

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 between fragments within each dataset and
17 compare the resulting distance distributions between sets. Even for 100mers, 95.4% of all
18 distances between natural fragments are in accordance with those of a random model
19 incorporating the natural residue composition. Hence, natural sequences are distributed almost
20 randomly in global sequence space. When further accounting for the specific residue composition
21 of domain-sized fragments, 99.2% of all distances between natural fragments can be modeled.
22 Local residue composition, which might reflect biophysical constraints on protein structure, is thus
23 the predominant feature characterizing distances between natural sequences globally, whereas
24 homologous effects are only barely 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

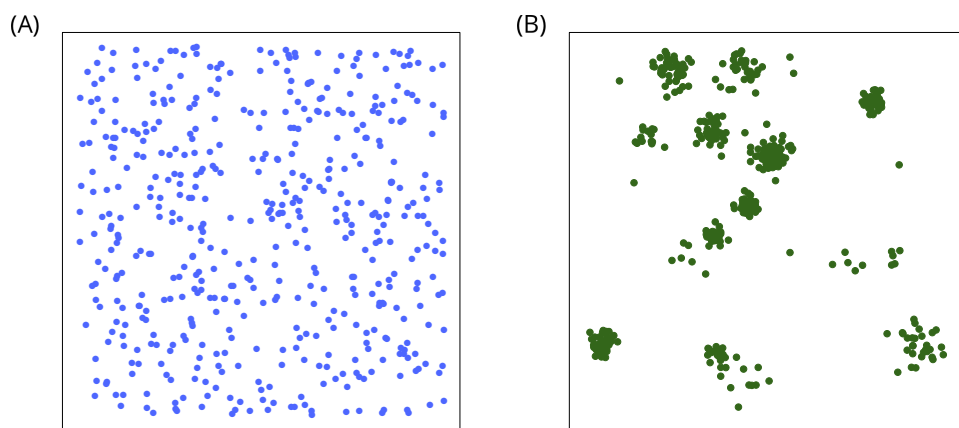


Figure 1. Sketches of sequence space occupation. (A) A random distribution has no structure. Purely random sequences are thus distributed homogeneously over the possible sequence space. With a blue circle, we represent a sequence in an abstract representation of space. (B) Natural sequences are known to frequently arise by replication and diversification. Recently duplicated sequences that have not diverged beyond obvious recognition form cluster in local areas of sequence space, indicated by green circles. Through significant similarity among multiple sequences, even very distant relatives can be assigned to a cluster of homologous sequences. This concept of sequence clusters, like islands in a gigantic ocean of possibilities, is commonly used.

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, relevant sequence lengths still mostly range above 80 residues. At a complexity of 20^{80} , it
45 is clear that this sequence space cannot be analyzed by an enumeration approach.

46 Although the sequence space of domain-sized fragments appears intractable due to its size, we have
47 nevertheless developed expectations about its occupation by natural sequences through decades
48 of bioinformatic research. This is because most proteins have arisen by descent and differentiation
49 from a set of domain prototypes, and can thus be classified into a hierarchy of domain families
50 and superfamilies. This points to the fact that the sequence space around domains is substantially
51 populated by their homologs, resulting in an image of local islands of natural sequences within a
52 global sea of virtual, unrealized possibilities (**Figure 1**). The extent to which this image is adequate
53 to describe the global sequence distribution is however unclear.

54 A first step to extend from local sequence islands to a more global view has been taken with
55 searches for variants close to existing proteins *Bershtein et al. (2017)* *Starr et al. (2017)* *Harms and*
56 *Thornton (2014)* *Urlinger et al. (2000)*. By testing exhaustively all mutations at certain sites, these
57 studies bypass intermediate mutants that would not have been viable in evolution. Contrasting the
58 abundance of possible functional variants to the small number of natural sequences demonstrates
59 how sparsely nature has explored sequence space, even locally. The high energy barriers, epistatic
60 effects, and functional dependencies prevent the establishment of random mutations and seem to
61 entrench already existing and functional forms *Starr and Thornton (2016)* *Shah et al. (2015)*.

62 Modern techniques of protein design allow to reach out further into the global sequence space
63 to find possible exemplars in unknown territory *Huang et al. (2016)* *Woolfson et al. (2015)*. Scaling
64 these scans up to the currently highest practicable level for a given structure or function has
65 uncovered viable solutions far from existing proteins *Stiffler et al. (2019)* *Chevalier et al. (2017)*
66 *Larson et al. (2002)*, showing that sequence similarity to existing proteins is not required for
67 functionality. This leads to the hypothesis that the usable part of sequence space is mostly randomly
68 structured, which has been proposed for unrelated natural sequences before *Lavelle and Pearson*
69 *(2009)*.

70 Apart from the seemingly random global structure, there are nevertheless biophysical requirements
71 for all usable protein sequences, natural as well as designed, such as foldability, hydrophobic
72 core formation and solubility. This indicates that these proteins may share some convergent
73 features, which would restrict a random drift away into unstructured space. Natural sequence
74 space could thus be characterized globally by sequences with the potential to fold, i.e. by convergent
75 features.

76 In this paper, we analyze the global structure of natural sequence space, aiming to identify general
77 features that characterize natural sequences and to evaluate the relative contributions of conver-
78 gence and homology to this space. We do this by contrasting natural data with a variety of random
79 models, in order to extract sequence features arising from different natural mechanisms.

80 **Results and Discussion**

81 **Natural sequence data and random sequence models**

82 Choice of a natural dataset

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

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

97 How random is random?

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

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

111 Having accounted for compositional effects resulting from environment, metabolism, and the need
112 to form a hydrophobic core, the remaining differences between natural and random sequences
113 must be attributed to sequence effects, due either to *divergence* from a common ancestor *Alva*
114 *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.

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

116 Sequence space has frequently been analyzed with a direct approach based on the exhaustive
117 enumeration of natural kmers, and the comparison of their frequencies to those derived from
118 a random model *Poznański et al. (2018)* *Lavelle and Pearson (2009)*. This approach is restricted
119 to kmers of length 5 or smaller, due to sequence space complexity and the data sparsity caused
120 thereby. It also does not represent the relative position among kmers within the global sequence
121 space.

122 We use an indirect approach to circumvent these problems. Our approach is built on the probability
123 mass function of pairwise distances between sequences of the same length, in the following referred
124 to as *distance distribution*. A distance distribution illustrates how often sequences are positioned
125 at a certain distance to each other and we use it to study the way sequences are spread across
126 the possible space. We built distance distributions for the natural dataset and for each dataset
127 of random sequences derived from specific models. By using lengths of up to 100 residues, our
128 sequences thus reach domain size *Wheeler et al. (2000)*.

129 As a metric for distance, we use the *normalized local alignment score* of a Smith-Waterman alignment,
130 since this metric is commonly used to capture similarities between natural sequences *Rost (1999)*
131 *Schneider et al. (1997)*. We note that the choice of distance metric is not of great relevance for the
132 main implications of our study; relative to each other, the distance distributions of the random
133 models deviate similarly from that of natural sequences irrespective of the chosen metric, as
134 outlined in the following Results sections. Details on the derivation of distance distributions and
135 the used distance metrics are provided in the Methods section. In this context, it is important
136 to note that our method differs from common approaches, as it only considers the pairwise
137 similarity between two sequences and thus their actual distance in sequence space. In contrast,
138 most bioinformatic methods that compare sequences to each other scale distances according to

139 their statistical significance and in many cases iterate comparisons in order to extract patterns
140 of conserved residues, as indicators of homologous relationships. These approaches result in
141 distances that reflect evolutionary relationships, visualized as islands of higher density in sequence
142 cluster maps *Alva et al. (2009)*.

143 Studying the layout of space through pairwise distances is common in other fields, such as protein
144 structure determination *Wüthrich (1986)*, spatial statistics *Diggle (2014)* and economics *Duranton*
145 *and Overman (2005)*, but has not, to our knowledge, been applied to investigate protein sequence
146 space. Such distance-based approaches do not preserve information about specific positions of
147 data points in space, but rather characterize their global distribution, which includes *global clustering*
148 *and dispersion*. A corollary of this is that distinct datasets become comparable through their distance
149 distribution, even if they do not share any specific data points.

150 **Comparing distance distributions**

151 For the comparison of the natural to a random distance distribution, we first subtract the fraction
152 of distances observed in the random dataset from that observed in the natural dataset for each
153 alignment score. We refer to this difference as the *residual*. Over all alignment scores, residuals sum
154 up to zero and may have values that are either positive (more natural distances) or negative (more
155 random distances). In order to obtain an overall measure of how different two distance distributions
156 are, we derive the *total residual*, which is the variational distance between two probability mass
157 functions. More precisely, the total residual is the sum over the absolute residuals, normalized to a
158 range between 0% and 100%.

159 If the two distance distributions are completely non-overlapping, the total residual assumes the
160 maximal value of 100%, indicating that no distance between natural fragments can be modeled
161 with the underlying random sequences. If they are identical, the total residual assumes a value of
162 0%, indicating that 100% of all distances in the natural distribution have a corresponding distance
163 in the random distribution. Thereby, the total residual represents the fraction of natural distances
164 that are not accounted for by the distance distribution of a random model.

165 **Global amino acid composition (A-model)**

166 We start our analysis by assessing to what extent the global amino acid composition, as captured in
167 the A-model, can describe natural sequences. We compare the distance distributions of the two
168 datasets for fragment lengths up to 100 residues, in increments of 10. At all fragment lengths, the
169 results are closely comparable. We show the results for 100mers as representative for domain-sized
170 sequences in *Figure 2* and provide the others in the supplementary figures.

171 The distance distributions of natural and A-model data overlap extensively (*Figure 2: A*). Both
172 are uni-modal with a peak at a low alignment score of 11%. Their minor differences only become
173 apparent, when their residuals are considered (*Figure 2: B*). These take the shape of a wave, with two
174 crests at alignment scores of 9% and 15% (reflecting an over-representation of the corresponding
175 distances in the natural dataset), and a trough at 11% (reflecting an under-representation). The
176 over-representation of distances both longer and shorter than expected from the random model,
177 suggests that natural sequences are less homogeneously distributed in space. We rationalize this
178 effect with the observation that natural sequences are enriched in certain parts of sequence space,
179 leading to an increase in shorter distances. This may occur both in regions with rare amino acids
180 (such as Cys, Trp and His in small proteins dominated by zinc-coordination and disulfide bonds
181 *Vallee and Auld (1990)*) and in regions with abundant amino acids (such as Leu, Ala and Glu in
182 all-alpha proteins, most extremely in coiled coils *Lupas et al. (1991)*). The compositional differences
183 in these enriched regions mean that their distance in sequence space will be larger than expected
184 from the A-model, and thus lead to a complementary increase in longer distances. Since residuals
185 add up to zero, the number of intermediate distances is correspondingly decreased.

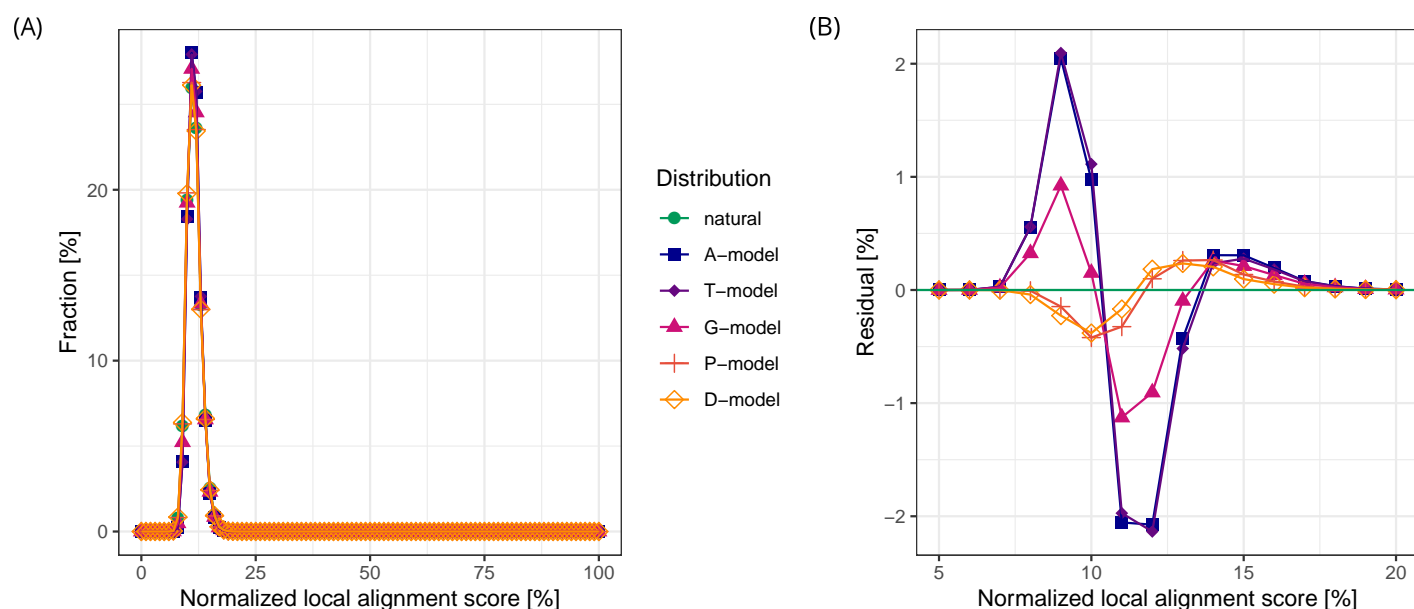


Figure 2. Comparing the sequence space occupation of random protein sequence models and natural sequence data. (A) Distance distributions are a descriptor of sequence space occupation. The distance between sequence fragments of the same length, defined as the sequence identity score obtained from a Smith-Waterman 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.

186 We note, however, that this discrepancy between natural sequences and the A-model is not very
187 pronounced, as the total residual has a value of only 4.6% for 100mers (*Figure 3: A*). It is even less
188 pronounced at smaller fragment lengths, reaching 0.4% for 10mers. We conclude that the A-model
189 becomes less accurate in describing the sequence space occupation of natural sequences at lengths
190 that are biologically relevant, but that it already achieves considerably higher accuracy than the
191 completely random model (E-model), which has a total residual of 30.4% for 100mers (data shown
192 in Methods).

193 We evaluated whether adding sequence information to the unified compositional bias of the A-
194 model could further improve it. Since nature favors certain amino acid combinations as neighboring
195 residues, a model that reflects the natural dipeptide frequency (T-model), has been proposed to
196 represent natural sequences better than the A-model *Lavelle and Pearson (2009)*. We implemented
197 the T-model by extracting the dipeptide frequencies from our natural dataset and using them to
198 generate random sequences with a Markov Chain Model. For all fragment lengths, we derived
199 the distance distribution of the T-model (*Figure 2: A*), its residuals (*Figure 2: B*) and the total
200 residual (*Figure 3: A*). By all these measures the T- and the A-model yielded essentially identical
201 results in modeling the natural distance distribution. This outcome was somewhat surprising, as
202 the addition of dipeptide frequencies to the A-model did produce a measurable improvement in
203 the enumeration study of 5mers *Lavelle and Pearson (2009)*. This may be due to the different
204 methodology in that study, which collated exact 5mer frequencies, corresponding to a position-wise
205 Hamming distance of zero, and thus being close to a global, not to a local alignment as used in our
206 study. In fact, when using the Hamming distance as metric, the T-model achieves a slightly better
207 accuracy over the A-model for sequences of 50 or less residues (*Figure 3: D*). From the results we
208 obtained with the A- and T-models, we conclude that global measures of composition and sequence
209 bias already approximate natural sequences fairly accurately, but that this accuracy decreases with
210 sequence length. Especially for longer fragments, we expect further improvement by including local
211 compositional biases as outlined in the previous section.

212 **Context-specific composition**

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

221 Since proteins are generally composed of domains, which are usually autonomous in structure and
222 also often in function, the next level of locality would be achieved by accounting for the compo-
223 sitional biases of individual domains. Producing such a D-model is however not straightforward,
224 as determining domain boundaries for proteins of unknown structure is fraught with errors and
225 many residues in our dataset cannot be assigned to a domain family. As a proxy for domains we
226 therefore derived all possible fragments of length 100 from our natural dataset and generated the
227 D-model by shuffling residues within each fragment (see Methods). Correspondingly, we considered
228 all natural sequences, whether or not they are part of a structured domain and thus included linker
229 sequences and intrinsically unstructured regions. The extent to which this model is an accurate
230 approximation of natural domains will be discussed below.

231 Comparing the G-model to the A- and T-models over the bacterial dataset shows a dampened wave
232 for the residuals, with the same shape, but a decreased amplitude (*Figure 2: B*). The total residual is
233 correspondingly smaller by a factor of about 2 for all fragment lengths (*Figure 3: A*), implying that

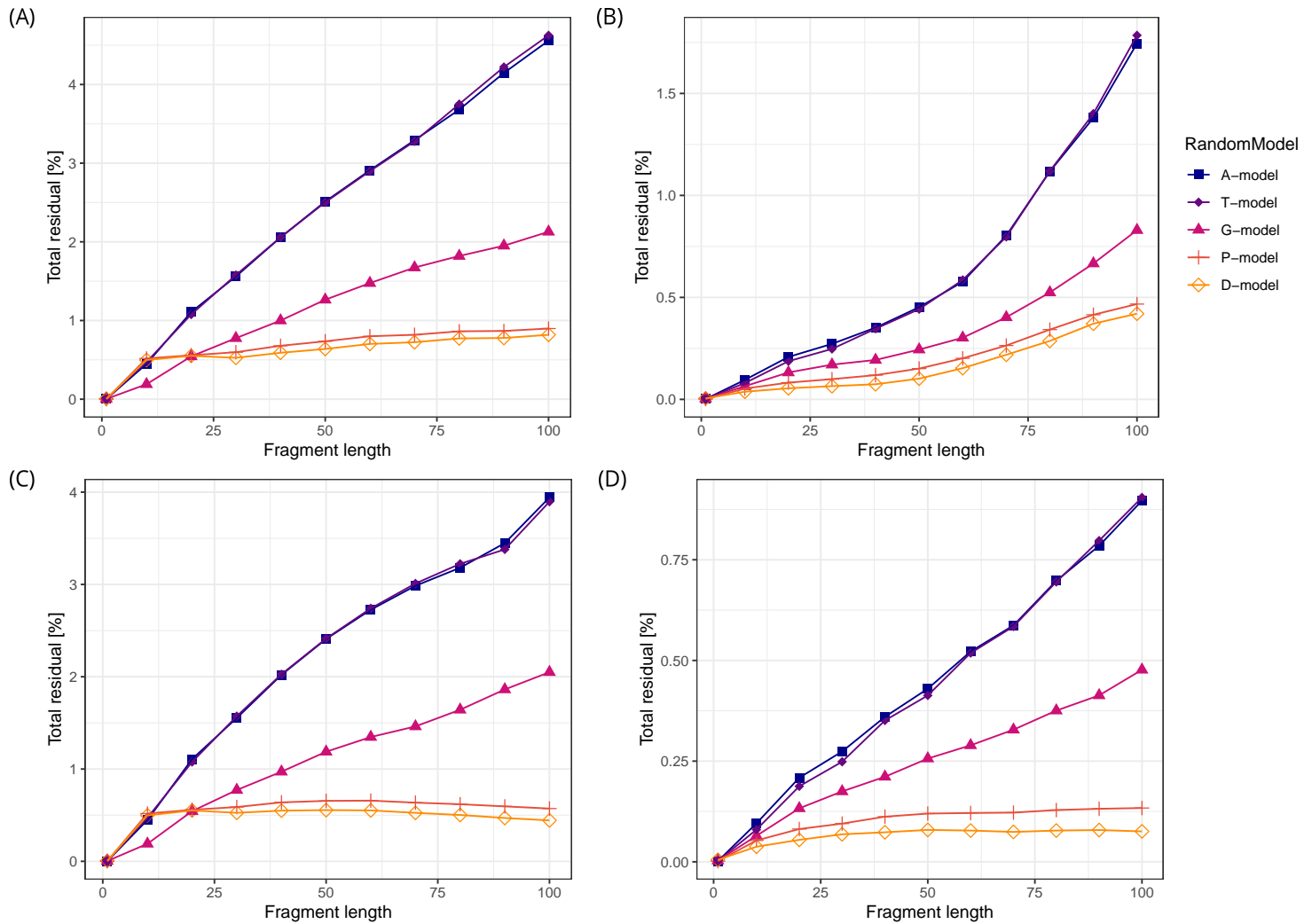


Figure 3. Deviation of random sequence models from natural sequences as a function of fragment length. (A) Total residuals when using a local Smith-Waterman alignment. The total residual indicates the extent to which the distance distribution of random sequence models deviates from the natural. It reflects 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 (overall amino acid composition) and the T-model (overall dipeptide frequency) deviate furthest followed by the G-model (residue composition of 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. A Shift alignment does not penalize beginning and end gaps and prohibits internal gaps. Similar to the Smith-Waterman alignment, the Shift alignment displays an inconsistency at fragment length around 10. (D) Total residuals when using a Hamming distance without 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.

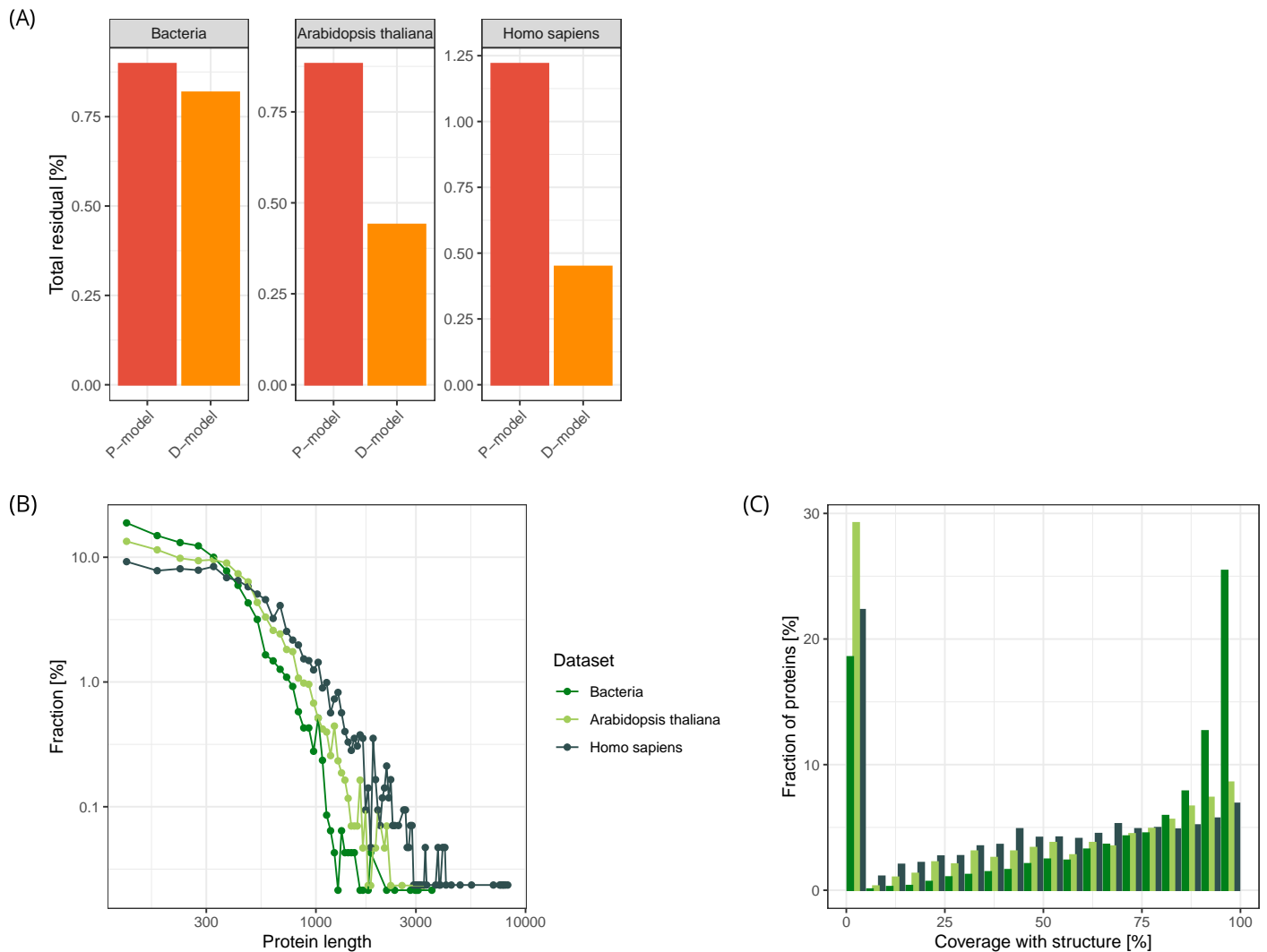


Figure 4. Contrasting the results of our bacterial dataset with those from two eukaryotic proteomes. (A) Total residuals of random models for bacterial dataset, the proteome of *Arabidopsis thaliana* and *Homo sapiens* of the P- and D-models. Relative to the total residual of the P-model, the total residuals 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 40% 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%).

234 controlling for genome composition provides a substantial improvement in modeling the natural
235 distance distribution. A further improvement is clearly achieved with the P-model, even though,
236 at sequence lengths below 20 residues, it produces minor inconsistencies in its total residuals
237 relative to the A-, T-, and G-models (**Figure 3: A**). We suspect that this is an artifact of using local
238 alignments (**Figure 3: A, C**) and, indeed, the effect disappears when using a global alignment as
239 distance metric over the same dataset (**Figure 3: B, D**). As for the A-, T-, and G-models, the residuals
240 of the P-model also have a wave shape, which is however qualitatively different from the shapes
241 for the less local random models, as it has only one crest at an alignment score of 13%. The crest
242 for the unexplained long-range distances is gone, which we attribute to the fact that accounting
243 for composition at the level of individual proteins has introduced the heterogeneity of natural
244 sequences into the random model. For 100mers the total residual of the P-model is 0.9% (**Figure 3:**
245 **A**), a value that is not improved remarkably by an even greater locality: The residuals of the D-model
246 have the same wave shape as those of the P-model and a comparable amplitude, providing only
247 a minor improvement with a total residual of 0.8%. This was somewhat surprising, as it is well
248 established that many proteins are composed of disparate parts such as domains of distinct fold
249 classes, intrinsically unstructured regions or fibrous parts, which are known to be characterized
250 by different residue compositions *Dosztányi et al. (2005)*. The composition of proteins that are
251 composed of heterogeneous parts should thus be scrambled in the P-model and preserved in the
252 D-model. We therefore expected that the D-model would provide a clearer improvement over the
253 P-model.

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

255 We see two reasons why the total residuals of the D- and the P-models are almost identical. One is
256 a technical reason, namely that there is no room for fluctuation of local residue composition in our
257 bacterial dataset, as it may comprise a large number of short and single-domain proteins. The other
258 is a potential qualitative characteristic of our dataset, namely that in long bacterial proteins the local
259 residue composition does not fluctuate remarkably. In order to distinguish how these two reasons
260 contribute to the comparable total residuals of the D- and P-models, we added two eukaryotic
261 datasets for comparison to the following analysis. We retrieved the highly curated proteomes of
262 *Homo Sapiens* and *Arabidopsis thaliana* from UniRef *Apweiler (2009)* and pruned them according to
263 the procedure used for our bacterial dataset. Comparisons of total residuals between the bacterial
264 and eukaryotic datasets show that, whereas the P- and D-models for the bacterial dataset are
265 essentially equivalent, the D-models for the eukaryotic datasets are roughly 2-fold smaller than
266 those of the P-models (**Figure 4: A**), and thus closer to our expectation.

267 In order to evaluate the technical reason, we analyzed sequence lengths in all three datasets and
268 estimated the number of single- and multi-domain proteins. The bacterial dataset has the shortest
269 proteins with a median length of 315 residues, the *Arabidopsis thaliana* dataset a median length
270 of 400 residues and the *Homo sapiens* dataset the longest proteins with a median length of 550
271 residues (**Figure 4: B**). To estimate the number of single and multi-domain proteins, we randomly
272 sampled each of the three datasets and used HHpred *Remmert et al. (2012)* for their domain
273 annotation against the ECOD database *Cheng et al. (2014)*, which represents the most recent and
274 comprehensive classification of domains of known structure (see Methods). ECOD is the current
275 "gold standard" in domain assignments and, at more than 13,000 families, provides a structural
276 basis for most known domains (as captured in databases such as Pfam *Punta et al. (2012)*, SMART
277 *Schultz et al. (1998)* or COGs *Tatusov et al. (2000)*). We considered proteins multi-domain if they
278 had at least 2 domains assigned to them, otherwise we considered them as single-domain proteins.
279 The predicted fraction of multi-domain proteins in our bacterial dataset is 30%, which is smaller in
280 *Arabidopsis thaliana* (25%) and greater in *Homo Sapiens* (35%). The overall length distribution thus
281 indeed correlates with the ratio between the total residuals of the P- and D-models, and potentially
282 contributes to the observed effect, whereas the number of domains per protein does not.

283 In order to evaluate the qualitative reason, namely that sequences of distinct composition are
284 combined within proteins, we assessed the fraction of structured and unstructured regions in the
285 used proteins. To that end, we estimated the fraction of structured regions for each protein with
286 HHpred against the ECOD database (**Figure 4: C**). For the bacterial dataset, 40% of all sampled
287 proteins are predicted to be structured over >90% of their sequence, a fraction that is smaller in
288 *Arabidopsis thaliana* (15%) and *Homo Sapiens* (13%). The structure content of proteins thus also
289 correlates with the ratio between the total residuals of the P- and D-models (**Figure 4: A**), possibly
290 because scrambling between structured and unstructured regions leads to greater compositional
291 disturbance than scrambling within these regions.

292 We conclude that the D-model approximates the natural distance distribution better than the
293 P-model in all cases, however in a more pronounced way for datasets containing heterogeneous
294 mixtures of long sequences combining structured with unstructured regions. In our analysis, these
295 effects were more pronounced in eukaryotic than in bacterial proteins.

296 **Sequence bias caused by homology**

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

305 The detection of homologous relationships requires advanced approaches, which are computa-
306 tionally much more expensive than the simple sequence alignments used to determine distances
307 in sequence space. We therefore only considered a small subset of our sequences and their rela-
308 tionships within this subset, which could be derived computationally in a reasonable amount of
309 time. For this, we randomly sampled our natural dataset to form 10 unbiased groups of 100mers,
310 containing approximately 650 sequences each. We used HHblits to generate profile Hidden Markov
311 Models (HMMs) for all individual sequences within these groups, then derived a set of relationships
312 by aligning the retrieved HMMs from one set of sequences to those of another. This we repeated
313 for arbitrary sets of 100mers, resulting in multiple unbiased samples of relationships. The likelihood
314 of homology between two HMMs was derived using the tool HHalign and required a strict threshold
315 of minimally 90% probability (see Methods). This process identified 0.11% of pairwise relationships
316 as homologous (**Figure 5: A, yellow**), with a standard error of the mean (SEM) of 0.0033% (**Figure 5:**
317 **A**).

318 For the remaining sequence pairs, we evaluated the likelihood of analogy by comparing their
319 HMMs to those of the ECOD database **Cheng et al. (2014)**. By virtue of containing only domains of
320 known structure, ECOD is the currently best resource for distinguishing between homology and
321 analogy in protein domains. For our analysis, we scored pairs of sequences as analogous if they
322 matched distinct X-groups in the ECOD hierarchy using the same probability cutoff of 90% as for the
323 homology assignment. In most cases, the X-level is the highest level at which homology still needs
324 to be considered as a possibility; requiring fragments to match different X-groups within this level
325 thus provided a conservative estimate of analogous relationships. This process identified 52.22% of
326 pairwise relationships as analogous (**Figure 5: A, purple**), with a SEM of 0.84%. We conclude from
327 this that the number of confident analogous pairs exceeds the number of confident homologous
328 pairs by more than 2 orders of magnitude. This already indicates that the influence of homology
329 on the global distance distribution in natural sequences will be dwarfed by analogy. All sequence
330 pairs that could not be confidently assigned to either group were considered to be of unknown

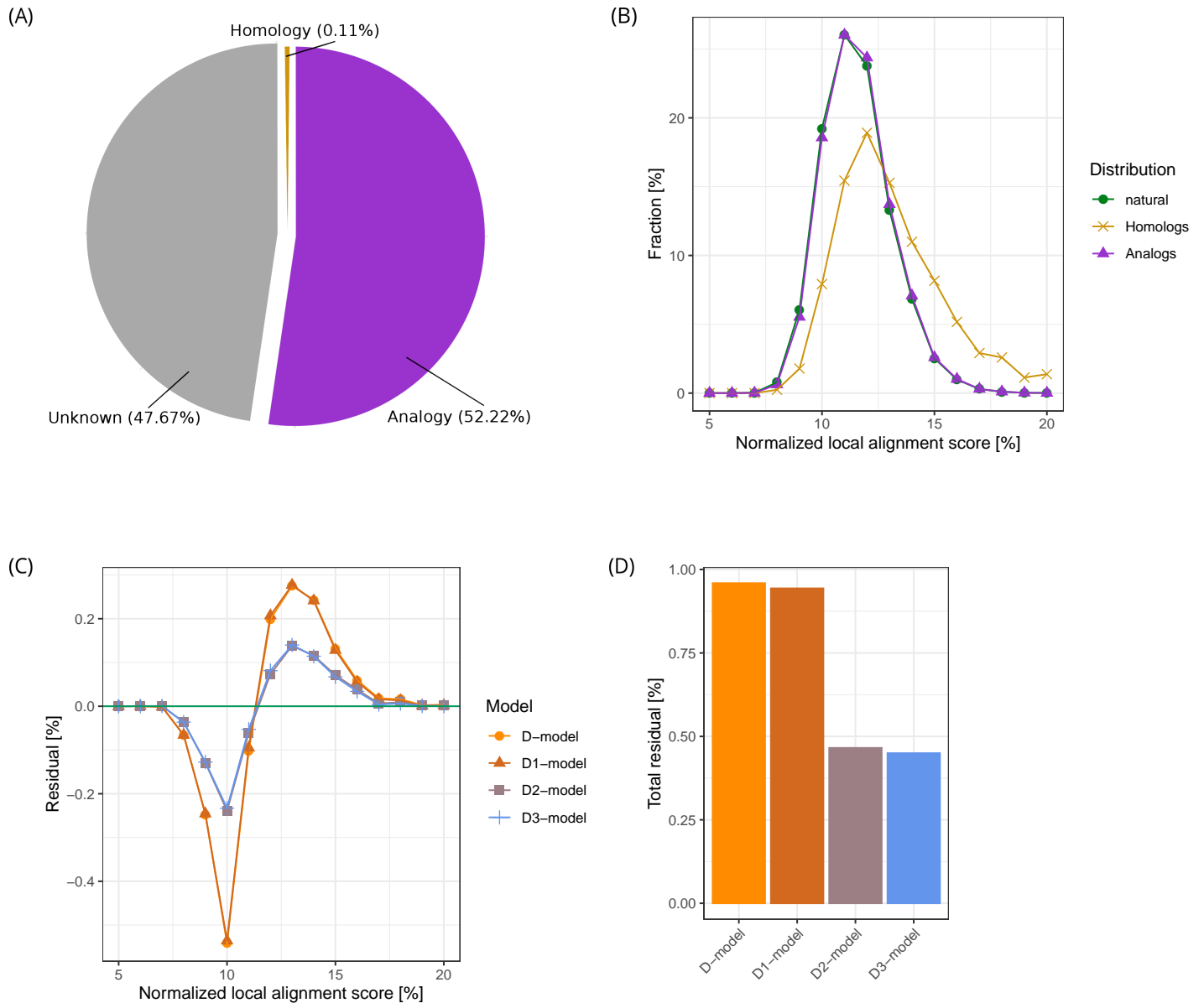


Figure 5. The contribution of homology and analogy to the global occupation of sequence space. (A) Decomposition of fragment pairs into their origins. We sampled 2 Million fragments pairs and analyzed if their relationship is confidently homologous or analogous. The fraction of analogous relationships 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.

331 relationship, amounting to 47.6% of the total with a SEM of 0.84% (**Figure 5: A**, grey).

332 Having decomposed sequence pairs into confident homologous and analogous relationships, we
333 analyzed to what extent the remaining total residual (0.8%) can be explained by incorporating
334 corresponding sequence biases into our D-model. Therefore, we generated three new hybrid
335 models in the following way: we omitted either homologous pairs, or analogous pairs, or both from
336 our set of assigned relationships, generated a D-model for the remaining fragment pairs through
337 the same shuffling procedure as used previously, and then added back the omitted pairs without
338 shuffling. In the following we refer to the hybrid model that adds the sequence bias of homologs to
339 the domain composition as the D1-model, the one that adds the sequence bias of analogs as the
340 D2-model, and the one that adds both biases as the D3-model.

341 The residuals of these three models are compared to that of the D-model in (**Figure 5: C**). Due to
342 the reduced sampling over only 2 million fragment pairs, instead of 500 million, the total residual of
343 the D-model in this analysis deviates slightly from that obtained over the entire dataset and has a
344 value of 0.96% instead of 0.8% (**Figure 5: D**).

345 Relative to this total residual of the purely compositional D-model, the D1-model, which includes
346 homologous sequence effects, is only minimally better (total residual reduced by 0.016%) at ap-
347 proximating the natural distance distribution (**Figure 5: D**). We assume that two reasons are mainly
348 responsible for this only minor improvement: First, the proportion of homologous relationships is
349 only 0.11%, giving them little leverage. Second, the distance distribution of homologs (**Figure 5: B**,
350 yellow) differs only to a small extent from the distance distribution of the natural dataset. This is
351 not entirely unexpected, given how difficult it is to distinguish distant homology from random fluc-
352 tuation in sequence comparisons. In fact, it has been recognized previously that most homologous
353 sequences share no significant similarity *Rost (1997)*.

354 In contrast, the total residual of the D2-model (0.46%), which includes analogous sequence effects, is
355 decreased about 2-fold relative to the D-model. Thus, although analogs have a distance distribution
356 that is very similar to the natural (**Figure 5: B**, purple and green), their leverage is 2 orders of
357 magnitude higher than that of homologs, causing these small differences to improve substantially
358 the fit of the D2-model to the natural distance distribution. This is again not entirely unexpected, as
359 most sequences in our natural dataset share the ability to form secondary structures (**Figure 4: C**),
360 resulting in a sequence bias that is not fully captured by residue composition *Pande et al. (1994)*
361 *Lavelle and Pearson (2009)*. As expected from the D1-model, adding the homologous sequence bias
362 to the D2-model did not really improve its ability to approximate the natural distance distribution.
363 We conclude that the sequence space of natural proteins is almost entirely shaped by compositional
364 effects and that the remaining sequence bias is almost entirely due to analogy, which we interpret
365 to result from secondary structure formation.

366 Conclusion

367 In this article we have undertaken a study of natural protein sequence space, using an approach
368 built on the probability mass function of pairwise distances between sequence fragments. With
369 this approach we were able to analyze the occupation of sequence space by fragments up to 100
370 residues in length, substantially exceeding previous efforts and for the first time characterizing
371 globally the relative position of sequences in space. Our results show that the global compositional
372 bias of natural proteins is already sufficient to approximate the distance distribution of natural
373 sequences by 95.4% and that accounting for local compositional bias down to the level of individual
374 100mers further improves this to 99.2%. The remaining 0.8% of unaccounted distances between
375 natural 100mers are almost entirely contributed by sequence effects arising from analogous
376 relationships, leaving only a negligible contribution to homology in the global characterization of
377 sequence space occupation.

378 This surprised us, as decades of bioinformatic work have mapped out an increasingly comprehen-
379 sive description of sequence space around protein families, based on the detection of ever more
380 remote homology. We therefore expected to find that homology also has a substantial role in
381 shaping the global structure of sequence space occupation, corresponding to the image of islands
382 formed by natural sequences within a global sea of possibilities. This expectation was not borne
383 out and in retrospect this might not seem as surprising, given that even within protein families,
384 the influence of homology is smaller than generally perceived. This is substantially due to the
385 way in which family relationships are represented, strongly emphasizing common features (such
386 as the generally few conserved residues) and omitting variable ones. This focus on biological
387 significance over raw sequence similarity leads to a perception of sequence space that is distorted
388 by *evolutionary distance* and does not reflect actual distances. Evidence for this can be seen for
389 example in the progressively more complex statistical methods needed to substantiate homology
390 across increasingly large evolutionary distances, the resulting difficulties to classify the detected
391 relationships into a hierarchy of protein families and superfamilies, and the remaining inability in
392 many cases to judge on the homologous or analogous nature of similarities, even in the presence
393 of extensive sequence and structure information *Rost (1997)*. These considerations show that even
394 at the level of protein families, many sequence relationships comprise a large random element,
395 substantially indistinguishable from random fluctuation and sequence convergence. This random
396 element not only results from our inability to detect homologs that have diverged strongly due to
397 low selective pressure, but also from the fact that in many families, a conserved core has been
398 elaborated in different ways with analogous sequences.

399 We find a much larger influence of analogous sequence biases on the global shape of naturally
400 occupied sequence space. The main common feature of proteins in our natural dataset is the ability
401 to fold, which translates into a propensity to assume secondary structure locally. We see this as the
402 main reason for the sequence bias that we observe between analogous sequences. Nevertheless,
403 the sequence biases of homology and analogy together account for only 0.8% of all distances
404 between natural sequences. We conclude that natural sequences stand out from randomness
405 primarily through their biased use of the 20 amino acids. Accounting for this bias at increasingly
406 local levels is largely sufficient to model the global structure of sequence space occupation. This
407 major relevance of composition has been acknowledged as it has been implemented into BLAST
408 *Schaffer (2002)* and been demonstrated to be key for the aggregation of intrinsically unstructured
409 proteins *Vymětal et al. (2019)*.

410 There seems to be no other striking feature of the primary structure in natural protein sequences
411 and in consequence there are also no other obvious features that distinguish natural from random
412 sequences. We conclude that viable proteins could be located anywhere in the sequence space
413 defined by natural residue compositions. The main reason why the proteome of nature currently
414 only comprises some 10^{12} proteins *Lupas and Koretke (2008)* and that these mainly fall into only
415 about 10^4 families *Punta et al. (2012)* is therefore not due to the limited availability of useful
416 sequence space, but rather to their evolutionary history. There is treasure everywhere.

417 **Methods**

418 **Natural data**

419 **Genome selection**

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

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

430 Genome curation

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

435 Low complexity filtering

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

446 Sequence adjustments

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

454 Complete statistics and data availability

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

- 457 • fasta-file of original genomes
- 458 • fasta-file of adjusted genomes
- 459 • overall amino acid composition

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

461 Fragment selection

462 We selected random fragments such that each character (amino acids and end-of-line-character) in
463 the dataset had the same probability of being chosen and that the same fragment pair would never
464 be chosen twice. We ensured this by implementing two linear congruential generator *Press et al.*
465 *(2007)* to enumerate all possible pairs of fragments. In detail, one linear congruential generator
466 was used for each member of the pair with multiplier $a = 1$ and moduli $m_1 = 223$ and $m_2 = 34,211$,
467 where both moduli are prime numbers relative to the total number of characters 1,168,754,000.
468 Depending on the starting points of the two generators, a different subset of index pairs can be
469 selected. This enabled us to calculate disjunctive fragment pairs in parallel. We selected $5 \cdot 10^8$ valid
470 pairs of fragments to accurately estimate the distance distributions and rejected fragments that
471 straddled protein boundaries or invalid regions, indicated by the end-of-line-character.

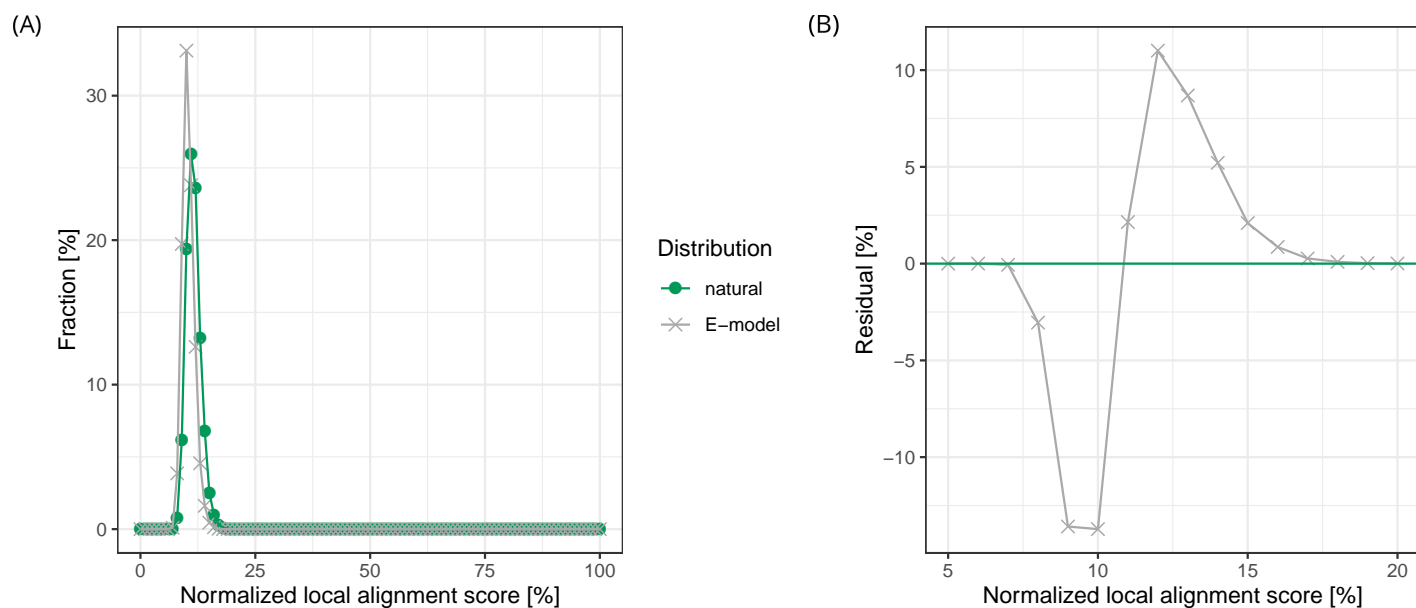


Figure 6. The E-model deviates severely from the natural dataset. (A) Distance distributions of natural dataset and E-model. In contrast to the primarily used models (A-, T-, G-, P-, and D-models) in this paper, the distance distribution of the E-model has an obvious deviation from that of the natural data. (B) Residuals of the E-model. Compared to the distances between fragments derived from the E-model, the distances between natural sequences have a strong tendency towards being shorter.

472 Models incorporating overall amino acid composition

473 The most standard random sequence model is based on the underlying amino acid composition of
474 a given dataset. We obtained randomized data for this A-model by randomly shuffling all amino
475 acids of the natural data. Thereby, protein length is maintained and the number of amino acids
476 stays exactly the same. As all our random models are based on random permutations, we used
477 the Mersenne Twister algorithm mt19337 of the C++ 14 std library with the standard seed value of
478 19650218. This algorithm is considered one of the best pseudo-random number generators and in
479 a test with a smaller dataset we found that our results did not depend on the type or seeding of the
480 random number generator.

481 For the E-model, we proceeded the same way as for the A-model. The only difference is that we
482 replaced the natural dataset, by writing over all valid amino acids with the 20 possible amino acids
483 in lexicographical order. When reaching the character Y for tryptophan, we started over with A for
484 alanine. The distance distribution of the E-model deviates severely from that of the natural dataset
485 (**Figure 6**) with a total residual of 30.4% for 100mers.

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

487 To account for genome or protein composition, we shuffled amino acids within the context of
488 genomes or proteins. For the G-model, we shuffled valid amino acids within each of the 1,307
489 genomes. For the P-model we shuffled valid amino acids within each protein. We used one instance
490 for genome and protein composition bias and stored them to generate the distance distribution for
491 the corresponding models. After shuffling, we divided proteins containing the invalid amino acid X
492 by replacing it with an end-of-line-character.

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

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

498 we select pairs of natural fragments as described above and consider the target fragment of length
499 N to be located in the middle of the domain. If the domain straddles any protein boundaries, we
500 adjust the domain boundaries such that the domain fits into the protein boundaries by shifting
501 to the right (starting sequences) or to the left (terminal sequences). Note that because of this
502 adjustment, the selection probability of amino acids into domains is not uniform but the selection
503 probability of amino acids in fragments is. The alternative would be a rejection procedure, where
504 we would reject fragments that are so close to protein boundaries that the domain of length of
505 100 would not fit. The downside of such a rejection procedure is that fragments close to protein
506 boundaries are not selected and hence the selection probability of fragments is not uniform
507 anymore, which differs from the selection of natural fragments or fragments for the A-, G-, and
508 P-models. The D1, D2, and D3-models, which incorporate the sequence bias of homologous and
509 analogous fragments, are presented further down.

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

511 Distance metric

512 We define the distance between two fragments of the same length N as the normalized rounded
513 score s from a Smith-Waterman alignment. In the alignment, an amino acid match is scored with
514 1, a mismatch with 0, gap opening penalty is equal to 3 and gap extension penalty is 0.1, which
515 are the same parameters for gaps as used in *Rost (1999)* *Schneider et al. (1997)*. Due to gaps, the
516 alignment scores p can rank between 0 and N in 0.1 steps; to obtain integer distances, we round
517 scores to the closest integer number. Distances exactly between two integers (such as 1.5) are
518 assigned to the smaller one. To compare the score p across different fragment lengths N , we
519 transform it into the normalized score s , scaling between 0-100%, as follows:

$$s = \frac{\text{round}(p)}{N}$$

520 This score s thereby reflects the number of dimensions in sequence space (positions in sequence),
521 which differ between two fragments, while allowing for gaps and insertions. In some cases, we use
522 the normalized global alignment score from a Needleman-Wunsch alignment with identical scoring
523 parameters, the Hamming distance or a Shift metric that allows for terminating and starting gap
524 without penalties, to illustrate differences to the used Smith-Waterman metric.

525 We also diversified gap penalties, leading to comparable results (data not shown). For all alignments,
526 we used the SeqAn C++ library, version 2.4 *Rahn et al. (2018)*, which enables many sequence
527 comparisons in parallel.

528 Comparing distance distributions

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

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

533 where s is the alignment score. For residuals $r(s)$ exceeding zero, there is a higher frequency of
534 these alignment scores in natural fragments relative to random fragments.

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

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

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

540 **Decomposition into homologous and analogous relationships**

541 Homology

542 We derived the fraction of sequence pairs that are confidently homologous using the tools of
543 HH-suite (version number 3.0.3) *Remmert et al. (2012)*. To derive this fraction, we systematically
544 sampled our dataset and extracted 10 sets of natural 100mers that are equally distributed over our
545 dataset, each containing approximately 650 fragments. With HHblits, we generated HMMs with the
546 standard settings for each of these fragments with two iterations, using uniclust30 as underlying
547 database (version August 2018).

548 Then, we pairwise aligned the generated HMMs with HHalign, in order to estimate whether two
549 fragments are homologous. We did this by aligning all fragments in one set to all of those in another
550 set, resulting in 90 possible directed combinations of which we chose 10 as representative sets of
551 pairwise relationships. Each set of fragments was considered twice in this comparison, once as the
552 set of query sequences and once as the set of target sequences in the alignment. This resulted
553 in 2 Million pairwise fragment comparisons divided into 10 disjunctive sets. Pairs of fragments
554 were considered to be homologous, if HHalign predicted them to be homologous with a probability
555 above 90%. In total 0.11% of the fragment pairs were found to be homologous; the standard error
556 of the mean (SEM) derived from the 10 sets of 0.0033%.

557 Analogy

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

562 For the assignment of structured domains, we used the ECOD classification *Cheng et al. (2014)*,
563 which is the currently best resource for distinguishing between homology and analogy in protein
564 domains. The HMMs of each 100mer (same as in the homology detection) were thereby compared
565 against all ECOD entries (retrieved on 9.4.2019) with HHsearch. We used HHsearch with the standard
566 parameter and assigned the best-scoring non-overlapping hits with a probability above 90% to the
567 corresponding fragment. Of all 100mers 70% could be assigned to a single domain and less than
568 1% to multiple domains, of which we considered each. Other 100mers were not assigned to any
569 domain, which we directly excluded to be analogous to any other sequence, since we are uncertain
570 about their origin.

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

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

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

580 In order to estimate the influence of homology and analogy to the natural distance distribution,
581 we generated mixed models that that account for their sequence bias. The D1-model includes the
582 homologous sequence bias by including the distances between all confidently homologous fragment
583 pairs without shuffling. We applied the D-model to the remaining fragment pairs and shuffled
584 the fragments of the corresponding pairs that are not homologous with the Unix command shuf
585 followed by deriving their distance. All distances combined resulted into the distance distribution
586 of the D1-model. The sequence bias between homologous fragments is therein preserved while for
587 other fragment pairs only their composition is accounted for. We proceeded the same way for the
588 D2-model by including the distances of unshuffled fragments that are confidently analogous, and
589 distances of the remaining pairs after shuffling the residues within each fragment. For the D3-model
590 we included both sequence bias of homologous and analogous natural fragments.

References

- 591
592 **Alva V**, Remmert M, Biegert A, Lupas AN, Söding J. A galaxy of folds. *Protein Science*. 2009 jan;
593 19(1). <http://www.ncbi.nlm.nih.gov/pubmed/19937658><http://www.pubmedcentral.nih.gov/articlerender.fcgi?>
594 [artid=PMC2817847](http://www.ncbi.nlm.nih.gov/pubmed/19937658)<http://doi.wiley.com/10.1002/pro.297>, doi: 10.1002/pro.297.
- 595 **Alva V**, Söding J, Lupas AN. A vocabulary of ancient peptides at the origin of folded proteins. *eLife*. 2015; 4. doi:
596 10.7554/eLife.09410.001.
- 597 **Apweiler R**. The universal protein resource (UniProt) in 2010. *Nucleic Acids Research*. 2009 jan; 38:D190–
598 5. <http://www.ncbi.nlm.nih.gov/pubmed/18045787><http://www.pubmedcentral.nih.gov/articlerender.fcgi?>
599 [artid=PMC2238893](http://www.ncbi.nlm.nih.gov/pubmed/18045787), doi: 10.1093/nar/gkp846.
- 600 **Bershtein S**, Serohijos AW, Shakhnovich EI. Bridging the physical scales in evolutionary biology: from protein
601 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.
- 602
- 603 **Cedano J**, Aloy P, Perez-Pons JA, Querol E. Relation between amino acid composition and cellular location of
604 proteins. *Journal of Molecular Biology*. 1997 feb; 266(3):594–600. [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S0022283696908049)
605 [article/pii/S0022283696908049](https://www.sciencedirect.com/science/article/pii/S0022283696908049), doi: 10.1006/jmbi.1996.0804.
- 606 **Cheng H**, Schaeffer RD, Liao Y, Kinch LN, Pei J, Shi S, Kim BH, Grishin NV. ECOD: An Evolutionary Classification
607 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)
608 [journal.pcbi.1003926](https://dx.plos.org/10.1371/journal.pcbi.1003926), doi: 10.1371/journal.pcbi.1003926.
- 609 **Chevalier A**, Silva DA, Rocklin GJ, Hicks DR, Vergara R, Murapa P, Bernard SM, Zhang L, Lam KH, Yao G, Bahl
610 CD, Miyashita SI, Goreshnik I, Fuller JT, Koday MT, Jenkins CM, Colvin T, Carter L, Bohn A, Bryan CM, et al.
611 Massively parallel de novo protein design for targeted therapeutics. *Nature*. 2017 oct; 550(7674):74–79.
612 <http://www.nature.com/articles/nature23912>, doi: 10.1038/nature23912.
- 613 **Chou KC**. Prediction of protein cellular attributes using pseudo amino acid composition. *PROTEINS Struct Funct*
614 *Genet* (Erratum *ibid*, 2001, Vol44, 60). 2001; 43(3):246–255. doi: 10.1017/CBO9781107415324.004.
- 615 **Deza MM**, Deza E. *Encyclopedia of Distances*. Springer Berlin Heidelberg; 2014. [https://books.google.de/books?](https://books.google.de/books?id=q_7FBAAAQBAJ)
616 [id=q_7FBAAAQBAJ](https://books.google.de/books?id=q_7FBAAAQBAJ).
- 617 **Diggle PJ**. *Statistical Analysis of Spatial and Spatio-Temporal Point Patterns*. CRC press; 2014.
- 618 **Dosztányi Z**, Csizsók V, Tompa P, Simon I. The pairwise energy content estimated from amino acid composition
619 discriminates between folded and intrinsically unstructured proteins. *Journal of Molecular Biology*. 2005;
620 347(4):827–839. doi: 10.1016/j.jmb.2005.01.071.
- 621 **Duranton G**, Overman HG. Testing for Localization Using Micro-Geographic Data. *The Review of Economic*
622 *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)
623 [00362](https://academic.oup.com/restud/article-lookup/doi/10.1111/0034-6527.00362), doi: 10.1111/0034-6527.00362.
- 624 **Fukuchi S**, Nishikawa K. Protein surface amino acid compositions distinctively differ between thermophilic and
625 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)
626 [S0022283601947187](https://www.sciencedirect.com/science/article/pii/S0022283601947187)[?via=ihub](https://www.sciencedirect.com/science/article/pii/S0022283601947187?via=ihub), doi: 10.1006/jmbi.2001.4718.
- 627 **Fukuchi S**, Yoshimune K, Wakayama M, Moriguchi M, Nishikawa K. Unique amino acid composition of proteins
628 in halophilic bacteria. *J Mol Biol*. 2003; 327(2):347–357. doi: 10.1016/S0022-2836(03)00150-5.
- 629 **Harms MJ**, Thornton JW. Historical contingency and its biophysical basis in glucocorticoid receptor evolu-
630 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)
631 [pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4447330](http://www.ncbi.nlm.nih.gov/pubmed/24930765), doi: 10.1038/nature13410.
- 632 **Huang PS**, Boyken SE, Baker D. The coming of age of de novo protein design. *Nature*. 2016; 537(7620):320–
633 7. <http://www.nature.com/doi/10.1038/nature19946>[http://www.ncbi.nlm.nih.gov/pubmed/](http://www.ncbi.nlm.nih.gov/pubmed/27629638)
634 [27629638](http://www.nature.com/doi/10.1038/nature19946), doi: 10.1038/nature19946.
- 635 **Larson SM**, England JL, Desjarlais JR, Pande VS. Thoroughly sampling sequence space: large-scale protein design
636 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)
637 [fcgi?artid=2373757](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2373757)[&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2373757&tool=pmcentrez&rendertype=abstract), doi: 10.1110/ps.0203902.
- 638 **Lavelle DT**, Pearson WR. Globally, unrelated protein sequences appear random. *Bioinformatics*. 2009; 26(3):310–
639 318. doi: 10.1093/bioinformatics/btp660.

- 640 **de Lucrezia D**, Slanzi D, Poli I, Polticelli F, Minervini G. Do natural proteins differ from random sequences polypep-
641 tides? natural vs. random proteins classification using an evolutionary neural network. *PLoS ONE*. 2012 may;
642 7(5):e36634. <https://dx.plos.org/10.1371/journal.pone.0036634>, doi: 10.1371/journal.pone.0036634.
- 643 **Luigi Luisi P**. Contingency and determinism. *Philos Trans R Soc A Math Phys Eng Sci*. 2003; 361(1807):1141–1147.
644 <http://rsta.royalsocietypublishing.org/cgi/doi/10.1098/rsta.2003.1189>, doi: 10.1098/rsta.2003.1189.
- 645 **Lupas A**, Koretke K. In: *Evolution of Protein Folds*; 2008. p. 131–151. doi: 10.1142/9789812778789_0006.
- 646 **Lupas AN**, Dyke MV, Stock J. Predicting Coiled Coils from Protein Sequences. *Science*. 1991; 252:1162–1164.
- 647 **Munteanu CR**, González-Díaz H, Borges F, de Magalhães AL. Natural/random protein classification mod-
648 els based on star network topological indices. *Journal of Theoretical Biology*. 2008; 254(4):775–783. doi:
649 [10.1016/j.jtbi.2008.07.018](https://doi.org/10.1016/j.jtbi.2008.07.018).
- 650 **Pande VS**, Grosberg aY, Tanaka T. Nonrandomness in protein sequences: evidence for a physically driven stage
651 of evolution? *Proc Natl Acad Sci U S A*. 1994; 91(26):12972–12975. doi: [10.1073/pnas.91.26.12972](https://doi.org/10.1073/pnas.91.26.12972).
- 652 **Parks DH**, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil PA, Hugenholtz P. A standardized
653 bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat Biotechnol*. 2018;
654 36(10):996. doi: [10.1038/nbt.4229](https://doi.org/10.1038/nbt.4229).
- 655 **Poznański J**, Topiński J, Muszewska A, Dębski KJ, Hoffman-Sommer M, Pawłowski K, Grynberg M. Global
656 pentapeptide statistics are far away from expected distributions. *Sci Rep*. 2018 dec; 8(1):15178. [http:](http://www.nature.com/articles/s41598-018-33433-8)
657 [://www.nature.com/articles/s41598-018-33433-8](http://www.nature.com/articles/s41598-018-33433-8), doi: 10.1038/s41598-018-33433-8.
- 658 **Press WH**, Teukolsky SA, Vetterling WT, Flannery BP. *Numerical Recipes 3rd Edition: The Art of Scientific*
659 *Computing*. 3 ed. New York, NY, USA: Cambridge University Press; 2007.
- 660 **Punta M**, Coggill PC, Eberhardt RY, Mistry J, Tate J, Boursnell C, Pang N, Forslund K, Ceric G, Clements J, Heger
661 A, Holm L, Sonnhammer ELL, Eddy SR, Bateman A, Finn RD. The Pfam protein families database. *Nucleic*
662 *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)
663 [gkr1065](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkr1065), doi: 10.1093/nar/gkr1065.
- 664 **Rahn R**, Budach S, Costanza P, Ehrhardt M, Hancox J, Reinert K. Generic accelerated sequence alignment in
665 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)
666 [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/bioinform-
667 [atics/bty380](https://academic.oup.com/bioinformatics/advance-article-abstract/doi/10.1093/bioinformatics/bty380/4992147).
- 668 **Remmert M**, Biegert A, Hauser A, Söding J. HHblits: lightning-fast iterative protein sequence searching by
669 HMM-HMM alignment. *Nature Methods*. 2012 feb; 9(2):173–175. <http://www.nature.com/articles/nmeth.1818>,
670 doi: [10.1038/nmeth.1818](https://doi.org/10.1038/nmeth.1818).
- 671 **Rost B**. Twilight zone of protein sequence alignments. *Protein Engineering Design and Selection*. 1999;
672 12(2):85–94. doi: [10.1093/protein/12.2.85](https://doi.org/10.1093/protein/12.2.85).
- 673 **Rost B**. Protein structures sustain evolutionary drift. *Folding and Design*. 1997 jun; 2(3):S19–S24. [https:](https://www.sciencedirect.com/science/article/pii/S135902789700059X)
674 [://www.sciencedirect.com/science/article/pii/S135902789700059X](https://www.sciencedirect.com/science/article/pii/S135902789700059X), doi: 10.1016/S1359-0278(97)00059-X.
- 675 **Schaffer AA**. Improving the accuracy of PSI-BLAST protein database searches with composition-based statistics
676 and other refinements. *Nucleic Acids Research*. 2002; 29(14):2994–3005. doi: [10.1093/nar/29.14.2994](https://doi.org/10.1093/nar/29.14.2994).
- 677 **Schneider R**, de Daruvar A, Sander C. The HSSP database of protein structure-sequence alignments. *Nucleic*
678 *Acids Res*. 1997 jan; 25(1):226–230. <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/25.1.226>,
679 doi: [10.1093/nar/25.1.226](https://doi.org/10.1093/nar/25.1.226).
- 680 **Schultz J**, Milpetz F, Bork P, Ponting CP. SMART, a simple modular architecture research tool: Identification of
681 signaling domains. *PNAS*. 1998; 95(May):5857–5864. doi: [10.1073/pnas.95.11.5857](https://doi.org/10.1073/pnas.95.11.5857).
- 682 **Shah P**, McCandlish DM, Plotkin JB. Contingency and entrenchment in protein evolution under purifying
683 selection. . 2015; 2015. <http://arxiv.org/abs/1404.4005>{%}0Ahttp://dx.doi.org/10.1073/pnas.1412933112, doi:
684 [10.1073/pnas.1412933112](https://doi.org/10.1073/pnas.1412933112).
- 685 **Starr TN**, Picton LK, Thornton JW. Alternative evolutionary histories in the sequence space of an ancient protein.
686 *Nat Publ Gr*. 2017; 549. <https://www.nature.com/articles/nature23902.pdf>, doi: 10.1038/nature23902.

- 687 **Starr TN**, Thornton JW. Epistasis in protein evolution. . 2016; 25(7):1204–1218. <http://www.ncbi.nlm.nih.gov/pubmed/26833806><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4918427>, doi: 10.1002/pro.2897.
- 688
689
- 690 **Stiffler MA**, Poelwijk FJ, Brock K, Stein RR, Teyra J, Sidhu S, Marks DS, Gauthier NP, Sander C. Protein structure from experimental evolution. bioRxiv. 2019; <https://www.biorxiv.org/content/early/2019/06/13/667790>, doi: 10.1101/667790.
- 691
692
- 693 **Tatusov RL**, Galperin MY, Natale DA, Koonin EV. The COG database: a tool for genome-scale analysis of protein functions and evolution. Nucleic Acids Research. 2000 jan; 28(1):33–36. <http://www.ncbi.nlm.nih.gov/pubmed/10592175><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC102395>, doi: 10.1093/nar/28.1.33.
- 694
695
- 696 **Urlinger S**, Baron U, Thellmann M, Hasan MT, Bujard H, Hillen W. Exploring the sequence space for tetracycline-dependent transcriptional activators: Novel mutations yield expanded range and sensitivity. Proc Natl Acad Sci. 2000 jul; 97(14):7963–7968. doi: 10.1073/PNAS.130192197.
- 697
698
- 699 **Vallee BL**, Auld DS. Zinc Coordination, Function, and Structure of Zinc Enzymes and Other Proteins. Biochemistry. 1990; 29(24):5647–5659. doi: 10.1021/bi00476a001.
- 700
- 701 **Vymětal J**, Vondrášek J, Hlouchová K. Sequence Versus Composition: What Prescribes IDP Biophysical Properties? Entropy. 2019 jul; 21(7):654. <https://www.mdpi.com/1099-4300/21/7/654>, doi: 10.3390/e21070654.
- 702
- 703 **Wheelan SJ**, Marchler-Bauer A, Bryant SH. Domain size distributions can predict domain boundaries. Bioinformatics. 2000; 16(7):613–618. doi: 10.1093/bioinformatics/16.7.613.
- 704
- 705 **Woolfson DN**, Bartlett GJ, Burton AJ, Heal JW, Niitsu A, Thomson AR, Wood CW. De novo protein design: How do we expand into the universe of possible protein structures? . 2015; 33:16–26. <http://dx.doi.org/10.1016/j.sbi.2015.05.0090959-440X/{#}>, doi: 10.1016/j.sbi.2015.05.009.
- 706
707
- 708 **Wootton JC**, Federhen S. Analysis of compositionally biased regions in sequence databases. Methods in Enzymology. 1996 jan; 266:554–571. <https://www.sciencedirect.com/science/article/pii/S0076687996660352>, doi: 10.1016/S0076-6879(96)66035-2.
- 709
710
- 711 **Wüthrich K**. NMR of Proteins and Nucleic Acids; 1986.