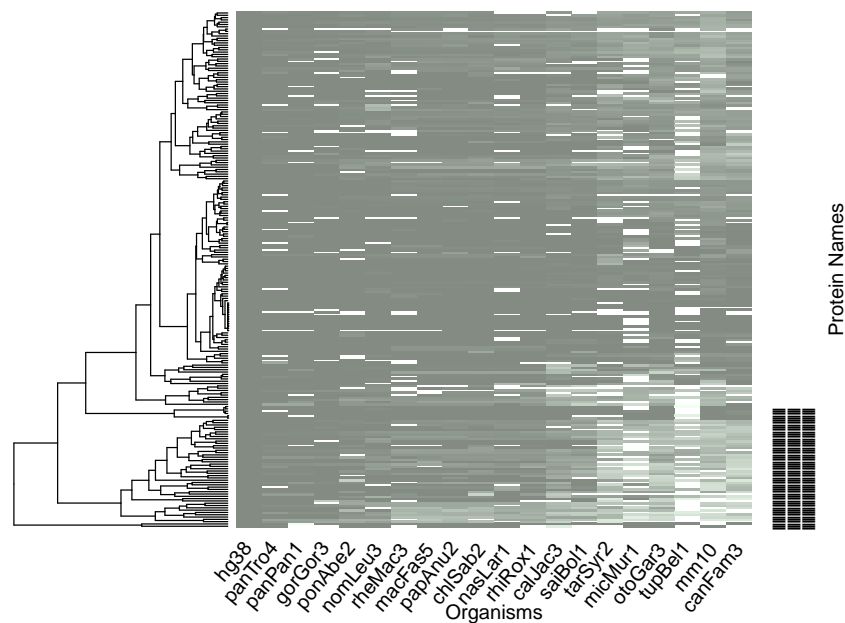


Figure 3: An illustration of the scoring of multiple alignments with the Euclidean distance matrix as a heatmap. Each protein (rows) from human (first column) is compared against the homologous proteins from another species (columns) using pairwise amino acid sequence alignments. We used only the alignment sites with amino acids  $A$  and  $B$ , where  $A \neq B$  and  $A \neq -$  and  $B \neq -$ . The total score (a cell in the heatmap) is calculated as the sum of all site scores. Proteins form two clusters based on the hierarchical tree on the left side of the figure. Proteins of the bottom cluster (also denoted with three dashes) are further analyzed using Gene Ontology terms.

212 The group of proteins clustered together by the Euclidean-based matrix (marked in the lower part  
213 of Figure 3) was further scrutinized by Gene Ontology (GO) terms. We used gProfiler (Reimand et al.,  
214 2007), to obtain the related GO terms for this set of proteins. These proteins are characterized by a  
215 distinct function compared to proteins found in other clusters. On the one hand, the proteins that are  
216 located in this cluster play a crucial role in the intestinal absorption of phytosterol and in cholesterol and  
217 lipid transportation. On the other hand, the remaining proteins, which are not found in this discrete  
218 cluster, are responsible for the homeostatic and transducer/receptor mechanisms of the cells.

## 219 Overall and Site Likelihoods

220 By converting PrInS distance matrices into rate matrices, it was possible to compute the overall likelihood  
221 and site likelihoods for every protein alignment. The resulting matrices from the *nEd*, *nMd* and *nPSm*  
222 are named as *nEs*, *nMs* and *nPSs*, respectively. For most of the proteins, 170 out of 223, the likelihoods  
223 computed by BLOSUM62 and PAM120 were better than our neighborhood Euclidean-based substitution  
224 matrices. This result is expected since BLOSUM62 and PAM120 substitution matrices were constructed  
225 based on alignment files, whereas the presented substitution matrices are alignment-unaware. For the  
226 remaining 54 protein alignments, our models and more specifically the Euclidean-based rate matrix  
227 resulted in higher likelihoods.

228 For every protein, we calculated the individual site likelihood for each amino acid site in the multiple  
229 sequence alignment using either the BLOSUM62 (or PAM120) or the *nEs* matrix. Results are illustrated  
230 in Figure 4) for the protein NCKX1, where the likelihood difference between BLOSUM62 and *nEs* is  
231 plotted for every amino acid. Thus, negative values are related to sites where the likelihood is greater  
232 for the *nEs* matrix, whereas positive values show sites where BLOSUM62 results in greater likelihoods.  
233 For the protein NCKX1, 240 sites are scored with a higher likelihood using the Euclidean substitution  
234 matrix derived from PrInS (negative values), whereas 390 sites are scored with a higher likelihood us-  
235 ing BLOSUM62 (positive values). The overall likelihood difference for this protein is positive (264.4)  
236 indicating that overall BLOSUM62 results in higher likelihood scores. NCKX1 was randomly selected  
237 to illustrate the site likelihood differences between the two approaches. It is a critical component of the  
238 visual transduction cascade, controlling the calcium concentration of outer segments during light and  
239 darkness McKiernan and Friedlander (1999).

## 240 Comparison between the average substitution likelihoods for different rate 241 matrices

242 The 20 amino acids form  $\binom{20}{2} = 190$  (unordered) pairs. For each of them, we have calculated the  
243 average differences between the BLOSUM62-based and the *nEs*-based likelihoods. For the calculations,  
244 we considered only the sites that consist of two amino acid states. Since all the parameters of the  
245 evolutionary model, but the substitution rate matrix are fixed (phylogenetic tree, equilibrium frequencies),  
246 then the matrix that results in the highest average likelihood for a certain amino acid pair describes the  
247 preferential model for this amino acid pair. Surprisingly, even though most of the protein alignments are  
248 scored with a higher likelihood using BLOSUM62, 111 (out of 190) amino acid pairs are characterized  
249 by negative average likelihood difference. This indicates that the *nEs*-based likelihood is greater than  
250 the BLOSUM62-based likelihood. The average likelihood is greater for the BLOSUM62 for only 50 pairs,  
251 whereas it cannot be assessed for 29 pairs because no occurrences of these 29 pairs were present in the

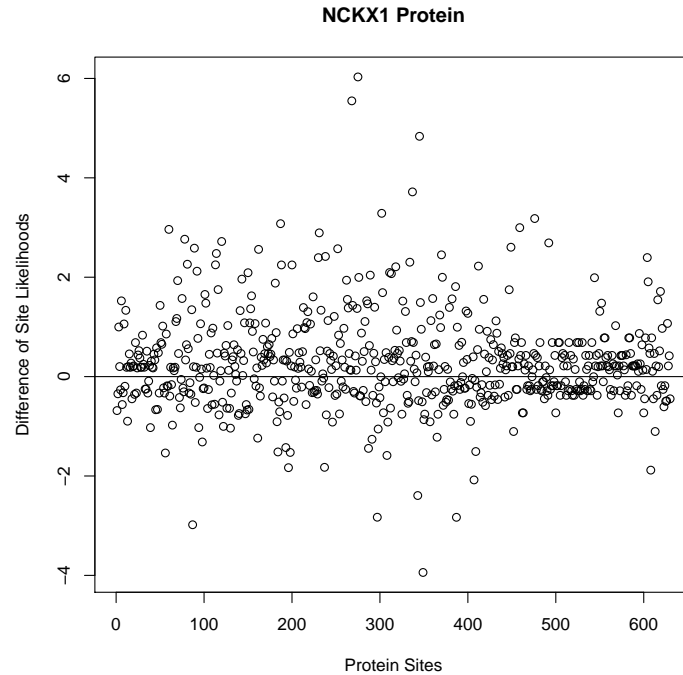


Figure 4: The differences between the BLOSUM62-based and the Euclidean-based likelihoods for each site of the NCKX1 protein is shown. On the x axis, the protein sites are shown, while on the y axis the differences between the likelihoods are depicted. Positive differences indicate higher likelihoods for the BLOSUM62 rate matrix, whereas negative differences indicate higher likelihoods for the Euclidean rate matrix.

252 alignment datasets. To scrutinize further this result, we plotted the average likelihood difference as a  
253 function of the occurrence frequency of the amino acid pair in the alignment (Figure 5). In Figure 5, it is  
254 apparent that the average likelihood difference is positively correlated with the frequency of occurrences  
255 of the amino acid pairs ( $r = 0.621$ ,  $r$  is the correlation coefficient, CI: (0.55, 0.67), CI denotes the 95%  
256 confidence intervals). In other words, the more frequent a substitution is in the alignments, the higher  
257 the difference of likelihoods, favoring the BLOSUM62 versus the *nEs*.

258 The heatmap in Figure 6 shows all amino acid pairs and the likelihood differences between the BLO-  
259 SUM62 and the *nEs* matrix for all amino acid pairs. Cells denoted by a ‘-’ are characterized by a higher  
260 likelihood for the Euclidean-based matrix, whereas a ‘+’ denotes cells with higher BLOSUM62-based  
261 likelihood. The darker the color of the cell the higher the difference between the two likelihoods. Cells  
262 with an ‘o’ illustrate pairs where no site with this pair of amino acids was found in any of the alignments.  
263 The greatest value for a ‘-’ cell is represented for the ‘Arginine-Glycine’ pair. This means that the *nEs*  
264 matrix is a preferential model for the specific amino acid pair. On the other hand, the pair with the  
265 highest ‘+’ value is the ‘Isoleucine-Valine’ pair.

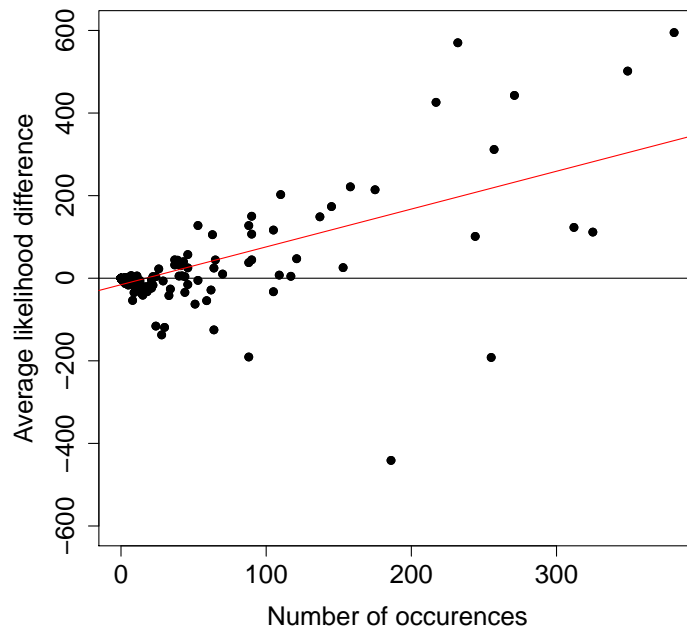


Figure 5: The average likelihood difference between the BLOSUM62-based and the Euclidean-based approaches that have been used to calculate likelihoods. As the figure indicates, for the majority of amino acid pairs the likelihood difference is close to 0, *i.e.*, both the Euclidean-based calculations and the BLOSUM62-based calculations result in similar outcomes. Positive values are fewer than negative values (see text), however the magnitude is much greater for positive values than for negative, indicating that BLOSUM62 results in much greater likelihoods than the Euclidean-based approach. Furthermore, as the occurrence frequency of the amino acid increases the difference of the likelihood increases as well, suggesting that BLOSUM62 outperforms the Euclidean-based approach for the amino acid pairs that occur frequently, either within the same protein or in different proteins.

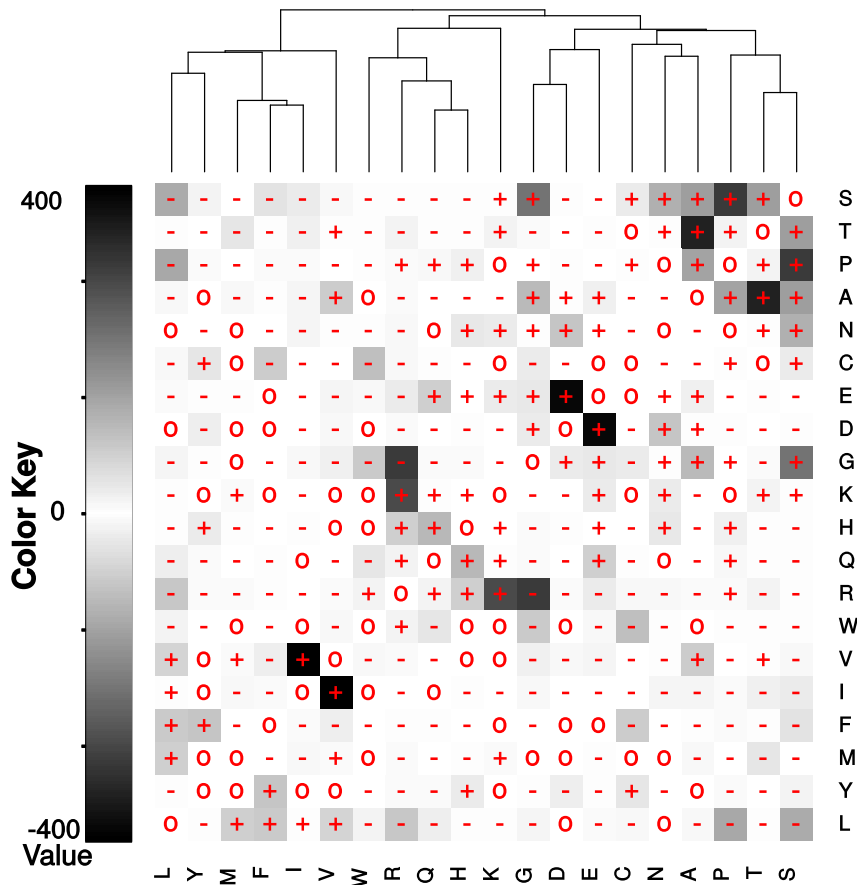


Figure 6: Comparison of likelihoods between BLOSUM62 and *nEs* for all amino acid pairs found in the protein alignments. Boxes with a '-' pinpoint to the specific amino acid pairs where *nEs* calculates a higher site likelihood. Contrarily, a '+' indicates amino acid pairs that the BLOSUM62 approach results in greater likelihood. Finally, 'o' cells depict the amino acid pairs that were not found in the multiple alignments, thus no comparison was possible. The darker the tone of the cell the greater the difference in the likelihood between the BLOSUM62 and the *nEs* approach.

## 266 Comparison to Grantham's and Sneath's distance matrices

267 Grantham's distance Grantham (1974) between two amino acids depends on three properties: composition  
268 (defined as the atomic weight ratio of noncarbon elements in end groups or rings to carbons in the side  
269 chain), polarity and molecular volume. Based on these three properties, distance  $D_{G[i,j]}$ , between amino  
270 acids  $i$  and  $j$ , is defined as:

$$D_{G[i,j]} = [\alpha(c_i - c_j)^2 + \beta(p_i - p_j)^2 + \gamma(v_i - v_j)^2]^{1/2} \quad (3)$$

271 where  $c$  is the composition,  $p$  the polarity and  $v$  the molecular volume. The three components of the  
272 distance ( $c, p, v$ ) are not independent. Thus, the constants  $\alpha, \beta, \gamma$  serve as normalizing factors and they  
273 can be calculated as a function of  $c, v, p$  Grantham (1974). Grantham (1974) provides the distances of  
274 Equation 3 for all amino acid pairs, thus a distance matrix  $D$ . Furthermore, Grantham demonstrated  
275 that the relative substitution frequency of amino acid pairs and their distances  $D_{G[i,j]}$ 's are correlated,  
276 underlying the physicochemical basis of amino acid substitution.

277 Similarly to Grantham's amino acid distances, Sneath's index Sneath (1966) takes into account 134  
278 categories of activity and structure. The dissimilarity index  $D_{S[i,j]}$  is the percentage of the sum of all  
279 properties not shared between two amino acids.

280 The amino acid neighborhood-based distance is correlated with both the Grantham's and Sneath's  
281 distances Figure 7, highlighting the fact that amino acid neighborhoods capture information related to  
282 the physicochemical properties of amino acids and also their relative substitution frequency.

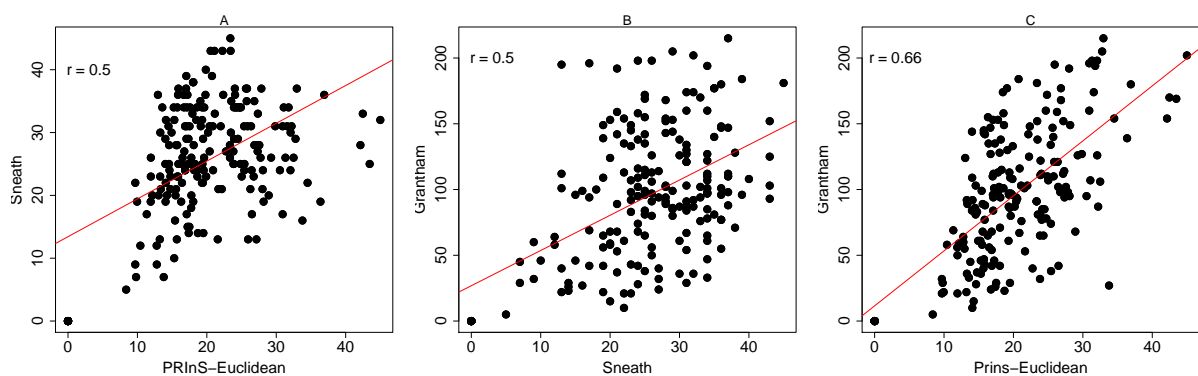


Figure 7: Scatterplots between amino acid distances for different distance methods. (A) Euclidean-based distance versus Sneath's index ( $r \approx 0.5$ ). (B) Sneath's index versus Grantham's distance ( $r \approx 0.5$ ). (C) Euclidean-based distance versus Grantham's distance ( $r \approx 0.66$ ).

## 283 Discussion

284 We studied the spatial properties of amino acid substitutions. We investigated if amino acid substitu-  
285 tions can be predicted from their neighborhood tertiary structure. We described statistically the amino  
286 acid residual neighborhoods using the software PrInS. Then, we converted PrInS output files to amino  
287 acid distance matrices and substitution rate matrices and evaluated their ability to model evolutionary  
288 changes.

289 Correlation analysis between the observed substitution frequencies (from multiple sequence align-  
290 ments) and the distance matrices indicated that residual neighborhoods capture evolutionary information  
291 and thus they can be useful in modeling evolution. In other words, substitutions in multiple sequence  
292 alignments can be predicted from the amino acid neighborhoods: residuals with similar neighbors can  
293 substitute each other, with a higher rate, during evolution. More specifically, from the three distance  
294 matrices (Euclidean, Manhattan, and Pearson Squared) that were compared against the substitution ma-  
295 trix, only five amino acids were discordant in at least two comparisons. These amino acids were proline,  
296 arginine, serine, tryptophan and tyrosine. From the substitution matrix, only a few substitutions involve  
297 tryptophan and tyrosine, whereas there is a multitude of substitutions involving proline, arginine and  
298 serine. In contrast, based on distance matrices, in the columns of tryptophan and tyrosine there are small  
299 distance values, while in the proline, arginine and serine columns distances are large. This can explain  
300 the discordance of these amino acids in these two matrices, because a large number of substitutions is  
301 associated with small distance values for an amino acid.

302 In addition, the Common Substitution Tool of NCBI Amino Acid Explorer Bulka et al. (2006) enhances  
303 our main result: Amino acid pairs with the lower distance values in the distance matrices, especially using  
304 the *nEd*, are found to substitute each other more frequently. Amino Acid Explorer returned a list of  
305 the amino acids from the most common to the the less common substitution for each amino acid. There  
306 are some differences in the order of the amino acids between the Common Substitution Tool and the  
307 distance matrices (e.g Euclidean). These differences are possibly based on the fact that the generation  
308 of the *nEd* matrix is based solely on the structural information of a protein family ( $\alpha$ -helical membrane  
309 proteins), while the Common Substitution Tool is based on BLOSUM62 substitution matrix that was  
310 created from alignments from multiple protein families. Thus, a plausible explanation is that the amino  
311 acid neighborhood approach we followed cannot capture adequately well substitutions of certain amino  
312 acids. Another plausible explanation is that we followed a linear approach to convert the amino acid  
313 distance matrix to the substitution rate matrix. Thus, if the distance between a pair of amino acids  $A_1$ ,  
314  $B_1$  is  $d_1$ , whereas the distance of amino acids  $A_2$  and  $B_2$  is  $d_2$ , then the relation between the substitution  
315 rates of  $A_1, B_1$  and  $A_2, B_2$  will be a linear function of the  $d_1$  and  $d_2$ . A non linear function between the  
316 distances and the substitution rates will perhaps result in more accurate substitution rates.

317 An interesting observation springs from Figure 5 where amino acid substitutions that occur very

318 frequently in binary sites of alignments (*i.e.* sites with two states) are better predicted with BLOSUM62  
319 than *nEs* matrix. The relation between the average likelihood difference between BLOSUM62 and *nEs* is  
320 positively correlated with the occurrence frequency of amino acid substitution. Thus, the neighborhood  
321 based approach cannot predict well substitutions that occur frequently in an alignment. A plausible  
322 solution would be to modify *nEs* taking into account the frequency of the involved amino acids (which  
323 is presumably related to the substitution frequency between them).

324 The neighborhood distance matrices are highly correlated with the Grantham's and Sneath's dis-  
325 tance matrices, illustrating that the formation of amino acid neighborhoods is determined by structural  
326 constraints. Grantham Grantham (1974) demonstrated that amino acid substitution is related to their  
327 structural properties. Thus, a plausible explanation of the amino acid substitution prediction ability of  
328 *nEd* approach is that neighborhood properties are defined by the structural properties of amino acids  
329 which also affect the substitutions between amino acids.



## References

- 330
- 331 Frances C Bernstein, Thomas F Koetzle, Graheme JB Williams, Edgar F Meyer, Michael D Brice, John R  
332 Rodgers, Olga Kennard, Takehiko Shimanouchi, and Mitsuo Tasumi. The protein data bank. *The FEBS*  
333 *Journal*, 80(2):319–324, 1977.
- 334 Blazej Bulka, Stephen J. Freeland, and others. An interactive visualization tool to explore the biophysical  
335 properties of amino acids and their contribution to substitution matrices. *BMC bioinformatics*, 7(1):  
336 329, 2006. URL [https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-](https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-7-329)  
337 [7-329](https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-7-329).
- 338 In-Geol Choi and Sung-Hou Kim. Evolution of protein structural classes and protein sequence families.  
339 103(38):14056–14061, 2006. URL <http://www.pnas.org/content/103/38/14056.short>.
- 340 MO Dayhoff, RM Schwartz, and BC Orcutt. 22 a model of evolutionary change in proteins. In *Atlas*  
341 *of protein sequence and structure*, volume 5, pages 345–352. National Biomedical Research Foundation  
342 Silver Spring, MD, 1978.
- 343 Martin Egli. Diffraction techniques in structural biology. In Serge L. Beaucage, Donald E. Bergstrom, Piet  
344 Herdewijn, and Akira Matsuda, editors, *Current Protocols in Nucleic Acid Chemistry*, pages 7.13.1–  
345 7.13.35. John Wiley & Sons, Inc., 2010. ISBN 978-0-471-14270-6. URL [http://doi.wiley.com/10.](http://doi.wiley.com/10.1002/0471142700.nc0713s41)  
346 [1002/0471142700.nc0713s41](http://doi.wiley.com/10.1002/0471142700.nc0713s41). DOI: 10.1002/0471142700.nc0713s41.
- 347 Eric Franzosa and Yu Xia. Structural perspectives on protein evolution. *Annu Rep Comput Chem*, 4(1):  
348 3–21, 2008.
- 349 R Grantham. Amino acid difference formula to help explain protein evolution. *Science*, 185(4154):  
350 862–864, 1974.
- 351 Leslie Hatton and Gregory Warr. Protein structure and evolution: are they constrained globally by a  
352 principle derived from information theory? 10(5):e0125663, 2015. URL [http://journals.plos.org/](http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0125663)  
353 [plosone/article?id=10.1371/journal.pone.0125663](http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0125663).
- 354 Steven Henikoff and Jorja G. Henikoff. Amino acid substitution matrices from protein blocks. *Proceed-*  
355 *ings of the National Academy of Sciences*, 89(22):10915–10919, 1992. URL [http://www.pnas.org/](http://www.pnas.org/content/89/22/10915.short)  
356 [content/89/22/10915.short](http://www.pnas.org/content/89/22/10915.short).
- 357 W. James Kent, Charles W. Sugnet, Terrence S. Furey, Krishna M. Roskin, Tom H. Pringle, Alan M.  
358 Zahler, and David Haussler. The human genome browser at UCSC. 12(6):996–1006, 2002. URL  
359 <http://genome.cshlp.org/content/12/6/996.short>.

- 360 Colleen J. McKiernan and Martin Friedlander. The retinal rod  $\text{Na}^+/\text{Ca}^{2+}$ ,  $\text{K}^+$  exchanger contains a  
361 noncleaved signal sequence required for translocation of the n terminus. 274(53):38177–38182, 1999.  
362 URL <http://www.jbc.org/content/274/53/38177.short>.
- 363 A. Nath Jha, S. Vishveshwara, and J. R. Banavar. Amino acid interaction preferences in helical membrane  
364 proteins. 24(8):579–588, 2011. ISSN 1741-0126, 1741-0134. doi: 10.1093/protein/gzr022. URL <https://academic.oup.com/peds/article-lookup/doi/10.1093/protein/gzr022>.  
365 <https://academic.oup.com/peds/article-lookup/doi/10.1093/protein/gzr022>.
- 366 Jri Reimand, Meelis Kull, Hedi Peterson, Jaanus Hansen, and Jaak Vilo. g:profiler web-based toolset for  
367 functional profiling of gene lists from large-scale experiments. 35:W193–W200, 2007. ISSN 1362-4962,  
368 0305-1048. doi: 10.1093/nar/gkm226. URL <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkm226>.  
369 <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkm226>.
- 370 Jessica Siltberg-Liberles, Johan A. Grahnen, and David A. Liberles. The evolution of protein structures  
371 and structural ensembles under functional constraint. 2(4):748–762, 2011. ISSN 2073-4425. doi:  
372 10.3390/genes2040748. URL <http://www.mdpi.com/2073-4425/2/4/748/>.
- 373 PHA Sneath. Relations between chemical structure and biological activity in peptides. *Journal of*  
374 *theoretical biology*, 12(2):157–195, 1966.
- 375 Catherine L. Worth, Sungsam Gong, and Tom L. Blundell. Structural and functional constraints in  
376 the evolution of protein families. 2009. ISSN 1471-0072, 1471-0080. doi: 10.1038/nrm2762. URL  
377 <http://www.nature.com/doi/10.1038/nrm2762>.