

1 Where Natural Protein Sequences 2 Stand out From Randomness

3 **Laura Weidmann¹, Tjeerd Dijkstra², Oliver Kohlbacher^{2,3,4,5,6}, Andrei Lupas¹**

*For correspondence:

lweidmann@tuebingen.mpg.de

4 ¹Department of Protein Evolution, Max Planck Institute for Developmental Biology,
5 Tübingen, Germany; ²Biomolecular Interactions, Max Planck Institute for Developmental
6 Biology, Spemannstr. 35, 72076 Tübingen, Germany.; ³Applied Bioinformatics,
7 Department for Computer Science, University of Tübingen, Sand 14, 72076 Tübingen,
8 Germany; ⁴Institute for Biomedical Informatics, University of Tübingen, Sand 14, 72076
9 Tübingen, Germany; ⁵Center for Quantitative Biology, University of Tübingen, Sand 14,
10 72076 Tübingen, Germany; ⁶Translational Bioinformatics, University Hospital Tübingen,
11 Hoppe-Seyler-Str. 9, 72076 Tübingen

12
13 **Abstract** Biological sequences are the product of natural selection, raising the expectation that
14 they differ substantially from random sequences. We test this expectation by analyzing all
15 fragments of a given length derived from either a natural dataset or different random models. For
16 this, we compile all distances in sequence space among fragments of each dataset and compare
17 the resulting distance distributions. Even for 100mers, 95.4% of all distances between natural
18 fragments are in accordance with those of a model based on the natural residue composition.
19 Hence, natural sequences are distributed almost randomly in global sequence space. When further
20 accounting for the specific residue composition of domain-sized fragments, 99.2% of all distances
21 between natural fragments can be modeled. Local residue composition, which might reflect
22 biophysical constraints on protein structure, is thus the predominant feature characterizing
23 distances between natural sequences globally whereas homologous effects are only barely
24 detectable.

26 Introduction

27 Natural proteins form the backbone of the complicated biochemical network that has given rise to
28 the great variety of life on Earth. This highly interwoven framework of reactions seems impossible
29 to have arisen by chance, simply because the great majority of random protein sequences fails to
30 form a specific structure, let alone possess chemical activity. Features that distinguish naturally
31 evolved from random sequences are therefore of great interest, both in order to understand protein
32 evolution *Shah et al. (2015)* *Luigi Luisi (2003)* and to guide the design of new proteins *Woolfson*
33 *et al. (2015)* *Pande et al. (1994)*.

34 Searches for such differences have hitherto focused on the exhaustive enumeration of short
35 peptides and their statistical analysis by exact occurrence *Poznański et al. (2018)* *Lavelle and*
36 *Pearson (2009)*. These studies showed that the natural frequency of most peptides is similar to
37 that expected from random sequences with the same composition. Nevertheless, the frequency of
38 some peptides was found to deviate substantially from random occurrence, an observation which
39 was variously discussed in terms of homologous descent and convergence due to structural and

40 functional constraints. This enumeration approach quickly reaches its limits at sequence lengths
41 above 5, due to the fact that there are simply not enough natural sequences to populate the
42 exponentially growing sequence space. Furthermore, pentapeptides are far from having a relevant
43 length for understanding protein sequences. Even if proteins are dissected into their constituent
44 domains, which form the units of folding and in general also of functional activity, relevant sequence
45 lengths still mostly range above 80 residues. Reaching sequences of this length, the complexity of
46 20^{80} needs to be drastically reduced to comprehend the global occupation of sequence space by
47 nature.

48 Nevertheless, decades of bioinformatic research have allowed us to form expectations about this
49 occupation of sequence space by domain-sized natural sequences. This is because most proteins
50 have arisen by descent and differentiation from a set of domain prototypes, and can thus be
51 classified into a hierarchy of domain families and superfamilies. This points to the fact that the
52 sequence space around existing domains is substantially populated by homologs. Is this image of
53 sequence space being populated by islands formed around domains families and superfamilies
54 representative for the global structure of naturally occupied sequence space?

55 Against our own initial expectations, we demonstrate in this paper, that this is not exactly the case.
56 The main indicators for this can be seen in the non-trivial task to identify homology and the often
57 only probabilistic estimates on the homologous relationship between sequences. Most homologs
58 are not obviously related to each other as they share a *randomly expected similarity Rost (1999)*,
59 making it hard to substantiate homologous relationships from random fluctuation or convergence.
60 Only advanced methods that estimate the *significance* of similarity among *multiple* sequences after
61 repeated searches [HHpred, PSIBlast], may succeed in detecting long-range common ancestry. Such
62 elaborated methods cluster sequences according to their estimated *evolutionary distance* and give
63 rise to the picture of islands formed by related sequences. This picture does however not reflect
64 the distribution of sequences across sequence space, as therein distances between sequences
65 represents a proxy for their evolutionary distance, not their actual distance in sequence space.
66 Although the impact of homology to the very local structure of sequence space is undoubtedly
67 significant, globally it may thus not be traceable. However, the contribution of homology to the
68 global space has not been studied and is substantially unclear.

69 A step towards a more in-depth understanding of the local space and the extent of homology has
70 been taken with searches for variants close to existing proteins *Bershtein et al. (2017) Starr et al.*
71 *(2017) Harms and Thornton (2014) Urlinger et al. (2000)* [Olivers paper from Sander]. By testing
72 exhaustively all mutations at certain sites, these studies bypass intermediate mutants that would
73 not have been viable in evolution. Contrasting the abundance of possible functional variants to
74 the small number of natural sequences demonstrates how sparsely nature has explored sequence
75 space, even locally. The high energy barriers, epistatic effects, and functional dependencies prevent
76 random mutations and seem to entrench already existing and functional forms *Starr and Thornton*
77 *(2016) Shah et al. (2015)*.

78 Modern techniques of protein design allow to reach out further into the global sequence space
79 to find possible exemplars in unknown territory *Huang et al. (2016) Woolfson et al. (2015)*. Scaling
80 these scans up to the currently highest practicable level for a given structure or function has
81 uncovered viable solutions far from existing proteins *Larson et al. (2002) Chevalier et al. (2017)*,
82 showing that sequence similarity is not required. This leads to the hypothesis that the usable part
83 of sequence space is mostly randomly structured, which has also been proposed for unrelated
84 natural sequences before *Lavelle and Pearson (2009)*.

85 Apart from the seemingly random global structure, there are nevertheless, biophysical requirements
86 for all usable protein sequences, natural as well as designed, such as foldability, hydrophobic core
87 formation and solubility. This indicates that, although randomly arrayed in global sequence space,

88 these proteins may still share some convergent features, which would restrict a random drift away
89 into unstructured space. Natural sequence space could thus be characterized globally by sequences
90 with the potential to fold, i.e. by convergent features. In contrast the contribution of homology to
91 this global space has not been studied and is substantially unclear.

92 In this paper, we analyze the global structure of natural sequence space, aiming to identify evolution-
93 ary footprints and general features that characterize natural sequences. We do this by contrasting
94 natural data with a variety of random models, in order to extract sequence features arising from
95 different natural mechanisms.

96 **Results and Discussion**

97 **Natural sequence data and random sequence models**

98 Choice of a natural dataset

99 For an adequate dataset that reflects the natural protein sequence space, we aimed to achieve a
100 reasonable coverage of deep phylogenetic branches with complete and well-annotated proteomes.
101 Given that the genome coverage for the archaeal and eukaryotic lineages is still sparser than for
102 bacteria and that particularly eukaryotic genomes are affected by issues of assembly, gene detection,
103 and intron-exon boundaries, we built our database from the derived bacterial proteomes collected
104 in UniProt *Apweiler (2009)*. To control for redundancy, we selected only one genome per genus and
105 filtered each for identical open reading frames and low-complexity regions. In total our dataset
106 comprises 1,307 genomes, $4.7 \cdot 10^6$ proteins, and $1.2 \cdot 10^9$ residues. We simplified complexities
107 arising from the use of modified versions of the 20 proteinogenic amino acids, which occurred
108 in a few hundred cases, by converting these to their unmodified precursors, thus maintaining an
109 alphabet of 20 characters throughout. Further details on the generation of our dataset and its
110 specific content are provided in the methods section.

111 In order to evaluate where our natural dataset differs from randomness, we developed a series of
112 increasingly specific models that account for compositional effects.

113 How random is random?

114 Our most basal random model considers completely random sequences of the 20 proteinogenic
115 amino acids, in which each occurs with an equal probability of 5% (E-model). This model is known
116 to approximate natural sequences only poorly *de Lucrezia et al. (2012) Munteanu et al. (2008)*.
117 This is hardly surprising as natural amino acid frequencies in fact range between 1% and 10%,
118 a bias which is associated with metabolic pathways, bio-availability, and codon frequency. We
119 therefore built models that factor in this compositional bias at increasingly local levels. The first
120 model incorporates the global amino acid composition of our natural dataset, which we refer to as
121 the A-model.

122 More specific models consider increasingly local fluctuations in composition. The composition of
123 different genomes, for example, varies with GC-content and environmental influences *Fukuchi and*
124 *Nishikawa (2001) Fukuchi et al. (2003)*. This effect can be factored in using the individual genome
125 composition (G-model). With an increasingly local focus, compositional bias can be accounted for
126 at the level of proteins (P-model) *Chou (2001) Cedano et al. (1997)*, domains (D-model) *Lavelle and*
127 *Pearson (2009)* and even sub-domain-sized fragments *Poznański et al. (2018)*.

128 Having accounted for compositional effects resulting from environment, metabolism, and the need
129 to form a hydrophobic core, the remaining differences between natural and random sequences
130 must be attributed to sequence effects, due either to *divergence* from a common ancestor *Alva*
131 *et al. (2015)* or *convergence* as a result of secondary structure formation *Pande et al. (1994)*.

Table 1. Random sequence models based on amino acid composition.

Model	natural feature	class of feature
E	natural amino-acid alphabet, equal propensity for each letter	single, overall descriptor
A	overall amino acid composition	
T	overall dipeptide frequency	
G	composition of individual genomes	context-specific composition
P	composition of proteins	
D	composition of domain-sized fragments	
D1	D-model + homology sequence bias	mixed models that incorporate sequence bias
D2	D-model + analogy sequence bias	
D3	D-model + homology and analogy sequence bias	

Table 1-source data 1. Random sequence models. Completely random sequences, where each amino acid occurs with the same probability of 5%, are represented by the E-model. The natural frequency of specific amino acids deviates remarkably from such an equal distribution, thus, random sequence models are usually based on the overall amino acid composition, represented by the A-model. The overall dipeptide frequency is considered by the T-model. The diversity of amino acid composition across genomes, is accounted for by the G-model. On a more specific level, the composition occurring in natural proteins or even domain-sized fragments can be used to generate random sequence models, here referred to as P- and D-models. In order to estimate the contribution of analogous and homologous relationships to the global occupation of sequence space, we generated models D1, D2 and D3 that include sequence bias in addition to the composition bias of the D-model. (These models will be explained in detail in the last section of the Results.) We compare our natural dataset to all of these models and illustrate to what extent they differ from the natural sequence space occupation. Our implementation of the models are described in the Methods section.

132 **Representing sequence space occupation based on pairwise distances**

133 Sequence space has frequently been analyzed with a direct approach based on the exhaustive
134 enumeration of natural kmers, and the comparison of their frequencies to those derived from
135 a random model *Poznański et al. (2018)* *Lavelle and Pearson (2009)*. This approach is restricted
136 to kmers of length 5 or smaller, due to sequence space complexity and the data sparsity caused
137 thereby. It also does not represent the relative position among kmers within the global sequence
138 space, given that it focuses on frequencies of exact 5mers, which correspond to single points in the
139 5mer sequence space.

140 We use an indirect approach to circumvent these problems. Our approach is built on the probability
141 mass function of pairwise distances between sequences of the same length, in the following referred
142 to as *distance distribution*. A distance distribution illustrates how often sequences are positioned
143 at a certain distance to each other and we use it to study the way sequences are spread across
144 the possible space. We built distance distribution for the natural dataset and for each dataset
145 of random sequences derived from specific models. with a length of up to 100 residues, thereby
146 covering the average domain size *Wheelan et al. (2000)*.

147 As a metric for distance, we use the *normalized local alignment score* of a Smith-Waterman alignment
148 since this metric is commonly used to capture similarities between natural sequences *Rost (1999)*
149 *Schneider et al. (1997)*. We note, that the choice of distance metric is not of great relevance for the
150 main implications of our study; relative to each other, the distance distributions of the random
151 models deviate similarly from that of natural sequences irrespective of the chosen metric, as
152 illustrated in *Figure 2*. Details on the derivation of distance distributions and the used distances
153 metrics are provided in the Methods.

154 Common methods that relate sequences to each other are mainly based on the significance
155 of specific similarities among all existing sequences *Alva et al. (2009)* to estimate evolutionary

156 distances. Our method differs from these approaches, as it only considers the pairwise similarity
157 between two sequences, reflecting their distance in sequence space without including external
158 information derived from other sequences. This has consequences for the way sequence clusters
159 are perceived, and we outline this in detail by analyzing the distance distribution of homologs in the
160 last chapter of the results.

161 Studying the layout of space through pairwise distances are common in other fields, such as protein
162 structure determination *Wüthrich (1986)*, spatial statistics *Diggle (2014)* and economics *Duranton*
163 *and Overman (2005)*, but have not, to our knowledge, been applied to investigate protein sequence
164 space. We note that, as all distance-based methods, this method loses all information about
165 the positions of considered sequences in space. This entails the effect that identical distance
166 distributions can be derived from multiple distinct sequences; different occupations of sequence
167 space, hence, can lead to the same distance distributions. However, in exchange for the positional
168 information, this method can characterize the global structure of how sequence space is being
169 occupied, which includes *global clustering and dispersion* of sequences. Grasping the global nature of
170 sequences demands to reduce the great complexity of sequence space drastically, and our methods
171 succeeds in this by losing positional information. Furthermore, this global consideration gives
172 less weight to local structures and we aim to detect general features that distinguish natural from
173 random sequences exclusively.

174 **Comparing distance distributions**

175 For the comparison of the natural to a random distance distribution, we first subtract the fraction of
176 distances observed in the random dataset from that observed in the natural dataset for each possi-
177 ble alignment score. We refer to this difference as the *residual*. Over all sequence identity scores,
178 residuals sum up to zero and may have values that are either positive (more natural distances)
179 or negative (more random distances). In order to obtain an overall measure of how different two
180 distances distributions are, we derive the variational distance between the distributions, referred
181 to as the *total residual*. More precisely, the total residual is the sum over the absolute residuals,
182 normalized to a range between 0% and 100%.

183 If the two distance distributions are completely non-overlapping, the total residual assumes the
184 maximal value of 100%, indicating that no distance between natural fragments can be modeled
185 with the underlying random sequences. If they are identical, the total residual assumes a value of
186 0%, indicating that 100% of all distances in the natural distribution have a corresponding distance
187 in the random distribution. Thereby, the total residual represent the fraction of natural distances
188 that are not accounted for by the distance distribution of a random model.

189 **Amino acid composition (A-model)**

190 We start our analysis by assessing to what extent the global amino acid composition, as captured
191 in the A-model, can describe natural sequences. We compare the distance distributions of the
192 two datasets for fragment lengths up to 100 residues, in increments of 10. At all fragment lengths,
193 the results are closely comparable. We present the results for 100mers as representative for
194 domain-sized sequences *Figure 1* and provide the others in the supplementary figures.

195 The distance distributions of natural and A-model data overlap extensively (*Figure 1: A*). Both
196 are uni-modal with a peak at a low alignment score of 11%. Their minor differences only become
197 apparent, when their residuals are considered (*Figure 1: B*). These take the shape of a wave, with two
198 crests at alignment scores of 9% and 15% (reflecting an over-representation of the corresponding
199 distances in the natural dataset), and a trough at 11% (reflecting an under-representation). The
200 over-representation of distances both longer and shorter than expected from the random model,
201 suggests that natural sequences are less homogeneously distributed in space. We rationalize this
202 effect with the observation that natural sequences are enriched in certain parts of sequence space,
203 leading to an increase in shorter distances. This may occur both in regions with rare amino acids

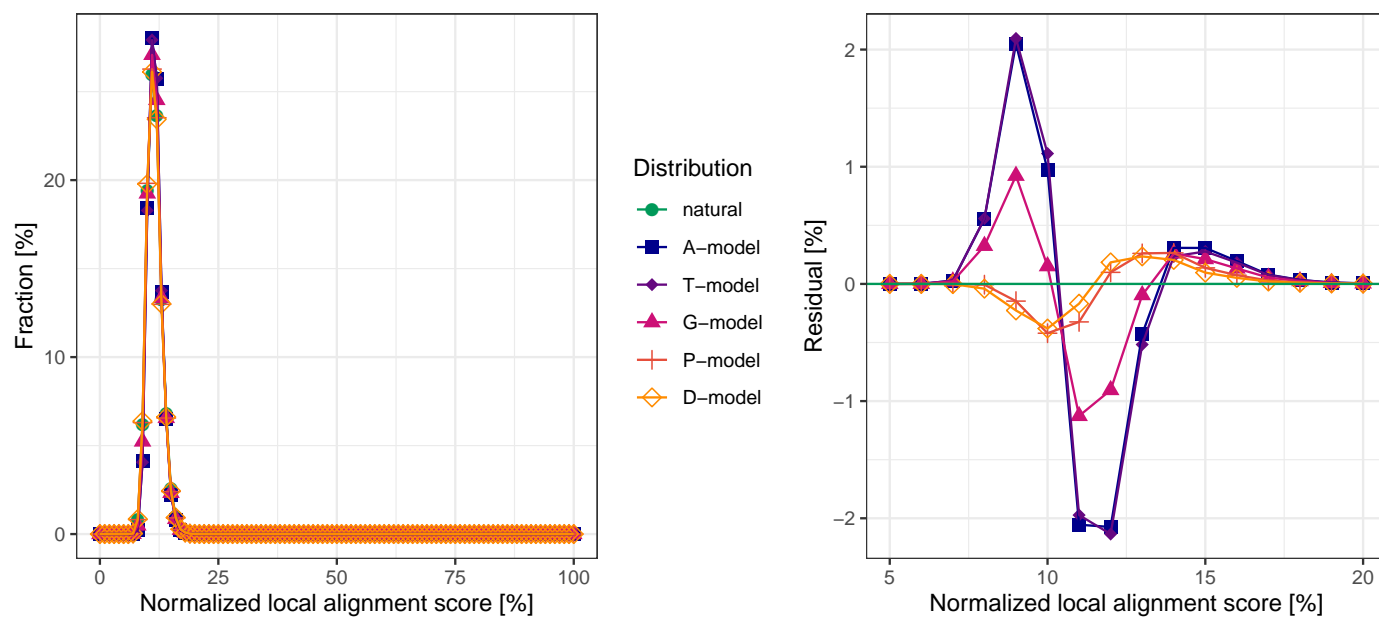


Figure 1. Comparing the sequence space occupation of random protein sequence models and natural sequence data. (A) Distance distributions are a descriptor of sequence space occupation. The distance between sequence fragments of the same length, defined as the sequence identity score obtained from a Smith-Waterman alignment, are plotted against the fraction of fragment pairs with the respective distance. We sampled 500 Million distances between fragments of length 100 for each model as well as for the natural sequence data. All distance distributions spike in the area of long-range distances with a mean sequence identity score around 11%. Both natural and random distance distributions are almost entirely overlapping. (B) Residuals represent the difference in sequence space occupation of random models compared to the natural sequences. We extract the distance-specific difference by subtracting the random from the natural distance distribution. The resulting residuals for each model indicate distances between natural fragments that are unaccounted for by the respective model (crests above zero). The A-, T- and G-model display a 2-peak behavior, associated with more long-range and short-range distances between natural fragments than modeled, reflecting an increased amount of both diversity and clustering in natural sequence space. The residuals of the P- and D-model possess only one peak for more short-range distances between natural sequences, hence an unexpected amount of clustering.

204 (such as Cys, Trp and His in small proteins dominated by zinc-coordination and disulfide bonds
205 *Vallee and Auld (1990)*) and in regions with abundant amino acids (such as Leu, Ala and Glu in
206 all-alpha proteins, most extremely in coiled coils *Lupas et al. (1991)*). The compositional differences
207 in these enriched regions mean that their distance in sequence space will be larger than expected
208 from the A-model, and thus lead to a complementary increase in longer distances. Since residuals
209 add up to zero, the number of intermediate distances is correspondingly decreased.

210 We note, however, that this discrepancy between natural sequences and the A-model is not very
211 pronounced, as the total residual has a value of only 4.6% for 100mers *Figure 2*. It is even less
212 pronounced at smaller fragment lengths, reaching 0.4% for 10mers. We conclude that the A-model
213 becomes less accurate in describing the sequence space occupation of natural sequences at lengths
214 that are biologically relevant, but that it already achieves considerably higher accuracy than the
215 completely random model (E-model), which has a total residual of 30.4% for 100mers (data shown
216 in Methods).

217 We evaluated whether adding sequence information to the unified compositional bias of the A-
218 model could further improve it. Since nature favors certain amino acid combinations as neighboring
219 residues, a model that reflects the natural dipeptide frequency (T-model), has been proposed to
220 represent natural sequences better than the A-model *Lavelle and Pearson (2009)*. We
221 implemented the T-model by extracting the dipeptide frequencies from our natural dataset and
222 using them to generate random sequences with a Markov Chain Model. For all fragment lengths,
223 we derived the distance distribution of the T-model (*Figure 1: A*), its residuals (*Figure 1: B*) and
224 the total residual (*Figure 2: A*, darkblue line). By all these measures the T- and the A-model
225 yielded essentially identical results in modeling the natural distance distribution. This outcome
226 was somewhat surprising, as the addition of dipeptide frequencies to the A-model did produce
227 a measurable improvement in an enumeration study of 5mers *Lavelle and Pearson (2009)*. This
228 may be due to the different methodology in that study, which collated exact 5mer frequencies,
229 corresponding to a position-wise Hamming distance of zero, and thus being close to a global, not
230 to a local alignment as used in our study. In fact, when using the Hamming distance as metric,
231 the T-model achieves a slightly better accuracy over the A-model for sequences of 50 or less
232 residues (*Figure 2: D*). From the results we obtained with the A- and T-models, we conclude that
233 global measures of composition and sequence bias already approximate natural sequences fairly
234 accurately, but that this accuracy decreases with sequence length. Especially for longer fragments,
235 we expect further improvement by including local compositional biases as outlined in the previous
236 section.

237 **Context-specific composition**

238 In order to capture context-dependent features, we investigated the effects of naturally occurring
239 local amino acid compositions. As a first step we considered a model that accounts for genome
240 diversity (G-model). Therein, the random dataset is produced by shuffling residues of the natural
241 dataset within the boundaries of each genome. Given that our natural dataset holds 1,307 genomes,
242 the derived sequences are thus sampled from 1,307 distinct compositions. Further locality is
243 achieved by accounting for the composition of individual proteins (P-model). Here, the random
244 dataset is produced by shuffling residues within each natural protein, corresponding to $4.7 \cdot 10^6$
245 compositions.

246 Since proteins are generally composed of domains, which are usually autonomous in structure
247 and also often in function, the next level of locality would be achieved by accounting for the com-
248 positional biases of individual domains. However, producing such a model is not straightforward,
249 as it is unclear how residues not assigned to a domain family should be taken into consideration.
250 Following upon this idea of the local composition defined by domains, we generated a D-model
251 that incorporates the local composition of natural fragments. For this, we used a fragment length

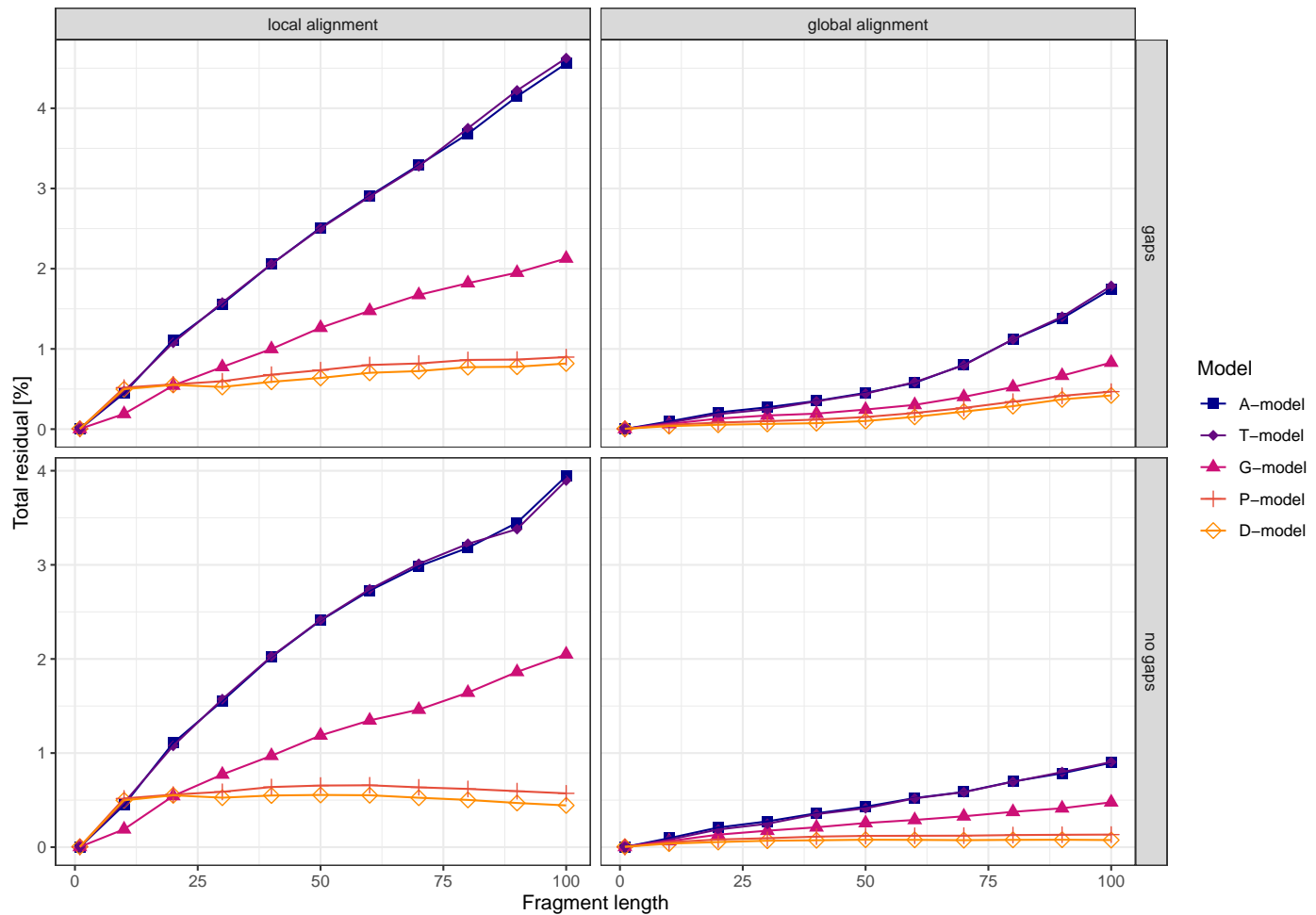


Figure 2. Deviation of random sequence models from natural sequences as a function of fragment length. (A) The total residual indicates the extent to which the distance distribution of random sequence models deviates from the natural. It reflects the fraction of distances between natural fragments that are unaccounted for in the random model. With increasing fragment length, the total residual of all models increases, implying that for longer fragments all models become worse in approximating similarities between natural fragments. The A-model (natural amino acid composition) and the T-model (dipeptide frequency) deviate furthest followed by the G-model (residue composition in genomes), the P-model (residue composition of proteins) and the D-model (residue composition of domain-sized fragments of length 100), which deviates the least. The intercept of the total residuals of the T- and D-model with the other models at fragment length around 10 is associated with edge effects of natural sequences and the usage of a local alignment as distance metric. (B) Total residuals when using a global Needleman-Wunsch alignment. The inconsistent continuation of the total residuals at sequence length 10 when using a local alignment has disappeared. Generally, the total residuals are reduced by 2.5-fold compared to the local alignment, reflecting that a global alignment captures less effects of natural sequences than a local alignment. (C) Total residuals when using a local Shift alignment. In contrast to the previously illustrated distance metrics, in the Shift alignment beginning and end-gaps are allowed without penalties but the possibility of internal gaps is excluded. Similar to the Smith-Waterman alignment, the Shift alignment displays an inconsistency at fragment length around 10. (D) Total residuals when using Hamming distance. The Hamming distance is similar to the Needleman-Wunsch based on a global alignment. It reflects the most stringent interpretation of similarity in sequence space, as the n-th position of one sequence is always compared with the n-th position of another sequence. It corresponds to a metric that considers the number of dimensions (positions in sequence) that are identical.

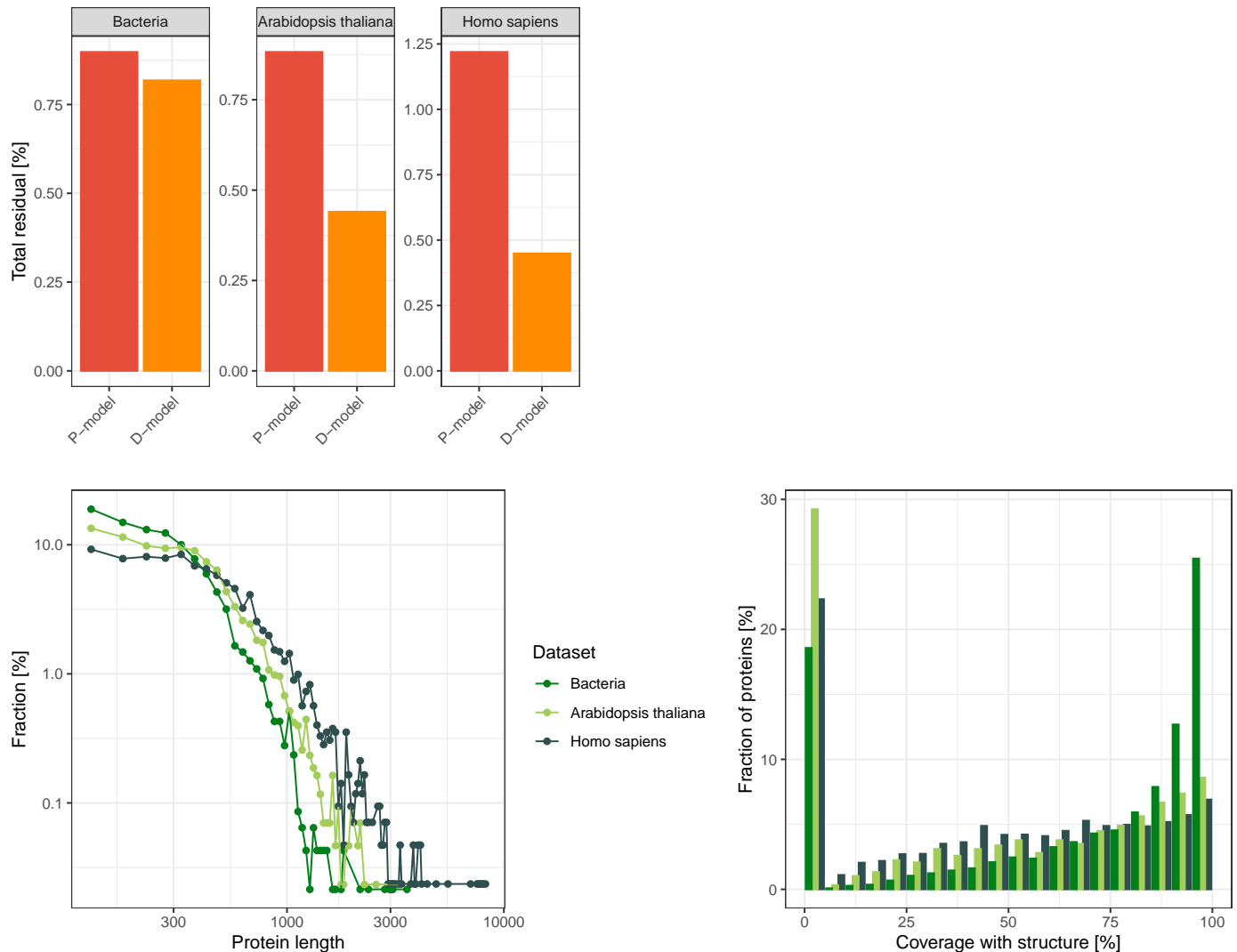


Figure 3. Contrasting the results of our bacterial dataset with those from two eukaryotic proteomes. (A) Total residuals of random models for bacterial dataset, the proteome of *Arabidopsis thaliana* and *Homo sapiens* of the P- and D-models. Relative to the total residual of the P-model, the total residuals of the D-models differ in the three presented datasets. In bacteria, both are almost identical, whereas for the eukaryotic datasets the D-models have a more than 2-fold increase in accuracy over the P-models. (B) Distribution of protein length. The median protein length is smallest for bacteria with 315 residues, 400 residues in the *Arabidopsis thaliana* dataset and 550 residues in the *Homo sapiens* dataset. The increase of median protein length correlates with the decrease in the total residual of the D-model relative to the P-model. (C) Coverage of proteins by structured domains. For each protein in the three datasets, an estimate of the coverage by structured domains was obtained by assigning ECOD families to regions in the protein. The fraction of residues within assigned domains compared to the protein length was obtained and plotted as a histogram over all sampled proteins. In bacteria 37% of the sampled proteins are almost completely structured (coverage of >90%), a fraction that is greater compared to that in *Arabidopsis thaliana* (15%) and *Homo sapiens* (13%).

252 of 100 residues, corresponding to an average domain length *Wheelan et al. (2000)*, that spans over
253 a major part of most protein sequences. We therefore considered the composition of all possible
254 fragments of length 100 from our natural dataset and shuffled their residues to derive sequence
255 fragments of the D-model (see Methods). Thus, we considered natural sequences that are not
256 exclusively part of a structured domain. This includes sequences that connect distinct domains, that
257 are part of non-globular regions (fibers, coiled-coils, amyloids) or that are intrinsically unstructured
258 regions.

259 Comparing the G-model to the A- and T-models over the bacterial dataset shows a dampened wave
260 for the residuals, with the same shape, but a decreased amplitude (*Figure 1: B*). The total residual
261 is correspondingly smaller by a factor of about 2 for all fragment lengths (*Figure 2: A*), implying
262 that controlling for genome composition provides a further substantial improvement in modeling
263 the natural distance distribution. A further improvement is clearly achieved with the P-model,
264 even though, at sequence lengths below 20 residues, it produces minor inconsistencies in its total
265 residuals relative to the A-, T-, and G-models (*Figure 2: A*). We suspect that this is an artifact of using
266 local alignments (*Figure 2: A,C*) and, indeed, the effect disappears when using a global alignment as
267 distance metric over the same dataset (*Figure 2: B,D*). As for the A-, T-, and G-models, the residuals
268 of the P-model also have a wave shape, which is however qualitatively different from the shapes
269 for the less local random models, as it has only one crest at an alignment score of 13%. The crest
270 for the unexplained long-range distances is gone, which we attribute to the fact that accounting
271 for composition at the level of individual proteins has introduced the heterogeneity of natural
272 sequences into the random model. For 100mers the total residual of the P-model 0.9% (*Figure 2: A*),
273 a value that is not improved remarkably by an even greater locality: The residuals of the D-model
274 have the same wave shape as those of the P-model and a comparable amplitude, providing only
275 a minor improvement with a total residual of 0.8%. This was somewhat surprising, as it is well
276 established that many proteins are composed of disparate parts such as domains of distinct fold
277 classes, intrinsically unstructured regions or fibrous parts, that are known to be characterized by
278 different residue compositions *Dosztányi et al. (2005)*. The local composition of proteins that is
279 composed of heterogeneous parts, should thus be scrambled in the P-model and preserved in the
280 D-model. We therefore expected that the D-model would provide a more pronounced improvement
281 over the P-model.

282 **Similar results of D- and P-models are associated to the dataset**

283 We see two reasons why the total residuals of the D- and the P-models are almost identical. One
284 is a technical reason, that there is no room for fluctuation of local residue composition in our
285 bacterial dataset, as it may comprise a large number of short and single-domain proteins. In order
286 to evaluate this, we analyzed their sequence lengths and estimated the number of single- and
287 multi-domain proteins. The second is a potential qualitative characteristic of our dataset, that in
288 long proteins the local residue composition does not fluctuate remarkably. We approached this
289 possibility by assessing the fraction of proteins that comprise a mixture of structured domains and
290 residues not assigned to a domain.

291 In order to distinguish how these two reasons contribute to the comparable total residuals of the
292 D- and P-models, we included two eukaryotic datasets in the following analysis. We retrieved the
293 proteomes of *Homo Sapiens* and *Arabidopsis thaliana* from UniRef *Apweiler (2009)*, pruned them
294 according to the procedure used for our bacterial dataset of low-complexity regions and fragments
295 shorter than 100 residues. We present the total residuals of the P- and D-models for these three
296 datasets (*Figure 3: A*) and find that the total residuals of the D-models for the eukaryotic datasets is
297 2-fold smaller than those of the P-models, which contrasts with the observation in the bacterial
298 dataset.

299 To investigate the first technical reason why the total residuals of the D- and P-models are almost

300 identical in the bacterial dataset, we determined the protein length distribution. The bacterial
301 dataset has the shortest proteins with an median length of 315 residues, in the *Arabidopsis thaliana*
302 dataset the median length is 400 residues and 550 residues in the *Homo sapiens* dataset (**Figure 3:**
303 B). Consequently, there is a tendency of bacterial proteins to be shorter, and the median protein
304 length of the three datasets correlates with the ratio in the total residual between the P- and
305 D-model.

306 To estimate if the number of single and multi-domain proteins has an impact, we randomly sampled
307 proteins for each of the three datasets and used HHpred *Remmert et al. (2012)* for their domain
308 annotation over the ECOD database *Cheng et al. (2014)*, which represents the most recent and
309 comprehensive classification of domains of known structure (see Methods). We considered proteins
310 as multi-domain proteins if they had at least 2 domains assigned to them, otherwise we considered
311 them as single-domain proteins. The predicted fraction of multi-domain proteins in our bacterial
312 dataset is 30%, which is smaller in *Arabidopsis thaliana* (25%) and greater in *Homo Sapiens* (35%).
313 The fraction of multi-domains does not correlate with the differences in the total residuals between
314 bacterial and eukaryotic datasets and there is hence no indication for it to effect the observed
315 similarity between the total residuals of the D- and P-models in the bacterial dataset.

316 We investigated the second possibility, that sequences of distinct compositions are combined within
317 proteins, by assessing the fraction of proteins that comprise a mixture of structured domains and
318 residues not assigned to a domain. To that end, we obtained the coverage for each protein by
319 residues assigned to a structured domains and identified the fraction of structured regions of the
320 protein (**Figure 3:** C). For the bacterial dataset, 40% of all sampled proteins are predicted to be
321 structured over >90% of their sequence, a fraction that is smaller in *Arabidopsis thaliana* (15%) and
322 *Homo Sapiens* (13%). We interpret the finding that there are more completely structured proteins in
323 bacteria to be a reason why local composition does not fluctuate over proteins as much as in the
324 eukaryotic examples. This observation also correlates with the ratio in the total residual between
325 the P- and D-model (**Figure 3:** A).

326 We conclude that the D-model approximates the natural distance distribution better than the
327 P-model for datasets containing local fluctuation in residue composition within proteins, which
328 is more enhanced in longer proteins in general. This may be due different reasons that lead to
329 proteins with a locally heterogeneous composition, such as random domain recombination or a
330 mixture of structured with unstructured parts. In our analysis, we found that heterogeneity of local
331 residue composition within proteins is more pronounced in eukaryotic than in bacterial proteins.
332 Overall, the D-model is in any considered case at least slightly better than the P-model,

333 **Sequence bias caused by homology**

334 Having accounted for compositional effects at increasingly local level, the remaining discrepancy
335 between the distance distributions of the D-model and the natural dataset must be related to
336 the actual sequence of amino acids. This discrepancy can arise either through divergence from
337 a common ancestor (homology) or convergence as a result of structural constraints, particularly
338 secondary structure formation (analogy). In order to evaluate the relative contribution of these
339 mechanisms to the natural distances between sequence fragments we aimed to identify what pro-
340 portion of distances could be assigned confidently to either homologous or analogous relationships
341 and evaluated their contribution to the natural distance distribution.

342 The detection of homologous relationships requires advanced approaches, which are computationally
343 much more expensive than the simple sequence alignments used to determine distances in
344 sequence space. We therefore only considered a small subset of our entire dataset and relationships
345 within this subset, which could be derived computationally in a reasonable amount of time.
346 Therefore, we systematically sampled our dataset into 10 unbiased groups of 100mers, containing
347 approximately 650 sequences each, and used HHblits to generate profile Hidden Markov Models

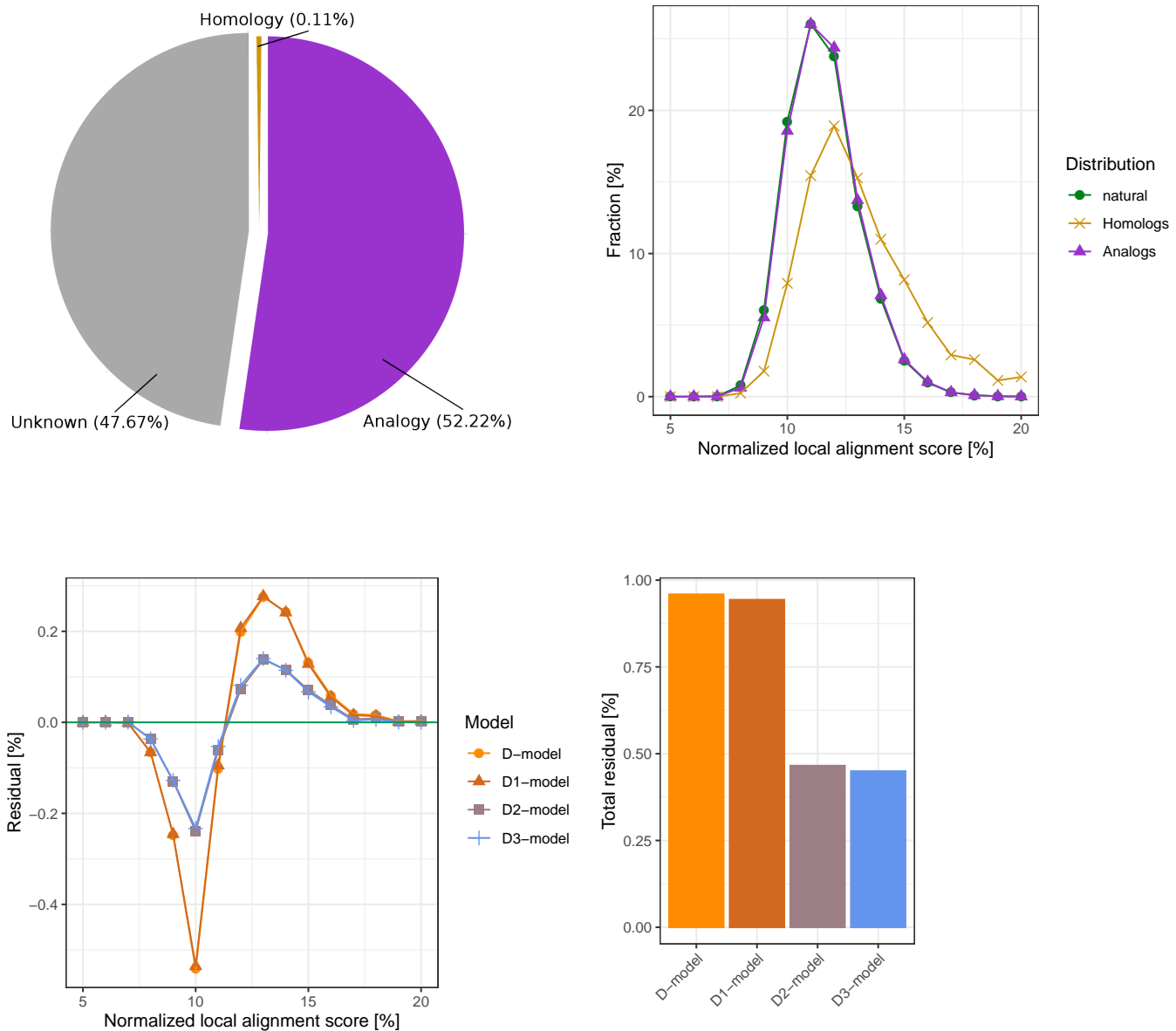


Figure 4. The contribution of homology and analogy to the global occupation of sequence space. (A) Decomposition of fragment pairs into their origins. We sampled 2 Million fragments pairs and analyzed if their relationship is confidently homologous or analogous. The fraction of analogous relationships was determined to be 52.22%, homologous relationships only 0.11% and the remaining fraction is labeled of unknown origin. Thus, the majority of relationships is generally analogous. (B) Distance distribution between homologs and analogs contrasted with the natural distance distribution. The qualitative difference between the distance distribution of analogs and that of all fragments is relatively small. Compared to this, the distance distribution of homologs displays a strong tendency towards a higher sequence identity score; it nevertheless has a major overlap with the natural distribution. (C) Residuals of the models incorporating the sequence bias of homology and analogy. We generated mixed models, that include the sequence bias of homology (D1-model), analogy (D2-model) and both (D3-model) into the D-model, which is only based on the composition of natural 100mers. The D1-model, which includes homologous sequence bias, displays almost the same residuals as the purely composition-based D-model. The residuals of the D2-model, which includes analogous sequence bias, deviate severely from that of the D-model. The D3-model yields similar results as the D2-model. (D) Total residuals of mixed models. The total residuals behave accordingly to the residuals. The D1-model has displays an only improvement in the total residual of 0.016% compared to that of the D-model. The D2-model reaches a total residual of 0.46% and is more than 2-fold more accurate than the D-model (0.96%). Adding the homology bias to the D2-model to obtain the D3-model has almost no effect.

348 (HMMs) for all individual sequences within these groups. We then derived a set of relationships by
349 aligning the retrieved HMMs from one set of sequences to those of another. This we repeated for
350 arbitrary sets of 100mers, resulting in multiple unbiased samples of relationships. The likelihood of
351 homology between two HHMs was derived using the tool HAlign and required a strict threshold of
352 minimally 90% probability (see Methods). This process identified 0.108% of pairwise relationships
353 as homologous, with a standard error of the mean (SEM) of 0.0033% (**Figure 4: A**).

354 For the remaining sequence pairs, we evaluated the likelihood of analogy by comparing their
355 HMMs to those of the ECOD database **Cheng et al. (2014)**. By virtue of containing only domains of
356 known structure, ECOD is the currently best resource for distinguishing between homology and
357 analogy in protein domains. For our analysis, we scored pairs of sequences as analogous if they
358 matched distinct X-groups in the ECOD hierarchy using the same probability cutoff of 90% as for the
359 homology assignment. In most cases, the X-level is the highest level at which homology still needs
360 to be considered as a possibility; requiring fragments to match different X-groups within this level
361 thus provided a conservative estimate of analogous relationships. We are aware that few proteins
362 of distinct X-groups may have a common ancestry, and acknowledge that using ECOD as golden
363 standard still may not be perfect in discriminating homology from analogy. This process identified
364 52.22% of pairwise relationships as analogous (**Figure 4: A**, cyan), with a SEM of 0.84%. We conclude
365 from this that the number of confident analogous pairs exceeds the number of homologous pairs
366 by more than 2 orders of magnitude. This already indicates that the influence of homology on the
367 global distance distribution in natural sequences will be dwarfed by analogy. All sequence pairs that
368 could not be confidently assigned to either group were considered to be of unknown relationship,
369 amounting to 47.6% of the total with a SEM of 0.84% (**Figure 4: A**, lightbrown).

370 Having decomposed pairwise sequence pairs into confident homologous and analogous rela-
371 tionships, we analyzed to what extent the remaining total residual (0.8%) can be explained by
372 incorporating corresponding sequence biases into our D-model. Therefore, we generated three
373 new hybrid D-models in the following way: we omitted either homologous pairs, or analogous pairs,
374 or both from our set of assigned relationships, generated a D-model for the remaining fragment
375 pairs through the same shuffling procedure as used previously, and then added back the omitted
376 pairs without shuffling. In the following we refer to the hybrid model that adds the sequence bias of
377 homologs to the domain composition as D1-model, the one that adds the sequence bias of analogs
378 as D2-model, and the one that adds both biases as D3-model.

379 The residuals of these three models in addition to that of the D-model are shown in (**Figure 4: C**).
380 Due to the reduced sampling over only 2 Million fragment pairs instead of 500 Million, the total
381 residual of the D-model deviates slightly from that of our main analysis and has a value of 0.96%
382 instead of 0.8% (**Figure 4: D**).

383 Relative to this total residual of the purely compositional D-model, the D1-model, which includes
384 homologous sequence effects, is only minimally better (0.016%) at describing the natural distance
385 distribution. We assume that two reasons are mainly responsible for this only minor improvement:
386 First, the proportion of homologous relationships is only 0.11%, giving them little leverage. Second,
387 the distance distribution of homologs (**Figure 4: B**, yellow distance distribution) differs only to a
388 small extent from the distance distribution of the D-model. This is not entirely unexpected, given
389 how difficult it is to distinguish distant homology from random fluctuation in sequence comparisons.
390 In fact, it has been recognized previously that most homologous sequences share no significant
391 similarity **Rost (1999)**.

392 In contrast, the total residual of the D2-model (0.46%), which includes analogous sequence effects,
393 is decreased about 2-fold relative to a D-model. Thus, although analogs have a distance distribution
394 that is very similar to that of the D-model (**Figure 4: B**, cyan), their leverage is 2 orders of magnitude
395 higher than that of homologs, causing these small differences to improve substantially the fit of

396 the D2-model to the natural distance distribution. This is again not entirely unexpected, as most
397 sequences in our natural dataset share the ability to form secondary structures (**Figure 3**: coverage
398 by structure), resulting in a sequence bias that is not fully captured by residue composition **Pande**
399 **et al. (1994)** **Lavelle and Pearson (2009)**. As expected from the D1-model, adding the homologous
400 sequence bias to the D2-model did not really improve its ability to approximate the spread of
401 natural sequences. We conclude that the sequence space of natural proteins is almost entirely
402 shaped by compositional effects and the remaining sequence bias is almost entirely due to analogy,
403 which we interpret to result from secondary structure formation.

404 **Conclusion**

405 In this article we have undertaken a study of natural protein sequence space, using an approach
406 built on the probability mass function of pairwise distances between sequence fragments. With
407 this approach we were able to evaluate the sequence space of fragments up to 100 residues in
408 length, substantially exceeding previous efforts and for the first time characterizing globally the
409 relative position of sequences in space. Our results show that the global compositional bias of
410 natural proteins is already sufficient to approximate the distance distribution of natural sequences
411 by 95.4% and that accounting for local compositional bias down to the level of individual 100mers
412 further improves this to 99.2%. The remaining 0.8% of unaccounted distances between natural
413 100mers are almost entirely contributed by sequence effects arising from analogous relationships,
414 leaving only a negligible contribution to homology in the global characterization of sequence space
415 occupation.

416 This surprised us, as decades of bioinformatic work have mapped out an increasingly comprehen-
417 sive description of sequence space around protein families, based on the detection of ever more
418 remote homology. We therefore expected to find that homology also has a substantial role in
419 shaping the global structure of sequence space occupation. This expectation was not borne out
420 and in retrospect this might not seem as surprising, given that even the space of single protein
421 families, can span over broad areas of sequence space. Other indicators for this can be seen for
422 example in the progressively more complex statistical methods needed to substantiate homology
423 across increasingly large evolutionary distances, the resulting difficulties to classify the detected
424 homologous relationships into a hierarchy of protein families and superfamilies, and the remaining
425 inability in many cases to judge on the homologous or analogous nature of similarities even in the
426 presence of extensive sequence and structure information **Rost (1999)**. These considerations show
427 that even at the level of protein families, many sequence relationships comprise a large random
428 element, substantially indistinguishable from random fluctuation and sequence convergence. This
429 random element not only results from our inability to detect homologs that have diverged strongly
430 due to low selective pressure, but also from the fact that in many families, a conserved core has
431 been elaborated in different ways with analogous sequences.

432 We find a much larger influence of analogous sequence biases on the global shape of naturally
433 occupied sequence space. The main common feature of proteins in our natural dataset is the ability
434 to fold, which translates into a propensity to assume secondary structure locally. We see this as the
435 main reason for the sequence bias that we observe between analogous sequences. Nevertheless,
436 the sequence biases of homology and analogy together are responsible for only 0.8% of distances
437 between natural sequences, that cannot be explained by a random model incorporating the natural
438 amino acid composition of domain-sized fragments. We conclude that natural sequences stand
439 out from randomness primarily through their biased use of the 20 amino acids. Accounting for
440 this bias at increasingly local levels is largely sufficient to model the global structure of sequence
441 space occupation. This major relevance of composition has been acknowledged as it has been
442 implemented into BLAST **Schaffer (2002)** and been demonstrated to be key for the aggregation of
443 intrinsically unstructured proteins **Vymětal et al. (2019)**.

444 There seems to be no other striking feature of the primary structure in natural protein sequences
445 and in consequence there are also no other obvious features that distinguish natural from random
446 sequences. We conclude that viable proteins could be located anywhere in the sequence space
447 defined by natural residue compositions. The main reason why the proteome of nature currently
448 only comprises some 10^{12} proteins [Andrei] and that these mainly fall into only about 10^5 families
449 *Punta et al. (2012)* is therefore not due to the limited availability of useful sequence space, but
450 rather to their evolutionary history. Thus, there is treasure everywhere.

451 **Methods**

452 **Natural data**

453 **Genome selection**

454 With the aim to achieve a reasonable coverage of deep phylogenetic branches with complete and
455 well-annotated proteomes, we selected the majority of bacterial genomes provided by UniRef
456 on 22.09.2017 *Apweiler (2009)*. Some genomes stood out as they possessed multiple replicas of
457 the same protein and were excluded, leaving 4,098 to remain. For each of the 1,307 genera we
458 randomly chose one representative for our natural data set. The genus was derived from the
459 full-length genome name via string matching.

460 We are aware of the general ambiguity of the definition of a genus *Parks et al. (2018)*. However,
461 with the genus selection we only aimed to reduce redundancy caused by some species that have
462 been sequenced many times. Lastly, we note that the bias towards bacteria that are easy to cultivate
463 prohibits a sampling of the true diversity among bacterial genomes.

464 **Genome curation**

465 Apart from redundancy at the genome level, we control for recent gene duplication events. For
466 each genome, we cluster its proteins using cd-hit (version 4.6 with 99% sequence identity and 90%
467 coverage). A representative protein sequence, as defined by cd-hit, was then selected for each
468 cluster; all other proteins were discarded.

469 **Low complexity filtering**

470 Low-complexity regions (LCRs) are a well-known feature of natural sequences, that do not occur as
471 frequently in random sequences. We first analyzed our data including LCRs and found that they
472 majorly contribute to the total residual and variance contrast between natural sequences and our
473 models (data not shown). Therefore, we pruned LCRs of our dataset using segmasker *Wootton*
474 *and Federhen (1996)* (version 2.3.0+ with the standard settings), to obtain differences between
475 natural and random sequences that are not due to this well-known feature. This pruning of LCRs
476 leads to sequences of slightly higher complexity than expected for short peptides (data not shown).
477 The pruning bias plays an insignificant role, especially for longer sequences, which are of most
478 interest in our study. Since, N-terminal methionines were sometimes included, we stripped them to
479 standardize our sequences.

480 **Sequence adjustments**

481 To simplify our analysis we changed a couple of hundred cases of uncommon amino acids to their
482 most similar proteinogenic amino acid. In order to use the exact same dataset for all sequence
483 lengths, we pruned our data set of sequences shorter than 100. Additionally, we removed the
484 invalid amino acid X by replacing it with an end-of-line-character, effectively dividing a protein
485 sequence into multiple parts. However, since some of our random models depend on shuffling
486 intact genomes or proteins, we performed this division into multiple parts after the shuffling (more
487 detail below).

488 Complete statistics and data availability

489 Taken together our dataset holds $1.2 \cdot 10^9$ valid amino acids of 1,307 genomes comprising $4.7 \cdot 10^6$
490 proteins. In the supplements we provide:

- 491 • fasta-file of original genomes
- 492 • fasta-file of adjusted genomes
- 493 • amino acid composition

494 **Fragment pair selection and random sequence models**

495 Fragment selection

496 We selected random fragments such that each character (amino acids and end-of-line-character) in
497 the dataset had the same probability of being chosen and that the same fragment pair would never
498 be chosen twice. We ensured this by implementing two linear congruential generator *Press et al.*
499 (2007) to enumerate all possible pairs of fragments. In detail, one linear congruential generator
500 was used for each member of the pair with multiplier $a = 1$ and moduli $m_1 = 223$ and $m_2 = 34,211$,
501 where both moduli are prime numbers relative to the total number of characters 1,168,754,000.
502 Depending on the starting points of the two generators, a different subset of index pairs can be
503 selected. This enabled us to calculate disjunctive fragment pairs in parallel. We selected $5 \cdot 10^8$ valid
504 pairs of fragments to accurately estimate the distance distributions and rejected fragments that
505 straddled protein boundaries or invalid regions, indicated by the end-of-line-character.

506 Model based on overall amino acid composition

507 The most standard random sequence model is based on the underlying amino acid composition of
508 a given dataset. We obtained randomized data for this A-model by randomly shuffling all amino
509 acids of the natural data. Thereby, protein length is maintained and the number of amino acids
510 stays exactly the same. As all our random models are based on random permutations, we used
511 the Mersenne Twister algorithm mt19337 of the C++ 14 std library with the standard seed value of
512 19650218. This algorithm is considered one of the best pseudo-random number generators and in
513 a test with a smaller dataset we found that our results did not depend on the type or seeding of the
514 random number generator.

515 For the E-model, we proceeded the same way as for the A-model. The only difference is that we
516 replaced the natural dataset, by writing over all valid amino acids with the 20 possible amino acids
517 in lexicographical order. When reaching the character Y for tryptophan, we started over with A for
518 alanine. The distance distribution of the E-model deviates severely from that of the natural dataset.
519 Here, we provide the distance distribution and its residuals in Figure X as we are referring in the
520 results section to these values.

521 Models based on the amino acid composition of genomes or proteins

522 To account for genome or protein composition, we shuffled amino acids within the context of
523 genomes or proteins. For the G-model, we shuffled amino acids within each of the 1,307 genomes.
524 For the P-model we shuffled amino acids within each protein. We used one instance for genome and
525 protein composition bias and stored them to generate the distance distribution for the correspond-
526 ing models. Multiple instances for genome and protein composition bias were found converge
527 to the same results after a large enough sampling of fragment pairs. After shuffling, we divided
528 proteins containing the invalid amino acid X by replacing it with an end-of-line-character.

529 Model based on the amino acid composition of domain-sized fragments

530 For the D-model, we randomly shuffled natural fragments of length 100. In contrast to the previous
531 random models, generating a single randomly shuffled dataset is not computationally feasible since
532 storing an instance of all shuffled 100mers would increase the data size approximately 100-fold. We
533 therefore shuffled 100mers on the fly during the calculation of the distance distribution. In detail,

534 we select pairs of natural fragments as described above and consider the target fragment of length
535 N to be located in the middle of the domain. If the domain straddles any protein boundaries, we
536 adjust the domain boundaries such that the domain fits into the protein boundaries by shifting.
537 Note that because of this adjustment, the selection probability of amino acids into domains is not
538 uniform but the selection probability of amino acids in fragments is. The alternative would be
539 a rejection procedure, where we would reject fragments that are so close to protein boundaries
540 that the domain of length of 100 would not fit. The downside of such a rejection procedure is
541 that fragments close to protein boundaries are not selected and hence the selection probability
542 of fragments is not uniform anymore, which differs from the selection of natural fragments or
543 fragments for the A-, G-, and P-models. The D1, D2, and D3-models, which incorporate the sequence
544 bias of homologous and analogous fragments are presented further down.

545 **Pairwise distances as descriptor for sequence space occupation**

546 Distance metric

We define the distance between two fragments of the same length N as the rounded score from a Smith-Waterman alignment s . An amino acid match is scored with 1, a mismatch with 0, gap opening penalty is equal to 3 and gap extension penalty is 0.1, which are the same parameters as used in *Rost (1999)*. Due to gaps, scores can rank between 0 and N in 0.1 steps; to obtain integer distances, we round scores to the closest integer number. Distances exactly between two integers (such as 1.5) are assigned to the smaller one. This distance metric thereby reflects the number of dimensions in sequence space (positions in sequence), which differ between two fragments, while allowing for gaps and insertions. To compare the score p across different fragment lengths N , we transform it into the alignment score s , scaling between 0-100%, as follows:

$$s = \frac{p}{N}$$

547 In some cases, we use the rounded score from a Needleman-Wunsch alignment with identical
548 scoring parameters to illustrate differences to the Smith-Waterman approach. We also diversified
549 gap penalties, leading to comparable results (data not shown). For all alignments, we used the
550 SeqAn C++ library, version 2.4 *Rahn et al. (2018)*, which enables many sequence comparisons in
551 parallel.

552 Comparing distance distributions

The residual corresponds to the variational distance at each possible sequence identity between two distance distributions. We use it to demonstrate the qualitative difference between the distance distribution of a random model and that natural sequences. Denoting the residual by r , the random model distance distribution by D_{rand} and the natural one by D_{nat} , we have:

$$r(s) = D_{nat}(s) - D_{rand}(s)$$

553 where s sequence identity. Regions of sequence identity s where the residual $r(s)$ exceeds zero
554 therefore indicate a higher frequency of these sequence identity scores in natural fragments.

To summarize the difference between natural and random model distance distributions in a single metric, we sum the absolute residuals over all sequence identities and normalize it to a range between 0 and 100%:

$$R = \sum_{0 \leq s \leq N} \frac{|r(s)|}{2}$$

555 We call R the total residual, which is variously called the variational distance, total variation distance
556 or Kolmogorov distance *Deza and Deza (2014)*.

557 **Decomposition into homologous and analogous relationships**

558 Homology

559 We derived the fraction of sequence pairs that are confidently homologous using the tools of
560 HH-suite (version number 3.0.3) *Remmert et al. (2012)*. To derive this fraction, we systematically
561 sampled our dataset and extracted 10 sets of natural 100mers that are equally distributed over our
562 dataset, each containing approximately 650 fragments. With HHblits, we generated HMMs with the
563 standard settings for each of these fragments with two iterations, using unclust30 as underlying
564 database (version August 2018).

565 Then, we pairwise aligned the generated HMMs with HHalign, in order to estimate whether two
566 fragments are homologous. We did this by aligning all fragments in one set to all of those in another
567 set, resulting in 90 possible directed combinations of which we chose 10 as representative sets of
568 pairwise relationships. Each set of fragments was considered twice in this comparison, once as the
569 set of query sequences and once as the set of target sequences in the alignment. This resulted in 2
570 Million pairwise fragment comparisons divided into 10 disjunctive sets. Fragments were taken to be
571 homologous, if HHalign predicted them to be homologous with a probability above 90%. In total
572 0.11% of the fragment pairs were found to be homologous; the standard error of the mean (SEM)
573 derived from the 10 sets of 0.006%.

574 Analogy

575 We derived the fraction of sequence pairs that are confidently analogous using a similar procedure
576 as used for the homology detection. We first assigned structured domains to each 100mer. We
577 then assumed a pair of 100mers to be of analogous origin, if the two 100mers matched only distinct
578 domains that are confidently not related to each other.

579 For the assignment of structured domains, we used the ECOD classification *Cheng et al. (2014)*,
580 which is the currently best resource for distinguishing between homology and analogy in protein
581 domains. The HMMs of each 100mer (same as in the homology detection) were thereby compared
582 against all ECOD entries (retrieved on 9.4.2019) with HHsearch. We used HHsearch with the standard
583 parameter and assigned the best-scoring non-overlapping hits with a probability above 90% to the
584 corresponding fragment. Of all 100mers 70% could be assigned to a single domain and less than
585 1% to multiple domains, of which we considered each. Other 100mers were not assigned to any
586 domain, which we directly excluded to be analogous to any other sequence, since we are uncertain
587 about their origin.

588 For the assignment of analogous relationships, we considered only pairs of 100mers that were
589 assigned to at least one domain. If their domains matched only distinct X-groups in the ECOD
590 hierarchy, the pair was assumed to have an analogous relationship. The X-group is the highest level
591 at which homology still needs to be considered as a possibility. All pairs of fragments that were
592 assigned to domains of only distinct X-levels were considered to be confidently analogous.

593 With this procedure 52.22% of the fragment pairs were found to be analogous; the standard error
594 of the mean derived from the 10 sets is 0.84%. The remaining 47.6% of the fragment pairs is of
595 unknown relationship.

596 **Mixed models containing sequence bias of homology or analogy**

597 In order to estimate the influence of homology and analogy to the natural distance distribution,
598 we generated mixed models that account for their sequence bias. For the D1-model we
599 included the homologous sequence bias by deriving the distances between the 1,309 confidently
600 homologous fragment pairs. We applied the D-model to the remaining fragment pairs and shuffled
601 the fragments of the corresponding pairs with the Unix command shuf followed by deriving their
602 distance. All distances combined resulted into the distance distribution of the D1-model. We
603 proceeded the same way for the D2-model by including the distances of unshuffled fragments that

604 are confidently analogous, and distances of the remaining pairs after shuffling the residues within
605 each fragment. Thereby, the sequence bias between analogous fragments was preserved while for
606 other fragment pairs only their composition is accounted for. For the D3-model we included both
607 sequence bias of homologous and analogous natural fragments.

References

- 608
609 **Alva V**, Remmert M, Biegert A, Lupas AN, Söding J. A galaxy of folds. *Protein Science*. 2009 jan;
610 19(1). <http://www.ncbi.nlm.nih.gov/pubmed/19937658>[http://www.pubmedcentral.nih.gov/articlerender.fcgi?](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC2817847)
611 <http://doi.wiley.com/10.1002/pro.297>, doi: 10.1002/pro.297.
- 612 **Alva V**, Söding J, Lupas AN. A vocabulary of ancient peptides at the origin of folded proteins. *eLife*. 2015; 4. doi:
613 10.7554/eLife.09410.001.
- 614 **Apweiler R**. The universal protein resource (UniProt) in 2010. *Nucleic Acids Research*. 2009 jan; 38:D190–
615 5. <http://www.ncbi.nlm.nih.gov/pubmed/18045787>[http://www.pubmedcentral.nih.gov/articlerender.fcgi?](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC2238893)
616 <http://doi.wiley.com/10.1093/nar/gkp846>, doi: 10.1093/nar/gkp846.
- 617 **Bershtein S**, Serohijos AW, Shakhnovich EI. Bridging the physical scales in evolutionary biology: from protein
618 sequence space to fitness of organisms and populations. *Curr Opin Struct Biol*. 2017; 42:31–40. <http://dx.doi.org/10.1016/j.sbi.2016.10.013>, doi: 10.1016/j.sbi.2016.10.013.
- 619
- 620 **Cedano J**, Aloy P, Perez-Pons JA, Querol E. Relation between amino acid composition and cellular location of
621 proteins. *Journal of Molecular Biology*. 1997 feb; 266(3):594–600. [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S0022283696908049)
622 [article/pii/S0022283696908049](https://www.sciencedirect.com/science/article/pii/S0022283696908049), doi: 10.1006/jmbi.1996.0804.
- 623 **Cheng H**, Schaeffer RD, Liao Y, Kinch LN, Pei J, Shi S, Kim BH, Grishin NV. ECOD: An Evolutionary Classification
624 of Protein Domains. *PLoS Computational Biology*. 2014 dec; 10(12):e1003926. [https://dx.plos.org/10.1371/](https://dx.plos.org/10.1371/journal.pcbi.1003926)
625 [journal.pcbi.1003926](https://dx.plos.org/10.1371/journal.pcbi.1003926), doi: 10.1371/journal.pcbi.1003926.
- 626 **Chevalier A**, Silva DA, Rocklin GJ, Hicks DR, Vergara R, Murapa P, Bernard SM, Zhang L, Lam KH, Yao G, Bahl
627 CD, Miyashita SI, Goreshnik I, Fuller JT, Koday MT, Jenkins CM, Colvin T, Carter L, Bohn A, Bryan CM, et al.
628 Massively parallel de novo protein design for targeted therapeutics. *Nature*. 2017 oct; 550(7674):74–79.
629 <http://www.nature.com/articles/nature23912>, doi: 10.1038/nature23912.
- 630 **Chou KC**. Prediction of protein cellular attributes using pseudo amino acid composition. *PROTEINS Struct Funct*
631 *Genet* (Erratum *ibid*, 2001, Vol44, 60). 2001; 43(3):246–255. doi: 10.1017/CBO9781107415324.004.
- 632 **Deza MM**, Deza E. *Encyclopedia of Distances*. Springer Berlin Heidelberg; 2014. [https://books.google.de/books?](https://books.google.de/books?id=q_7FBAAAQBAJ)
633 [id=q_7FBAAAQBAJ](https://books.google.de/books?id=q_7FBAAAQBAJ).
- 634 **Diggle PJ**. *Statistical Analysis of Spatial and Spatio-Temporal Point Patterns*. CRC press; 2014.
- 635 **Dosztányi Z**, Csizsók V, Tompa P, Simon I. The pairwise energy content estimated from amino acid composition
636 discriminates between folded and intrinsically unstructured proteins. *Journal of Molecular Biology*. 2005;
637 347(4):827–839. doi: 10.1016/j.jmb.2005.01.071.
- 638 **Duranton G**, Overman HG. Testing for Localization Using Micro-Geographic Data. *The Review of Economic*
639 *Studies*. 2005 oct; 72(4):1077–1106. [https://academic.oup.com/restud/article-lookup/doi/10.1111/0034-6527.](https://academic.oup.com/restud/article-lookup/doi/10.1111/0034-6527.00362)
640 [00362](https://academic.oup.com/restud/article-lookup/doi/10.1111/0034-6527.00362), doi: 10.1111/0034-6527.00362.
- 641 **Fukuchi S**, Nishikawa K. Protein surface amino acid compositions distinctively differ between thermophilic and
642 mesophilic bacteria. *J Mol Biol*. 2001 jun; 309(4):835–843. [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0022283601947187)
643 [S0022283601947187](https://www.sciencedirect.com/science/article/pii/S0022283601947187)?via{&}3Dihub, doi: 10.1006/jmbi.2001.4718.
- 644 **Fukuchi S**, Yoshimune K, Wakayama M, Moriguchi M, Nishikawa K. Unique amino acid composition of proteins
645 in halophilic bacteria. *J Mol Biol*. 2003; 327(2):347–357. doi: 10.1016/S0022-2836(03)00150-5.
- 646 **Harms MJ**, Thornton JW. Historical contingency and its biophysical basis in glucocorticoid receptor evolu-
647 tion. *Nature*. 2014 aug; 512(7513):203–207. <http://www.ncbi.nlm.nih.gov/pubmed/24930765>[http://www.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4447330)
648 [pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4447330](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4447330), doi: 10.1038/nature13410.
- 649 **Huang PS**, Boyken SE, Baker D. The coming of age of de novo protein design. *Nature*. 2016; 537(7620):320–
650 7. <http://www.nature.com/doi/10.1038/nature19946>?%5Cn[http://www.ncbi.nlm.nih.gov/pubmed/](http://www.ncbi.nlm.nih.gov/pubmed/27629638)
651 [27629638](http://www.ncbi.nlm.nih.gov/pubmed/27629638), doi: 10.1038/nature19946.
- 652 **Larson SM**, England JL, Desjarlais JR, Pande VS. Thoroughly sampling sequence space: large-scale protein design
653 of structural ensembles. *Protein Sci*. 2002; 11(12):2804–13. [http://www.pubmedcentral.nih.gov/articlerender.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2373757)
654 [fcgi?artid=2373757](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2373757){&}&tool=pmcentrez{&}&rendertype=abstract, doi: 10.1110/ps.0203902.
- 655 **Lavelle DT**, Pearson WR. Globally, unrelated protein sequences appear random. *Bioinformatics*. 2009; 26(3):310–
656 318. doi: 10.1093/bioinformatics/btp660.

- 657 **de Lucrezia D**, Slanzi D, Poli I, Polticelli F, Minervini G. Do natural proteins differ from random sequences polypep-
658 tides? natural vs. random proteins classification using an evolutionary neural network. *PLoS ONE*. 2012 may;
659 7(5):e36634. <https://dx.plos.org/10.1371/journal.pone.0036634>, doi: 10.1371/journal.pone.0036634.
- 660 **Luigi Luisi P**. Contingency and determinism. *Philos Trans R Soc A Math Phys Eng Sci*. 2003; 361(1807):1141–1147.
661 <http://rsta.royalsocietypublishing.org/cgi/doi/10.1098/rsta.2003.1189>, doi: 10.1098/rsta.2003.1189.
- 662 **Lupas AN**, Dyke MV, Stock J. Predicting Coiled Coils from Protein Sequences. *Science*. 1991; 252:1162–1164.
- 663 **Munteanu CR**, González-Díaz H, Borges F, de Magalhães AL. Natural/random protein classification mod-
664 els based on star network topological indices. *Journal of Theoretical Biology*. 2008; 254(4):775–783. doi:
665 10.1016/j.jtbi.2008.07.018.
- 666 **Pande VS**, Grosberg aY, Tanaka T. Nonrandomness in protein sequences: evidence for a physically driven stage
667 of evolution? *Proc Natl Acad Sci U S A*. 1994; 91(26):12972–12975. doi: 10.1073/pnas.91.26.12972.
- 668 **Parks DH**, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil PA, Hugenholtz P. A standardized
669 bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat Biotechnol*. 2018;
670 36(10):996. doi: 10.1038/nbt.4229.
- 671 **Poznański J**, Topiński J, Muszewska A, Dębski KJ, Hoffman-Sommer M, Pawłowski K, Grynberg M. Global
672 pentapeptide statistics are far away from expected distributions. *Sci Rep*. 2018 dec; 8(1):15178. <http://www.nature.com/articles/s41598-018-33433-8>, doi: 10.1038/s41598-018-33433-8.
- 673
- 674 **Press WH**, Teukolsky SA, Vetterling WT, Flannery BP. *Numerical Recipes 3rd Edition: The Art of Scientific*
675 *Computing*. 3 ed. New York, NY, USA: Cambridge University Press; 2007.
- 676 **Punta M**, Coggill PC, Eberhardt RY, Mistry J, Tate J, Boursnell C, Pang N, Forslund K, Ceric G, Clements J, Heger
677 A, Holm L, Sonnhammer ELL, Eddy SR, Bateman A, Finn RD. The Pfam protein families database. *Nucleic*
678 *Acids Research*. 2012 jan; 40(D1):D290–D301. [https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkr1065)
679 [gkr1065](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkr1065), doi: 10.1093/nar/gkr1065.
- 680 **Rahn R**, Budach S, Costanza P, Ehrhardt M, Hancox J, Reinert K. Generic accelerated sequence alignment in
681 SeqAn using vectorization and multi-threading. *Bioinformatics*. 2018; p. 1–9. [https://academic.oup.com/](https://academic.oup.com/bioinformatics/advance-article-abstract/doi/10.1093/bioinformatics/bty380/4992147)
682 [bioinformatics/advance-article-abstract/doi/10.1093/bioinformatics/bty380/4992147](https://academic.oup.com/bioinformatics/advance-article-abstract/doi/10.1093/bioinformatics/bty380/4992147), doi: 10.1093/bioinfor-
683 matics/bty380.
- 684 **Remmert M**, Biegert A, Hauser A, Söding J. HHblits: lightning-fast iterative protein sequence searching by
685 HMM-HMM alignment. *Nature Methods*. 2012 feb; 9(2):173–175. <http://www.nature.com/articles/nmeth.1818>,
686 doi: 10.1038/nmeth.1818.
- 687 **Rost B**. Twilight zone of protein sequence alignments. *Protein Engineering Design and Selection*. 1999;
688 12(2):85–94. doi: 10.1093/protein/12.2.85.
- 689 **Schaffer AA**. Improving the accuracy of PSI-BLAST protein database searches with composition-based statistics
690 and other refinements. *Nucleic Acids Research*. 2002; 29(14):2994–3005. doi: 10.1093/nar/29.14.2994.
- 691 **Schneider R**, de Daruvar A, Sander C. The HSSP database of protein structure-sequence alignments. *Nucleic*
692 *Acids Res*. 1997 jan; 25(1):226–230. <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/25.1.226>,
693 doi: 10.1093/nar/25.1.226.
- 694 **Shah P**, McCandlish DM, Plotkin JB. Contingency and entrenchment in protein evolution under purifying
695 selection. . 2015; 2015. <http://arxiv.org/abs/1404.4005>{%}0Ahttp://dx.doi.org/10.1073/pnas.1412933112, doi:
696 10.1073/pnas.1412933112.
- 697 **Starr TN**, Picton LK, Thornton JW. Alternative evolutionary histories in the sequence space of an ancient protein.
698 *Nat Publ Gr*. 2017; 549. <https://www.nature.com/articles/nature23902.pdf>, doi: 10.1038/nature23902.
- 699 **Starr TN**, Thornton JW. Epistasis in protein evolution. . 2016; 25(7):1204–1218. [http://www.ncbi.nlm.](http://www.ncbi.nlm.nih.gov/pubmed/26833806)
700 [nih.gov/pubmed/26833806](http://www.ncbi.nlm.nih.gov/pubmed/26833806)<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4918427>, doi:
701 10.1002/pro.2897.
- 702 **Urlinger S**, Baron U, Thellmann M, Hasan MT, Bujard H, Hillen W. Exploring the sequence space for tetracycline-
703 dependent transcriptional activators: Novel mutations yield expanded range and sensitivity. *Proc Natl Acad*
704 *Sci*. 2000 jul; 97(14):7963–7968. doi: 10.1073/PNAS.130192197.

- 705 **Vallee BL**, Auld DS. Zinc Coordination, Function, and Structure of Zinc Enzymes and Other Proteins. *Biochemistry*.
706 1990; 29(24):5647–5659. doi: 10.1021/bi00476a001.
- 707 **Vymětal J**, Vondrášek J, Hloučová K. Sequence Versus Composition: What Prescribes IDP Biophysical Proper-
708 ties? *Entropy*. 2019 jul; 21(7):654. <https://www.mdpi.com/1099-4300/21/7/654>, doi: 10.3390/e21070654.
- 709 **Wheelan SJ**, Marchler-Bauer A, Bryant SH. Domain size distributions can predict domain boundaries. *Bioinform-*
710 *atics*. 2000; 16(7):613–618. doi: 10.1093/bioinformatics/16.7.613.
- 711 **Woolfson DN**, Bartlett GJ, Burton AJ, Heal JW, Niitsu A, Thomson AR, Wood CW. De novo protein design: How
712 do we expand into the universe of possible protein structures? . 2015; 33:16–26. [http://dx.doi.org/10.1016/j.](http://dx.doi.org/10.1016/j.sbi.2015.05.0090959-440X/{#})
713 [sbi.2015.05.0090959-440X/{#}](http://dx.doi.org/10.1016/j.sbi.2015.05.009), doi: 10.1016/j.sbi.2015.05.009.
- 714 **Wootton JC**, Federhen S. Analysis of compositionally biased regions in sequence databases. *Methods in*
715 *Enzymology*. 1996 jan; 266:554–571. <https://www.sciencedirect.com/science/article/pii/S0076687996660352>,
716 doi: 10.1016/S0076-6879(96)66035-2.
- 717 **Wüthrich K**. *NMR of Proteins and Nucleic Acids*; 1986.