## Accelerated Article Preview

# Uncovering new families and folds in the natural protein universe

Janani Durairaj, Andrew M. Waterhouse, Toomas Mets, Tetiana Brodiazhenko, Minhal Abdullah, Gabriel Studer, Gerardo Tauriello, Mehmet Akdel, Antonina Andreeva, Alex Bateman, Tanel Tenson, Vasili Hauryliuk, Torsten Schwede & Joana Pereira

This is a PDF file of a peer-reviewed paper that has been accepted for publication. Although unedited, the content has been subjected to preliminary formatting. Nature is providing this early version of the typeset paper as a service to our authors and readers. The text and figures will undergo copyediting and a proof review before the paper is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers apply.

# Uncovering new families and folds in the natural protein universe

Janani Durairaj[1,2], Andrew M. Waterhouse[1,2], Toomas Mets[3,4], Tetiana Brodiazhenko[3], Minhal Abdullah[3,4], Gabriel Studer[1,2], Gerardo Tauriello[1,2], Mehmet Akdel[5], Antonina Andreeva[6], Alex Bateman[6], Tanel Tenson[3], Vasili Hauryliuk[3,4,7,8], Torsten Schwede[1,2], Joana Pereira[1,2]

[1] Biozentrum, University of Basel, Basel, Switzerland
[2] SIB Swiss Institute of Bioinformatics, University of Basel, Basel, Switzerland
[3] Institute of Technology, University of Tartu, Tartu, Estonia
[4] Department of Experimental Medical Science, Lund University, Lund, Sweden
[5] VantAI, New York, USA
[6] European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Hinxton, United Kingdom
[7] Science for Life Laboratory, Lund, Sweden
8 Virus Centre, Lund University, Lund, Sweden

**Correspondence to:** Joana Pereira (joana.pereira@unibas.ch), Torsten Schwede (torsten.schwede@unibas.ch)

**We are now entering a new era in protein sequence and structure annotation, with hundreds of millions of predicted protein structures made available through the AlphaFold database[1]. These models cover nearly all proteins that are known, including those challenging to annotate for function or putative biological role using standard homology-based approaches. In this study, we examine the extent to which the AlphaFold database has structurally illuminated this "dark matter" of the natural protein universe at high predicted accuracy. We further describe the protein diversity that these models cover as an annotated interactive sequence similarity network, accessible at https://uniprot3d.org/atlas/AFDB90v4. By searching for novelties from sequence, structure, and semantic perspectives, we uncovered the β-flower fold, added multiple protein families to Pfam database[2], and experimentally demonstrate that one of these belongs to a new superfamily of translation-targeting toxin-antitoxin systems, TumE-TumA. This work underscores the value of large-scale efforts in identifying, annotating, and prioritising novel protein families. By leveraging the recent deep learning revolution in protein bioinformatics, we can now shed light into uncharted areas of the protein universe at an unprecedented scale, paving the way to innovations in life sciences and biotechnology.**

Since the sequencing of the first protein, large-scale efforts brought about by faster and cheaper genome sequencing techniques have shed light into some of the sequences that nature has sampled so far. Currently, there are over 350 million unique protein coding sequences deposited in UniProt and over 3 billion in MGnify[3,4]. The rate at which this data is growing is much faster than experimental functional characterization. To close the gap, functional information is gathered for a subset of proteins and the findings extrapolated to close homologs.

45  Manual curation is carried out by those assembling the genomes and by biocurators[5] and
46  incorporated into automated annotation pipelines such as InterPro[6].
47  Despite the great success of such approaches, only 83% of UniProt sequences are covered by
48  InterPro, and many correspond to Domains of Unknown Function (DUF). Thus, numerous
49  protein sequences remain functionally unannotated and unclassified. Some of these may just
50  correspond to divergent forms of known protein families that lie beyond the detection horizon
51  of automated, homology-based methods; others could belong to so far undescribed protein
52  families with yet-to-be determined molecular or biological functions[7].
53  The 3D structure of a protein is intrinsically linked with its molecular function. Experimental
54  structure determination is an expensive and time-consuming process, and homology-based
55  computational prediction loses its power for proteins without close homologs[8].
56  Notwithstanding, deep learning based approaches have recently achieved unprecedented
57  accuracy, with AlphaFold2 at the forefront. Its success drove the establishment of the
58  AlphaFold database (AFDB), which contains predicted structural models for about 215 million
59  natural protein sequences from UniProt, including many of the unannotated proteins. At the
60  same time, deep learning-based approaches have also recently been employed for predicting
61  functional properties from structure[9] and protein names from sequence[10].
62  In this work, we combine sequence similarities and structure features with deep learning-based
63  function prediction tools to shed light on "functionally dark" proteins in UniProt. We revised
64  their proportion, evaluated how many of them now have high confidence structural models that
65  can be leveraged for downstream analysis, and constructed for the first time an annotated and
66  interactive sequence similarity network with millions of proteins. By exploring this network,
67  we discovered 290 putative new protein families, identified at least one novel protein fold, and
68  defined a new superfamily of translation-targeting toxin-antitoxin systems which we
69  experimentally validated and dubbed TumE-TumA. This work demonstrates that functional
70  annotation of proteins, even from a purely computational perspective, requires a combination
71  of data sources and approaches, which become increasingly available and attainable due to the
72  rapid and ongoing advances at the interface between life sciences and deep learning.
73
**Functional darkness in UniProt and AFDB**
75  As of August 2022, there were more than 350 million unique protein sequences in UniProt (i.e.,
76  UniRef100 clusters[11]). We focus our analysis on these as they have a higher confidence than
77  those deposited in metagenomics databases such as MGnify. These sequences correspond to
78  circa 50 million non-redundant proteins when clustered to a maximum sequence identity of
79  50% (UniRef50). Starting from these clusters, we define the "functional brightness" of a given
80  protein as the full-length coverage with annotations of its close homologs, and a UniRef50
81  cluster is as "bright" as the "brightest" sequence it encompasses (Fig. 1a). For that, we only
82  considered those annotations that correspond to domains and families whose title does not
83  include "Putative", "Hypothetical", "Uncharacterised" and "DUF", but considered predicted
84  coiled coil and intrinsically disordered segments in order to focus our analysis solely on
85  functionally dark proteins with a potential for a globular (or other) fold type.
86  We found that 34% of all UniRef50 clusters (10% of UniRef100, ~34 million unique proteins)
87  are dark as they do not reach a functional brightness higher than 5% (Extended data Fig. 1a).
88  While the brightness of a cluster is not directly proportional to the number of sequences within

89    it (Pearson correlation coefficient of 0.0), bright clusters (functional brightness ≥ 95%) tend to

90    be larger than those whose members are poorly annotated (mean 19 $\pm$ 123 unique sequences

91    in bright clusters compared to 2 $\pm$ 7 in dark).

92    While UniRef50 clusters encompass sequences from the UniProt Knowledgebase (UniProtKB)

93    and the UniProt Archive (UniParc)[12], the latest version of AFDB (version 4) covers only

94    UniProtKB and excludes both long and viral sequences. Consequently, 78% of all UniRef50

95    clusters have members with a predicted structure in AFDB (Extended data Fig. 1b). Of these,

96    29% are functionally dark, a proportion that drops with an increase in predicted model accuracy

97    (Extended data Fig. 1c,d) while retaining a similar proportion of DUFs (Extended data Fig. 1e).

98    Thus, there is a considerable proportion of proteins in UniProt that can not be automatically

99    annotated, but that high confidence structural information can now be leveraged to gain insights

100   about a substantial number of these.

101

**Sequence similarity network of AFDB90**

103   While UniRef50 provides groups of sequences that are overall similar at the sequence level,

104   they do not reach the family and superfamily levels and do not account for local similarities.

105   To reach these levels and put functionally dark clusters into evolutionary context, we

106   constructed a large-scale sequence similarity network of all clusters where structural

107   information can be confidently leveraged to support functional annotations. This corresponds

108   to the 6'136'321 UniRef50 clusters (circa 53 million unique protein sequences) which have

109   structural representatives with an average pLDDT > 90 in AFDB (the AFDB90 dataset).

110   We employed MMseqs2[13] for all-against-all sequence searches (Fig. 1b), connecting two

111   sequences if they have an alignment that covers at least 50% of one of the proteins with E-

112   value < $1 \times 10^{-4}$. The resulting network has over 4 million connected nodes and 10 million edges,

113   which includes 43% of all dark UniRef50 clusters (Fig. 2). Remarkably, 40% of these dark

114   clusters connect to bright UniRef50 clusters, revealing potential evolutionary relationships for

115   over 700'000 unique proteins.

116   The network is composed of 242'876 connected components with at least 2 nodes, with the

117   largest encompassing about 50% of all AFDB90 (Fig. 2a). Of these components, 19% have an

118   average brightness content below 5% ("fully dark") (Fig. 2d). Only 25% of the components are

119   "fully bright" (i.e., average functional brightness >95%). The percentage of UniRef50 clusters

120   in fully dark components decreases with the component's size (Fig. 2b,c), highlighting that the

121   lower the number of homologs the harder a protein is to annotate. Still, and while the

122   distribution is skewed towards smaller sizes in both fully dark and fully bright components

123   (Fig. 2e,f), the largest dark component in our network has over 800 nodes. These fully dark

124   components are fertile ground for novel family discovery, as exemplified by the two new

125   families we describe below.

126

*A new glycosyltransferase family*

128   The largest functionally dark connected component in our set is component 27, with 836

129   UniRef50 clusters (4'889 unique bacterial protein sequences, average brightness 2±13%, Fig.

130   3a). Their representatives have a median length of 665 ± 169 amino acids, most are predicted

131   to be transmembrane, and are annotated as "Uncharacterised YfhO" in InterPro. Indeed, the

132   proteins in this component that are not called "Uncharacterised protein" mostly have the title
133   "YfhO family protein", which corresponds to a family involved in lipoteichoic acid or wall
134   teichoic acid glycosylation[14]. However, the predicted structural model superposes poorly to the
135   YfhO family (TM-score 0.58, Fig. 3b), prompting a more in-depth investigation.
136   HHPred[15] and Foldseek[16] find multiple, medium-to-high confidence matches in the PDB
137   (Probability > 95% and TM-score ~0.6, Fig. 3b), including the eukaryotic Dolichyl-
138   diphosphooligosaccharide-protein glycosyltransferase subunit STT3 and its bacterial homolog
139   oligosaccharyltransferase PglB[17,18], absent from our network because their representatives have
140   an average pLDDT < 90. We collected sequences for all four groups of proteins (YfhO, STT3,
141   PglB, and component 27) and built a sequence similarity network in order to investigate how
142   they may relate at the sequence level (Fig. 3a). This network highlighted that most dark proteins
143   in component 27 cluster separately from the reference YfhO, forming a single YfhO-like
144   protein family that is linked to the STT3/PglB groups by multiple hypothetical proteins, mostly
145   of prokaryotic origin, often annotated as "Glycosyltransferase family 39 protein".
146   These results support the notion that component 27 belongs to the well-studied superfamily of
147   transmembrane oligosaccharyl- and glycosyltransferases, but also indicate that it is a hitherto
148   undescribed bacterial protein family. In this case, inspecting the AlphaFold model revealed
149   possible inconsistencies in their automated annotation, illustrating the added value of structural
150   models to guide sequence-based family classification.
151   

## *A new toxin-antitoxin superfamily*

153   Component 159 is composed of 327 UniRef50 clusters, corresponding to 1'222 unique protein
154   sequences, mostly annotated as "Domain of Unknown Function 6516" (i.e. DUF6516, Fig. 4).
155   These proteins are predicted to adopt a conserved α+β fold, where two α-helices pack against
156   an antiparallel β-sheet with 7 strands (Extended data Fig. 2). Contrary to component 27,
157   HHPred and Foldseek searches found no confident matches in the PDB. A high resolution
158   similarity network unravelled 7 distinct classes of DUF6516-containing proteins (Fig. 4a).
159   Based on the AFDB models, structure-based function predictor DeepFRI[9] proposed that they
160   may bind DNA or other nucleic acids and carry a hypothetical catalytic site with a hydrolase
161   activity over ester bonds (Fig. 4c, Supplementary file 1). Genomic context analysis with
162   GCsnap[19] highlighted that DUF6516-coding genes are commonly found in a conserved two-
163   gene (bicistronic) genomic arrangement, with DUF6516 predominantly located downstream of
164   the conserved bicistronic "partner" (clusters 1, 2, 4 and 6).
165   While most of the "partner" genes associated with DUF6516 code for "hypothetical proteins"
166   of unknown function, one in cluster 1 is a remote homolog of RelB, a well-characterised
167   antitoxin[20]. Indeed, the bicistronic arrangement is typical for toxin-antitoxin (TA) systems[21].
168   When active, the TA toxin proteins abolish bacterial growth, and the control of this toxicity is
169   executed by the antitoxin, which, in the case of "type II TA systems", is a protein that acts by
170   forming an inactive complex with the toxin. DeepFRI predictions for DUF6516 partners
171   suggests they may also bind DNA (Supplementary file 1), an activity characteristic for diverse
172   antitoxins[21], and co-folding prediction with AlphaFold-Multimer generated high confidence
173   models (93 average pLDDT, 0.902 iPTM) that support the interaction between the two proteins
174   as a dimer of dimers (Fig. 4b), as commonly observed for type II TAs. Therefore, we
175   hypothesised that DUF6516 is a novel toxic TA effector that is neutralised either *in trans* by

176  diverse unrelated antitoxins (subclusters 1-4, 6 and 7) or *in cis* by a fused unknown antitoxin

177  domain (UnkD, subcluster 5).

178  To validate the putative TAs experimentally and gain insights into the mechanism of

179  DUF6516-mediated toxicity, we used our established toolbox for TA studies[22]. We targeted

180  TA from six Gammaproteobacterial species for testing in *E. coli* surrogate host, and all the

181  putative toxins dramatically abrogated *E. coli* growth (Fig. 4d) while the putative antitoxins

182  had no effect (Extended data Fig. 3). Neutralisation assays showed full suppression of toxicity

183  when the toxins were co-expressed with cognate antitoxins (Fig. 4d), thus directly validating

184  that these gene pairs are, indeed, *bona fide* TA systems.

185  To probe the mechanism of DUF6516-mediated toxicity, we carried out metabolic labelling

186  assays with $^{35}$S methionine (a proxy for translation), or $^{3}$H uridine (a proxy for transcription)

187  or $^{3}$H thymidine (a proxy for replication). Expression of *Allochromatium tepidum* strain NZ

188  DUF6516 toxin resulted in a decrease in efficiency of $^{35}$S methionine incorporation (Fig. 4e),

189  indicative of  the inhibition of protein synthesis. We hypothesise that the effect could be

190  mediated by the yet-unproven RNase activity of the DUF6516 toxin.

191  We conclude that DUF6516 is a *bona fide* translation-targeting toxic effector of a novel TA

192  family, and propose renaming it TumE (for "dark" in Estonian), with the antitoxin components

193  dubbed as TumA, with A for "antitoxin". This example illustrates the difficulty of automating

194  functional annotation for proteins from completely novel superfamilies. Here, the combination

195  of genomic context information, remote homology searches on genomic neighbours, and deep

196  learning-based structure-guided function prediction helped formulate a testable functional

197  hypothesis.

198

199  **Semantic consistency across the network**

200  Recently, the ProtNLM[10] large language model was implemented as an approach to

201  automatically name proteins in UniProtKB titled as "Uncharacterised protein". Given that

202  language models have the tendency to "hallucinate" predictions when faced with an

203  unknown[23], we hypothesise that such an approach would generate a wide diversity of predicted

204  names for completely novel protein families. To investigate this hypothesis, we compared the

205  diversity of names predicted by the first release of ProtNLM for proteins in fully dark

206  components and those in fully bright.

207  In both cases, the distributions of names and words (collectively referred to as "semantic

208  diversity") were highly skewed towards extremely low diversities, but the fully dark set was

209  significantly different from the fully bright (Kolmogorov–Smirnov two-sided test statistic

210  0.2915, P-value = $8.882 \times 10^{-16}$, Extended data Fig. 4a,b). Most bright components had a low

211  semantic diversity, indicating a coherent and consistent naming. The maximum word diversity

212  in these was 37%, corresponding to cases with variations of the same name (e.g. multiple

213  "Cytotoxins" with different labels for component 100'340). On the other hand, fully dark

214  components tended to have a higher semantic diversity, with a name diversity of 19%

215  (compared to 10% in fully bright) and a word diversity of 7% (compared to 4%). The more

216  consistently named dark components were those with previously submitted names, such as

217  "DUF6516".

218  The dark component with the highest semantic diversity (45%) was component 3'314,

219  composed of 53 proteins with a wide variety of unrelated predicted names, including

220 "Integrase","NADH-quinone oxidoreductase subunit F", "Dynein light chain", "Prophage
221 protein", etc. Despite this, proteins in component 3'314 share a common fold (Extended data
222 Fig. 5a) but FoldSeek found no hits in the PDB. HHPred searches highlighted a small local
223 match to the tubulin-binding domain of *Chlamydomonas reinhardtii* TRAF3-interacting
224 protein 1 (Probability 71%), but when clustered together at sequence-level these two groups of
225 proteins only formed a few weak connections (Extended data Fig. 5a). Though small,
226 component 3'314 is dispersed throughout bacteria and bacteriophages, and the members do not
227 share a conserved genomic context (Extended data Fig. 5b). Together with the presence of
228 prophage-associated protein encoding genes in these genomic contexts, such as "Host-nuclease
229 inhibitor protein Gam"[24], these data support the "Prophage protein" title.

230 Another example with a high semantic diversity (35%), and where structure information aided
231 function assignment, is component 6'732. It consists of 54 entries, some of which are annotated
232 inconsistently as "AbiEi_1 domain-containing protein", "Transposase", "Acyl-CoA
233 dehydrogenase" and "TetR family transcriptional regulator". HHpred searches found no hits in
234 the PDB, but structure-based searches using AFDB models yielded matches to a number of
235 type II restriction endonucleases. The most similar was EndoMS, a mismatch restriction
236 endonuclease[25] that superposes with an RMSD of 2.3-2.6 Å. Within the structural alignment,
237 the most conserved residues are those constituting the EndoMS active site (Extended data Fig.
238 5c), which are invariant in all members of component 6'732. This suggests that they share a
239 similar active site architecture that has a common restriction endonuclease active site motif
240 (E/D)-Xn-(E/D)XK[26,27], and that component 6'732 may represent a new family of putative
241 restriction endonucleases whose precise function is unknown.

242 These results highlight that ProtNLM when presented with families with no homologs was
243 indeed hallucinating a diverse range of names. By setting a word diversity cutoff of >20% for
244 components with >50 proteins, we identified 290 such functionally dark components, covering
245 4'618 UniRef50 clusters and 37'211 unique protein sequences, and are defining Pfam[2] families
246 for each of them (133 new families available in the next Pfam releases 36.0 and 37.0;
247 Supplementary file 2). This includes component 3'314 as the PF21779 family and whose
248 members are now titled DUF6874, and component 6'732, which is now PF22187 and its
249 members named DUF6946.

250 Overall, pooling predictions across the network can help assess the consistency of automated
251 annotation methods, especially in data-driven approaches. As we define new Pfam families,
252 their naming should become consistent as future versions of ProtNLM consume this data.
253 Starting from UniProt release 2023_01, the criteria for displaying ProtNLM names has changed
254 to include an ensemble approach, an increased confidence threshold, and an automatic
255 corroboration pipeline (https://www.uniprot.org/help/ProtNLM), thus many of these
256 hallucinated names have now reverted to "Uncharacterised protein".

257

258 **Structural outliers across the network**
259 Just as semantic diversity revealed novelties in protein sequence space, we also investigated
260 how different the predicted structural characteristics of proteins in our network are from the
261 structures in the PDB. For this, we introduced the concept of "structural outliers" by using an
262 alphabet of substructure representations covering 1'024 local structural contexts (16 residues
263 in sequence and 10Å spatial neighbourhood, Extended data Fig. 6). We trained an outlier

264  detector on PDB structures and predicted that 699'084 AFDB90 structures have substructure
265  compositions that are rare or absent in the PDB, giving us a measure of plausibility that can
266  help prioritise protein family classification.
267  While the examples described in the previous section are all structural inliers, we found that
268  30% of outliers are in dark UniRef50 clusters (Fig. 5a) and that they tend to be shorter and
269  more repetitive than inliers (Fig. 5a,b). Proteins may be structural outliers for a variety of
270  reasons, including novel folds as in the next section. Short outliers typically represent
271  fragments of existing families (Fig. 5c), likely due to frameshift errors introduced during
272  whole-genome sequencing. Long outliers tend to be highly repetitive proteins (6'791 clusters,
273  with >500 residues and shape-mer diversity fraction <0.1, of which 4'948 are bright), which
274  are rare or absent in the PDB (Fig. 5d). Proteins that require conditions to fold that are not
275  modelled by AlphaFold2, such as binding partners (Fig. 5e), sometimes have models in AFDB
276  that do not resemble the single chain of the complex as found in the PDB, i.e the predicted
277  monomeric fold may not always be functionally meaningful.
278  While most fully dark and fully bright components do not contain structural outliers, the outlier
279  content is significantly different between the two sets (Kolmogorov–Smirnov two-sided test
280  statistic 0.0586, P-value = $5.245 \times 10^{-81}$, Extended data Fig. 4c). Fully dark components have on
281  average a higher outlier content (21%) than fully bright (15%), but these only correspond to
282  about half of the structural outliers. Indeed, 44% of outliers are singletons, i.e UniRef50 clusters
283  which do not form a component with at least 2 nodes, giving us a measure to prioritise even
284  these cases for further analysis, as in the example below.
285
286  ### *The β-flower fold*
287  UniRef50_A0A494VZL1 is an example of a structural outlier which is a singleton in the
288  network. It folds as a shallow, symmetric β-barrel with 96 residues, made of 10 short
289  antiparallel β-strands that form a hydrophobic channel. On one side of the β-barrel, the loops
290  connecting each strand are much longer (9 residues) than those on the other side (4 residues),
291  and some are enriched with positively charged arginine and lysine residues with phenylalanines
292  at the tips pointing towards the exterior of the β-barrel (Fig. 5f). Overall, it looks like a flower
293  (Fig. 5g) and hence we named it the "β-flower" fold.
294  Foldseek searches found hits to 43 AFDB90 clusters (TM-score >0.6, most from bacteria)
295  across 13 different components, some of which are bright because they are annotated as "Cell
296  wall-binding protein" or "MORN repeat variant". There are at least three globally different
297  folds (Fig. 5f), differing in the number of strands (8, 10, or 12), with their "petals" comprising
298  β-hairpins that are arranged in four-, five- or six-fold symmetry. Some of the hits resemble half
299  of a flower, perhaps corresponding to fragments of longer domains, and many enclose a C-
300  terminal hydrophobic α-helix. Some β-flowers also contain N-terminal lipoprotein attachment
301  motifs[28,29], suggesting they may be associated with the bacterial inner membrane or transferred
302  to the inner leaflet of the outer membrane.
303  Although no similarity to the PDB was highlighted by Foldseek or HHpred searches, the β-
304  flower folds with six-fold symmetry are reminiscent of the Tubby C-terminal domain[30], which
305  adopts a twelve-stranded β-barrel fold enclosing a hydrophobic α-helix (Fig. 5f,g). Tubby-like
306  proteins either bind to phosphoinositides or function as phospholipid scramblases[30]. β-flowers
307  and Tubby-like proteins share a network of aromatic hydrophobic residues that flank the edges

308 of the β-strands and point toward the interior of the β-barrel, thus engaging in tight contacts
309 with the central hydrophobic helix. Interestingly, the N-terminal strand of Tubby is circularly
310 permuted in β-flowers (Fig. 5g), which leads to a different entry point of the α-helix into the β-
311 barrel channel, and to a difference in its directionality. Additionally, the length of the β-strands
312 and the connecting loops in the β-flower proteins are significantly shorter.
313 Based on their global structural similarity and the presence of a semi-conserved [DNEQ]XXG
314 sequence motif at the tip of the β-hairpin, and the repeat unit of both β-flowers and Tubby-like,
315 the diversity of these proteins has been added to Pfam as the new entries PF21784, PF21785
316 and PF21786, which together with the Tubby C-terminal domain now form the CL0395 clan.
317 This, together with the different types of structural outliers described, highlights that the 3D
318 context provided by the models in AFDB is highly informative for protein analysis efforts and
319 that the structural space covered needs to be put into a coherent evolutionary, functional, and
320 local structural context before any model, even with high predicted accuracy, is used as a
321 reference.
322
323 **Towards large-scale function annotation**
324 In this work, we carried out a large-scale analysis of the UniProt protein sequence space
325 covered by high confidence predicted structural models, as made available through AFDB
326 version 4. In order to aid functional annotation of this space, we constructed an interactive
327 sequence similarity network accounting for about 53 million proteins enriched with predicted
328 name diversity and structural plausibility scores, the first network at such a large scale. We
329 demonstrate that this network is a rich source of putative novel protein folds, families and
330 superfamilies, providing multiple starting points for further downstream studies.
331 We find that many functionally unannotated proteins are remote homologs of annotated ones,
332 relationships which can now be easily explored. Additionally, over 1 million proteins belong
333 to completely unannotated connected components, many of which cannot be named
334 consistently using the most recent deep learning-based approaches or contain proteins with
335 structural features distinct from what is seen in the PDB. When combined with traditional
336 protein evolution approaches, structure-based comparisons, genomic context information,
337 structure-based function prediction, and the conservation of local features such as active sites,
338 we could gather support for common evolutionary origins, gain valuable insights into putative
339 functions and put forward concrete testable hypotheses for experimental characterisation.
340 Indeed, the functional annotation of dark proteins, even from a purely computational
341 perspective, requires a combination of data sources and approaches. It is crucial to combine
342 individual predictions across connections in the network to increase the confidence of any
343 hypothesis. Most of our examples had such support from both sequence and structure, and even
344 for the novel β-flower fold, a singleton in our network, the presence of a semi-conserved
345 sequence motif captured only due to local structural similarities allowed us to generate an initial
346 classification. This information can now help guide further validation experiments, such as
347 those carried out for TumE.
348 Our study has some caveats and limitations, however. All alignments required coverage across
349 the entire protein sequence, while a domain-based exploration would provide a possible
350 complementary solution. Our functional brightness definition excluded predicted intrinsically
351 disordered and coiled-coil proteins, and misclassifies some functionally uncharacterised

proteins as bright due to ambiguous annotations (e.g "transmembrane" or "repeat"), or characterised ones as dark due to "Putative" annotations. Furthermore, we focus only on proteins with high confidence predicted structures from AFDB, setting aside the wealth of potential darkness in metagenomic data for which structural models are also now available through the ESM Metagenomic Atlas[31]. Though we could already highlight a significant proportion of novelty, in-depth exploration combining multiple sources of evidence could only be carried out for a small number of families and folds. Thus, the examples we discuss are the low-hanging fruit of uncharacterised or unannotated protein families, and they are only the tip of the iceberg.

Similarity networks are a common representation of protein space[32,33] and recent approaches to categorise protein diversity and uncover novelties have showcased the importance of incorporating multiple perspectives and methods in protein annotation[31,34–36]. Our work combines these concepts by providing the first annotated similarity network model of protein sequence space at such a large scale, which we make available as an interactive and accessible web resource. We anticipate that further advances in deep learning-based methods for function prediction[9], remote homology detection[37,38] and protein structure prediction[31] will allow for analyses on an even larger scale, incorporating more diverse data sources with greater confidence. As such advances continue, we as a community are closer than ever to harnessing the full potential of the protein universe, from unknown biology to new biomedical, pharmaceutical and biotechnological applications.

**Main text references**

1.  Varadi, M. *et al.* AlphaFold Protein Structure Database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Res.* **50**, D439–D444 (2022).
2.  Mistry, J. *et al.* Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412–D419 (2020).
3.  UniProt Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Res.* **51**, D523–D531 (2023).
4.  Richardson, L. *et al.* MGnify: the microbiome sequence data analysis resource in 2023. *Nucleic Acids Res.* **51**, D753–D759 (2023).
5.  Boutet, E. *et al.* UniProtKB/Swiss-Prot, the Manually Annotated Section of the UniProt KnowledgeBase: How to Use the Entry View. *Methods Mol. Biol.* **1374**, 23–54 (2016).
6.  Paysan-Lafosse, T. *et al.* InterPro in 2022. *Nucleic Acids Res.* **51**, D418–D427 (2023).
7.  Levitt, M. Nature of the protein universe. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 11079–11084 (2009).
8.  Bienert, S. *et al.* The SWISS-MODEL Repository-new features and functionality. *Nucleic Acids Res.* **45**, D313–D319 (2017).
9.  Gligorijević, V. *et al.* Structure-based protein function prediction using graph convolutional networks. *Nat. Commun.* **12**, 1–14 (2021).
10. Gane, A. *et al.* ProtNLM: Model-based Natural Language Protein Annotation. (2022).
11. Suzek, B. E. *et al.* UniRef clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics* **31**, 926–932 (2015).
12. Leinonen, R. *et al.* UniProt archive. *Bioinformatics* **20**, 3236–3237 (2004).
13. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat. Biotechnol.* **35**, 1026–1028 (2017).

398   14.  Rismondo, J., Percy, M. G. & Gründling, A. Discovery of genes required for
399        lipoteichoic acid glycosylation predicts two distinct mechanisms for wall teichoic acid
400        glycosylation. *J. Biol. Chem.* **293**, 3293–3306 (2018).
401   15.  Söding, J. Protein homology detection by HMM-HMM comparison. *Bioinformatics* **21**,
402        951–960 (2005).
403   16.  van Kempen, M. *et al.* Fast and accurate protein structure search with Foldseek. *Nat.*
404        *Biotechnol.* (2023) doi:10.1038/s41587-023-01773-0.
405   17.  Kelleher, D. J. & Gilmore, R. An evolving view of the eukaryotic
406        oligosaccharyltransferase. *Glycobiology* **16**, 47R–62R (2006).
407   18.  Szymanski, C. M. & Wren, B. W. Protein glycosylation in bacterial mucosal pathogens.
408        *Nature Reviews Microbiology* vol. 3 225–237 Preprint at
409        https://doi.org/10.1038/nrmicro1100 (2005).
410   19.  Pereira, J. GCsnap: Interactive Snapshots for the Comparison of Protein-Coding
411        Genomic Contexts. *J. Mol. Biol.* **433**, 166943 (2021).
412   20.  Gotfredsen, M. & Gerdes, K. The Escherichia coli relBE genes belong to a new toxin-
413        antitoxin gene family. *Mol. Microbiol.* **29**, 1065–1076 (1998).
414   21.  Jurėnas, D., Fraikin, N., Goormaghtigh, F. & Van Melderen, L. Biology and evolution
415        of bacterial toxin-antitoxin systems. *Nat. Rev. Microbiol.* **20**, 335–350 (2022).
416   22.  Kurata, T. *et al.* A hyperpromiscuous antitoxin protein domain for the neutralization of
417        diverse toxin domains. *Proc. Natl. Acad. Sci. U. S. A.* **119**, (2022).
418   23.  Ziwei Ji Hong Kong University of Science and Technology, Hong Kong *et al.* Survey of
419        Hallucination in Natural Language Generation. *ACM Computing Surveys* (2023)
420        doi:10.1145/3571730.
421   24.  Akroyd, J. E., Clayson, E. & Higgins, N. P. Purification of the gam gene-product of
422        bacteriophage Mu and determination of the nucleotide sequence of the gam gene.
423        *Nucleic Acids Res.* **14**, 6901–6914 (1986).
424   25.  Nakae, S. *et al.* Structure of the EndoMS-DNA Complex as Mismatch Restriction
425        Endonuclease. *Structure* **24**, 1960–1971 (2016).
426   26.  Aggarwal, A. K. Structure and function of restriction endonucleases. *Curr. Opin. Struct.*
427        *Biol.* **5**, 11–19 (1995).
428   27.  Pingoud, A. & Jeltsch, A. Structure and function of type II restriction endonucleases.
429        *Nucleic Acids Res.* **29**, 3705–3727 (2001).
430   28.  Klein, P., Somorjai, R. L. & Lau, P. C. Distinctive properties of signal sequences from
431        bacterial lipoproteins. *Protein Eng.* **2**, 15–20 (1988).
432   29.  Hayashi, S. & Wu, H. C. Lipoproteins in bacteria. *J. Bioenerg. Biomembr.* **22**, 451–471
433        (1990).
434   30.  Bateman, A. *et al.* Phospholipid scramblases and Tubby-like proteins belong to a new
435        superfamily of membrane tethered transcription factors. *Bioinformatics* **25**, 159–162
436        (2009).
437   31.  Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a
438        language model. *Science* **379**, 1123–1130 (2023).
439   32.  Nepomnyachiy, S., Ben-Tal, N. & Kolodny, R. Global view of the protein universe.
440        *Proceedings of the National Academy of Sciences* **111**, 11691–11696 (2014).
441   33.  Alva, V., Remmert, M., Biegert, A., Lupas, A. N. & Söding, J. A galaxy of folds.
442        *Protein Sci.* **19**, 124–130 (2010).
443   34.  Bordin, N. *et al.* AlphaFold2 reveals commonalities and novelties in protein structure
444        space for 21 model organisms. *Commun Biol* **6**, 160 (2023).
445   35.  Akdel, M. *et al.* A structural biology community assessment of AlphaFold2 applications.
446        *Nat. Struct. Mol. Biol.* **29**, 1056–1067 (2022).
447   36.  Barrio-Hernandez, I. *et al.* Clustering predicted structures at the scale of the known

448  protein universe. *bioRxiv* 2023.03.09.531927 (2023) doi:10.1101/2023.03.09.531927.
449  37.  Kaminski, K., Ludwiczak, J., Alva, V. & Dunin-Horkawicz, S. pLM-BLAST – distant
450  homology detection based on direct comparison of sequence representations from
451  protein language models. Preprint at https://doi.org/10.1101/2022.11.24.517862.
452  38.  Pantolini, L., Studer, G., Pereira, J., Durairaj, J. & Schwede, T. Embedding-based
453  alignment: combining protein language models and alignment approaches to detect
454  structural similarities in the twilight-zone. *bioRxiv* 2022.12.13.520313 (2022)
455  doi:10.1101/2022.12.13.520313.
456  39.  Lomize, A. L., Todd, S. C. & Pogozheva, I. D. Spatial arrangement of proteins in planar
457  and curved membranes by PPM 3.0. *Protein Sci.* **31**, 209–220 (2022).
458  40.  Berisio, R. & Delogu, G. PGRS domain structures: Doomed to sail the mycomembrane.
459  *PLoS Pathog.* **18**, e1010760 (2022).
460

461  **Figure legends**

462  **Figure 1. General workflow for the collection, classification and mapping of functionally**
463  **dark proteins in UniProt and AlphaFold database.** (a) Starting from the clusters in
464  UniRef50, we collected all the functional annotations for all included UniProtKB and UniParc
465  entries, including coiled coil and intrinsically disordered (IDPs) predictions and excluding all
466  of those with "Putative", "Hypothetical", "Uncharacterised" and "DUF" in their names. We
467  selected the protein with the highest full-length annotation coverage (i.e., brightness) as the
468  functional representative of each cluster. (b) From the collected UniRef50 clusters, we selected
469  those with a structural representative with pLDDT >90 in the AlphaFold database v4, and
470  constructed a large-scale sequence similarity network by all-against-all MMseqs2 searches,
471  representing the sequence landscape of more than 6 million UniRef50 clusters.
472

473  **Figure 2. Large-scale sequence similarity network for over 6 million UniRef50 cluster**
474  **representatives with high predicted accuracy models in AFDB (AFDB90).** (a) Layout of
475  the resulting network, as computed with Cosmograph (https://cosmograph.app/). The network
476  contained 4'270'404 nodes connected by 10'339'158 edges, reduced for simplicity to a set of
477  688'852 communities connected by a total of 1'488'764 edges (see *Methods Section Large-*
478  *scale Sequence Similarity Network* for details). The 1'865'917 UniRef50 clusters that did not
479  connect to any other in the MMseqs2 searches were excluded. Only the 473'612 communities
480  that have at least one inbound or outbound edge (degree of 1) are displayed in the figure. Nodes
481  are coloured by the average functional brightness of the UniRef50 clusters included in the
482  corresponding community. An interactive version is available at
483  https://uniprot3d.org/atlas/AFDB90v4. (b) Histograms of functional brightness content for
484  connected components with more than 50'000 and with only 5 to 2 nodes (UniRef50 clusters),
485  highlighting their different darkness content. (c) Scatter plot of the component size (i.e. number
486  of UniRef50 clusters) cut-off and the percentage of functionally dark UniRef50 clusters. (d)
487  Histogram of the average brightness per component. Size distribution for (e) fully dark
488  connected components (average brightness <5%) and (f) fully bright connected components
489  (average brightness >95%).
490

491 **Figure 3. Connected component 27 is a new family in a well-studied superfamily of**
492 **transmembrane glycosyltransferases.** (a) High resolution sequence similarity network for
493 7'004 homologs of the sequences in component 27, computed with CLANS at an E-value
494 threshold of $1x10^{-20}$. Points represent individual proteins and grey lines BLASTp matches at
495 an E-value better than $1x10^{-20}$. Individual clusters are coloured and labelled accordingly to their
496 representative members. Only YfhO-like and STT3/PglB sequences are highlighted, with grey
497 dots depicting other homologous groups. AglB corresponds to the PglB/STT3-like sequences
498 from archaea. Black dots depict those sequences that make component 27 in our network, and
499 white dots mark those that are bright. (b) Predicted structural models as in AFDBv4 for the
500 representative of component 27 (C27, UniProt ID A0A7X7MB17), and YfhO (UniProt ID
501 YFHO_BACSU), and experimental structures of the PglB (PDB ID 6GXC, chain A) and STT3
502 (PDB ID 7OCI, chain F) cluster representatives. Models are coloured according to the colour
503 of their corresponding cluster in (a). The membrane regions, as predicted with PPM 3.0
504 server[39], are marked by dashed lines.
505
506 **Figure 4. Connected component 159 is a novel toxin in the hitherto undescribed toxin-**
507 **antitoxin superfamily TumE-TumA. (**a) High resolution sequence similarity network for
508 2'453 homologs of the sequences in component 159, computed with CLANS (E-value $1x10^{-10}$
509 $^{10}$). Points represent proteins and grey lines BLASTp matches (E-value $<1x10^{-4}$). Individual
510 subclusters are labelled 1-7, and subclusters a-c. The consensus genomic contexts, as identified
511 by GCsnap, are displayed with different flanking families coloured from blue to red. (b) 3D
512 model of the complex between the putative toxin and antitoxin from *Allochromatium tepidum*
513 strain NZ, modelled with AlphaFold-Multimer, highlighting the regions where DNA is
514 predicted to interact with the antitoxin. (c) Structural model of *A. tepidum* TumE/DUF6516
515 toxin (EntrezID WP_213381069.1) coloured according to the two most frequent molecular
516 functions predicted for 100 homologs with DeepFRI. Residues responsible for the predictions
517 are highlighted in red. The percentage reflects the frequency of the highlighted prediction. (d)
518 Validation of *tumE-tumA*. Plasmids for expression of putative toxins (pBAD33 derivates) were
519 co-transformed into *E. coli* BW25113 cells with antitoxin expression plasmids or the empty
520 pMG25 vector. Bacteria were grown for five hours in liquid LB media supplemented with
521 appropriate antibiotics and 0.2% glucose. The cultures were normalised to $OD_{600} = 1.0$, serially
522 diluted and spotted on LB plates containing appropriate antibiotics and 0.2% arabinose for
523 toxin induction and 500 µM IPTG for antitoxin induction. The plates were scored after an
524 overnight incubation at 37 °C. For source data, see Supplementary figure 1. (e) Metabolic
525 labelling assays with *E. coli* BW25113 expressing *A. tepidum* TumE/DUF6516 toxin. Error
526 bars indicate the standard error (SE) of the arithmetic mean. All experiments shown on (d) and
527 (e) were performed as n=3 biologically independent replicates (individual independent
528 cultures). All repetitions of the experiments shown on (d) yielded similar results.
529
530 **Figure 5. Structural outliers can represent fragments, repetitive proteins, proteins**
531 **requiring folding conditions out of the scope of AlphaFold2, or novel folds.** (a-b)
532 Distribution of brightness, shape-mer diversity and length of the (a) structural outliers and (b)
533 the same number of structural inliers with the most positive outlier scores. Shape-mer diversity
534 is defined as the number of unique shape-mers by the length of the protein. (c) An AFDB model

12

535  of "TonB-dependent receptor-like" protein that is a fragment of the β-barrel domain. Over
536  16'500 proteins across 1'258 components have this annotation, of which 86% are fully bright.
537  From these, 82% have less than the required number of β-sheet shape-mers, despite 55% not
538  being explicitly annotated as fragments in UniProtKB. (d) Two long repetitive outliers, one
539  belonging to the PE-PGRS superfamily (G0TGH8), thought to be novel folds and found widely
540  in mycobacteria[40], and one to the Tetratricopeptide-like helical domain superfamily
541  (A0A015IZK3) where the median PDB structure length of structures with resolution < 3Å is
542  only 370. (e) AFDB model annotated as containing "Putative type VI secretion system, Rhs
543  element associated Vgr domain" (A0A377W562), a trimeric PDB structure (PDB ID 6SK0)
544  also containing this domain, and an AlphaFold-Multimer model of the A0A377W562 trimer
545  which has 1.1Å RMSD to the PDB structure. The AFDB model does not resemble the PDB
546  structure because these proteins form obligate complexes and adopt a trimeric β-solenoid fold.
547  (f) AlphaFold models of different variations of the β-flower, with positively charged residues
548  in red and phenylalanine in green for A0A494VZL1, and PDB structures of the human Tubby
549  C-terminal domain (PDB ID 2FIM). Black arrows indicate the circularly permuted loop in
550  A0A0S7BXY3 and PDB ID 1ZXU. (g) AlphaFold model of A0A0S7BXY3 and PDB structure
551  of *Arabidopsis thaliana* putative phospholipid scramblase (PDB ID 1ZXU). Black arrows
552  indicate the circularly permuted loop.
553
554  **Methods**
555
556  **Data collection**
557  We started from the 53'625'855 UniRef50[11] clusters as of August 2022 (UniRef version
558  2022_03) and the 214'683'829 structural models for most UniProtKB entries available via the
559  AlphaFold database (version 4, AFDBv4). For each Swiss-Prot[5], TrEMBL[3] and UniParc[12]
560  entry in each UniRef50 cluster we collected their sequence, taxonomy and functional and
561  structural annotations from UniProt and InterPro[6] using custom Python 3.6 code. Redundant,
562  overlapping annotations were continuously merged (Fig. 1a), selecting as the preferential name
563  the first occurrence that did not include "Putative", "Hypothetical", "Uncharacterised" and
564  "DUF". Each entry in AFDBv4 was mapped to their UniRef50 cluster, selecting as the
565  structural representative the longest protein with an average pLDDT[41] > 70.
566
567  **Darkness estimation**
568  We define functional brightness of a given protein as the full-length coverage with annotations
569  of its close homologs, with 0% meaning "dark" and 100% meaning "bright". We first computed
570  the full-length coverage with annotations for all entries in all UniRef50 clusters, and considered
571  a cluster as "bright" as the "brightest" sequence it encompasses (Fig. 1a). Annotations
572  considered were: domains annotated in InterPro, and families, predicted disorder and predicted
573  coiled coil regions annotated in UniProtKB and UniParc. All those with "Putative",
574  "Hypothetical", "Uncharacterised" and "DUF" in their name were given a coverage of 0.
575  Pearson correlation was computed using SciPy (v1.5.4).
576
577  **Large-scale sequence similarity network**

578 To model the sequence landscape covered by all UniRef50 clusters with a high confidence
579 structural model, we built a large-scale sequence similarity network of 6'136'321 clusters
580 having a structural representative with pLDDT > 90 (AFDB90 dataset). All-against-all
581 MMseqs2[13] (release 13-45111) comparisons were carried out with the UniRef50 cluster
582 representatives of all selected clusters, connecting two sequences if they have a match that
583 covers at least 50% of their full length sequences with an E-value better than $10^{-4}$. Each edge
584 was given a weight proportional to the E-value of the match, and a maximum of 4 outbound
585 edges were considered per node (Fig. 1b). The direction of the edges was not further
586 considered.

587 To visualise the graph, each connected component was simplified to a set of connected
588 communities, detected using the asynchronous label propagation algorithm, as implemented in
589 the *asyn_lpa_communities* method in networkx (v2.5.1)[42]. This reduced the graph to a total of
590 688'852 communities (hereafter referred to as the AFDB90Communities set) connected by
591 1'488'764 edges, whose layout could then be computed with Cosmograph
592 (https://cosmograph.app/) with the following settings: maximum space allowed = 8192,
593 gravity = 0.5, repulsion = 1.4, repulsion theta = 1.71, link strength = 2, minimum link distance
594 = 1, friction = 1. For each community, we collected the longest and median-length
595 representatives, whose structures were used in our analyses. Individual connected components
596 were visualised in figures with Datashader (v0.12.1, https://datashader.org/index.html).

597 The interactive, annotated and searchable web version of this network was created using the
598 Cosmograph library (https://github.com/cosmograph-org/cosmos, v1.3.0) for network
599 visualisation and the Mol* toolkit (v3.35.0) [43] for 3D macromolecular visualisation of
600 individual structure representatives. Sequence searches over the interactive network are carried
601 out with a simple *k*-mer search to rapidly identify close homologues in the AFDB (>70%
602 sequence identity) and structure searches with Foldseek (3Di method[16], E-value better than $10^{-1}$
603 ) through its API over the AFDBv4 database filtered to 50% sequence identity (UniProt50).
604 Returned matches are mapped back to their corresponding communities.

605

606 **Sequence-based prioritisation of dark connected components and their semantic name**
607 **diversity**
608 Each node in a connected component was attributed a functional brightness value, and
609 components were sorted by their average brightness and their overall size (i.e., number of
610 nodes), so that the top ranking were the largest and darkest. To analyse UniProt name diversity,
611 we extracted names as of UniProt version 2022_04 (December 2022, which includes the initial
612 release of ProtNLM[10] predictions) for all UniRef100 representatives included in clusters of
613 fully dark (average functional brightness $\leq$ 5%) and fully bright (average functional brightness
614 $\geq$ 95%) connected components with at least 50 unique protein sequences. We computed the
615 proportion of unique names (i.e., name diversity) as well as the proportion of unique words
616 (i.e., word diversity), in order to account for small variations of the same name. Kolmogorov–
617 Smirnov statistical test (two-sided) was computed using SciPy (v1.5.4).

618

619 **Protein substructure decomposition**
620 To represent and analyse 3D substructure composition, we built upon Geometricus (v0.5.0,
621 Python 3.9)[44], and use 16 rotation invariant moments[45–47] and one chiral invariant moment[48].

622 These moments were calculated on α-carbon coordinates for overlapping *k*-mers of size 8 and
623 16, and overlapping spheres of radii 5Å and 10Å; for a total of 68 moments for each central
624 residue in a protein, using ProDy (v2.2.0). We trained a neural network using PyTorch
625 (v1.12.0)[49] with these 68 moments as input, 2 linear hidden layers of size 32, a sigmoid output
626 layer of size 10, and with contrastive loss to reduce the output distance between equivalent
627 pairs of central residues and increase the distance between non-equivalent pairs in a training
628 set. The output of the network for each residue, 10 floating point numbers between 0 and 1,
629 was discretized into 10 bits based on whether the value was greater than or less than 0.5,
630 resulting in 1024 shape-mers.
631 The training set was created from structures from the CATH database (v4.2.0) having less than
632 40% sequence identity (CATH40) that could be assigned to a CATH functional family
633 (FunFam[50]) with an E-value better than $1\text{x}10^{-6}$. From these 8'333 structures, US-align (version
634 20220924)[51] was used to align and superpose all pairs within each FunFam cluster and three
635 randomly chosen pairs for each protein across clusters. Aligned pairs of residues from two
636 same FunFam proteins with TM-score > 0.8 were considered as positive pairs. Aligned or
637 random pairs of residues from two proteins belonging to different CATH superfamilies, with
638 TM-score < 0.6 were considered as negative pairs. In addition, using all 31,883 CATH40
639 proteins, we sampled up to 50 pairs of central residues from each protein, where positive pairs
640 had <2 sequence distance and negative pairs had 5-20 sequence distance. In total, this resulted
641 in 6 million residue pairs for training, of which 42% were positive pairs. This dataset could be
642 used for training and/or refining any kind of residue-level contrastive learning task. Training
643 took 30 mins on 1 RTX-3080TI with the ADAM optimizer, a batch size of 1024, and a learning
644 rate of $10^{-3}$ over 5 epochs.
645 Shape-mers were calculated for ProteinNet CASP12 proteins in the 100% sequence identity
646 set[52] with over 20 amino acids. Extended data Fig. 6 shows an example protein with its 6 most
647 common shape-mers highlighted. We trained a FastText model[53] on the shape-mer bit
648 representations using Gensim[54] (v4.2.0, window size of 16, embedding size of 1024). Extended
649 data Fig. 7a shows the sensitivity of SCOPe family retrieval on the SCOPe40 dataset of 11'211
650 structures for all-vs-all Smith-Waterman alignment with FastText shape-mer similarities used
651 as the score matrix (runtime: 12 mins on 10 threads). Shape-mer FastText alignment scores are
652 compared to three structure aligners, Dali[55], Foldseek[16], and TM-align[56]; one sequence aligner,
653 MMseqs2[13]; and 2 other structure alphabet-based structural sequence aligners, 3D-BLAST[57]
654 and CLE-SW[58], using the scripts and benchmark data provided in van Kempen *et al.*[16]. Protein-
655 level embeddings are obtained by averaging across normalised FastText embeddings using the
656 *get_sentence_vector* function. Extended data Fig. 7b shows the distributions of cosine distances
657 of these embeddings within the same SCOPe family and across SCOPe folds.
658
659 **Structural outlier detection**
660 The benchmarking and comparison results (Extended data Fig. 7) demonstrate that the learned
661 structural alphabet and FastText similarities still have discriminative power in distinguishing
662 protein families, despite being much less "local" than approaches such as Foldseek and TM-
663 align which work on individual coordinates of up to 2 residues. We don't explore further
664 alignment optimization, such as compositional bias correction or penalty optimization to
665 increase sensitivity, as more local structural aligners will still have the advantage of higher

666 resolution alignment. However, for the task at hand, our substructure representations give us a
667 good compromise - a discriminative structural alphabet for representing a protein structure as
668 a structural sequence; and substructure decomposition at the level of whole secondary-
669 structural elements, allowing for a broader exploration of substructure composition across the
670 AlphaFold database.

671 For this, we trained the Isolation Forest outlier detection algorithm[59] as implemented in scikit-
672 learn (v1.1.1)[60] on the ProteinNet CASP12 FastText sentence embeddings with 1%
673 contamination rate. Shape-mers for all AFDB90 structural representative AlphaFold models
674 were calculated following the approach described in the analysis of AFDBv1[35] to split each
675 protein into segments with Gaussian smoothed plDDT > 70, after first splitting into domains
676 based on a combination of pLDDT and the predicted aligned error (PAE) matrix, and
677 concatenating shape-mers across each segment in each domain. A shape-mer diversity fraction
678 was defined for each protein as the number of unique shape-mers divided by the total number
679 of residues for which shape-mers are calculated. The trained outlier detection model was used
680 to predict structural outlier scores for AFDB90 proteins. Proteins with negative scores are
681 labelled as outliers. Kolmogorov–Smirnov statistical test (two-sided) was computed using
682 SciPy (v1.5.4).

683

684 **Computational investigation of selected examples**
685 For the analysis of all examples, we combined data from the sequence-based network and its
686 functional brightness annotations, as well as from structural searches with Foldseek and the
687 outlier scores. Structural homologs for selected representatives (those with a length close to the
688 median length in the component) in the PDB or the AFDB90Communities set were searched
689 with Foldseek (v7.04e0ec8) using the TM-align mode[16]. Remote sequence homologs were
690 detected for selected representatives by HHPred searches over the PDB, ECOD and Pfam
691 databases through the MPI Bioinformatics toolkit using default settings[61,62]. AlphaFold-
692 Multimer[63] version 3 was used for protein complex prediction when required, with default
693 settings and relaxation, and the model with the best predicted TM score (pTM) and interface
694 pTM score was selected. PyMol (v2.5.0) was used to visualise selected examples. Further case-
695 by-case analyses were carried out as below.

696 *Component 27*
697 All UniRef100 representatives represented by the nodes of connected component 27 were
698 collected and filtered to a maximum sequence identity of 50% with MMseqs2. The reduced set
699 of sequences was aligned with MUSCLE[64] (v5.1) and the resulting MSA used as input for three
700 independent BLASTp[65] searches over the eukaryotic, archaea and bacterial sequences in *nr*
701 filtered to 70% sequence identity (nr_euk70, nr_arc70, nr_bac70) through the MPI-
702 Bioinformatics toolkit as of January 2023. The same BLAST searches were carried out for
703 Swiss-Prot representatives of the PglB, STT3 and YfhO families (UniProt IDs PGLB_CAMJR,
704 STT3_YEAST and YFHO_BACSU). The full-length sequences matched in all searches were
705 then combined with those representatives of connected component 27 and filtered to a
706 maximum sequence identity of 30% with MMseqs2. The resulting set of 7'004 sequences was
707 clustered based on BLASTp all-against-all searches with CLANS[66] at an E-value of $1\times10^{-20}$
708 until equilibrium.

709

### *Component 159*

710 **Component 159**

711 Ninety-four randomly selected sequences from component 159 were aligned with MUSCLE.
712 The resulting alignment was used for three independent PSI-BLAST[65] searches over the
713 eukaryotic, archaea and bacterial sequences in *nr* (nr_euk, nr_arc, nr_bac) with 8 rounds
714 through the MPI-Bioinformatics toolkit as of October 2022[61,62]. All collected sequences were
715 filtered to a maximum sequence identity of 95% with MMseqs2 and clustered based on
716 BLASTp all-against-all pairwise searches with CLANS until equilibrium at an E-value of $1 \times 10^{-10}$.
717

718 The resulting sequence similarity network was used as input for GCsnap (v1.0.17)[19] for the
719 analysis of the conservation of the genomic contexts encoding for each of the proteins in the
720 individual clusters. A window of four flanking genes was used, MMseqs2 was employed for
721 protein family clustering at an E-value better than $1 \times 10^{-4}$ and clusters of similar genomic
722 contexts were detected using the *operon_cluster_advanced* method, which employs PaCMAP
723 (v0.7.0)[67] to project genomic contexts in 2D based on their family composition and DBSCAN[68]
724 (as implemented in scikit-learn v1.2.2) to identify clusters of similar genomic contexts. Only
725 families that were found in at least 30% of all genomic contexts were considered. For each
726 cluster in the sequence similarity network and each identified neighbour family, up to 100
727 structure representatives were selected from AFDBv4 and used as input to DeepFRI (v1.0.0)[9]
728 with default settings. The top 10 most common predictions per cluster/context family were
729 retrieved. The highest average scoring and most frequently predicted molecular functions were
730 considered the most likely for each case.

731 We generated the 3D structure of a tetramer consisting of two chains of the *Allochromatium*
732 *tepidum* TumE toxin (EntrezID: WP_213381069.1) and two of its putative, cognate TumA
733 antitoxin (EntrezID: WP_213381068.1) using AlphaFold-Multimer.

734

### *Component 3314*

735 **Component 3314**

736 All non-redundant protein sequences represented by the nodes of connected component 3314
737 were collected and filtered as for component 27, but over *nr* filtered to 90% sequence identity
738 (nr_euk90, nr_arc90, nr_bac90, nr_vir90). The same BLAST searches were carried out for the
739 tubulin-binding domain of *Chlamydomonas reinhardtii* TRAF3-interacting protein 1 (UniProt
740 ID A8JBY2_CHLRE, residues 1-131). The full-length sequences matching component 3314
741 homologs and the local sequence matching the TRAF3-interacting protein 1 tubulin binding
742 domain were then combined with representatives of component 3314 and filtered to a
743 maximum sequence identity of 90% with MMseqs2. The resulting set of 890 sequences was
744 clustered based on BLASTp all-against-all searches with CLANS at an E-value of $1 \times 10^{-5}$ until
745 equilibrium. The 141 sequences making subcluster 1 in the resulting network, which included
746 the component 3314-like proteins, were extracted, filtered to a maximum sequence identity of
747 50% with MMseqs2 and used as input for GCsnap (v1.0.17), where a window of four flanking
748 genes was used and MMseqs2 employed for protein family clustering at an E-value better than
749 $1 \times 10^{-4}$.

750

### *Component 6732*

751 **Component 6732**

752   We have built the Pfam family PF22187 (named DUF6946) using component 6732 sequences
753   and iteratively searching for homologs using HMMER (v3.3)[69]. Selected members of this Pfam
754   family were subjected to HHpred searches (HHblits[70] against UniRef30, 3 iterations with cutoff
755   for inclusion $1x10^{-3}$ for multiple alignment generation and PDB70 search database). Foldseek
756   and Dali server (DaliLite v.5)[55] were subsequently used for structure similarity searches, using
757   AFDB models as queries. The obtained structural alignments were manually inspected and
758   compared with the Pfam family alignment. PF22187 was assigned to clan CL0236 that includes
759   diverse families of nucleases.
760
761   *β-flower fold*
762   We constructed three new Pfam families to cover the sequence space of β-flower proteins. To
763   do this we selected example proteins with 4,5 and 6-fold rotational symmetry and iteratively
764   searched for homologs using HMMER's hmmsearch. In general, we used an inclusion
765   threshold of 27 bits, but manually lowered the threshold to identify more homologs or raised it
766   to exclude false matches as identified by AlphaFold2 models. These three families were added
767   to Pfam with accession numbers: PF21784, PF21785 and PF21786 and Pfam clan CL0395,
768   which includes the Tubby C-terminal domain.
769

**770 Experimental validation and characterisation of a predicted toxin-antitoxin family**
**771 (component 159)**
772   Six Proteobacteria TumE examples from subcluster 1a in the CLANS sequence similarity
773   network produced for component 159. and their cognate TumA antitoxins were selected for
774   experimental characterization (Supplementary file 3). The plasmids were constructed using the
775   Circular Polymerase Extension Cloning (CPEC)[71] approach with synthetic DNA procured from
776   Integrated DNA Technologies. ORFs were synthesised with added strong Shine-Dalgarno
777   sequence (AGGAGGAATTAA) and flanking sequences overlapping with multicloning sites
778   of pBAD33[72] (toxin genes) or pMG25[73] (antitoxin genes). The DNA fragments were amplified
779   with Phusion polymerase (Thermo Scientific™) using pBAD_SD_TOX_fwd and
780   pBAD_TOX_MCS_rev or pMG25_insert_fwd and pMG25_insert_rev primer pairs. pBAD33
781   was linearized using primers pBAD_lin_1 and pBAD_lin_2 and pMG25 was linearized using
782   pMG25_lin_from_BlpI and pMG25_lin_from_HindIII. CPEC with Phusion polymerase
783   (Thermo Scientific™) was performed to clone the genes into the vector backbone (25 cycles
784   with 5 min 30 s extension). The CPEC reaction mixture was transformed into DH5α *E. coli*
785   cells and colony PCR with HOT FIREPol® Blend Master Mix (Solis Biodyne) was used to
786   identify colonies with correctly sized inserts. Plasmids were extracted from the overnight
787   cultures using FavorPrepTM Plasmid Extraction Mini Kit (Favorgen) and sequenced. The
788   cognate antitoxin plasmid or empty pMG25 was co-transformed with the toxin plasmids into
789   BW25113 *E. coli* cells. DNA fragments and DNA oligonucleotides used for plasmid
790   construction are provided in Supplementary file 3.
791   Validation of toxicity and metabolic labelling experiments with $^{35}$S methionine, $^{3}$H uridine and
792   $^{3}$H thymidine were performed as described earlier by Kurata *et al.*[22]. Briefly, *E. coli* BW25113
793   strains were transformed with a plasmid pair that allowed for controllable co-expression of
794   putative TumE toxins (pBAD33 derivatives, the toxin is expressed under the control of L-
795   arabinose-inducible P$_{BAD}$ promotor) and TumA antitoxins (pMG25 derivatives[73], IPTG-

796    inducible expression of the antitoxin is driven by $P_{Tac}$ promotor) and pregrown in liquid
797    Lysogeny broth (LB) medium (Lennox) supplemented with 100 μg/mL carbenicillin
798    (AppliChem) and 25 μg/mL chloramphenicol (AppliChem) as well as 0.2% glucose (for
799    repression of toxin expression). Serial 10-fold 5 µL dilutions were spotted on LB plates
800    supplemented with antibiotics (carbenicillin and chloramphenicol) as well as either 0.2%
801    glucose (repressive conditions) or 0.2% arabinose and 1 mM IPTG (induction conditions).
802    Plates were scored after an overnight incubation at 37 °C.
803    For metabolic labelling experiments with TumE toxins, *E. coli* BW25113 strains co-
804    transformed with pBAD33 derivatives (for L-arabinose-inducible expression of toxins) as well
805    as the empty pMG25 vector were first plated out on LB plates supplemented with 100 μg/ml
806    carbenicillin, 25 μg/ml chloramphenicol and 0.2% glucose (to suppress the leaky expression
807    of the toxin). Using fresh, individual *E. coli* colonies for inoculation, 2 mL liquid cultures were
808    prepared in defined Neidhardt MOPS minimal media[74] supplemented with 100 μg/ml
809    carbenicillin, 25 μg/ml chloramphenicol, 0.1% of casamino acids, and 0.2% glucose, and
810    grown overnight at 37 °C with shaking. Next, experimental 15-mL cultures were prepared in
811    125 mL conical flasks in MOPS medium supplemented with 0.5% glycerol, 100 μg/ml
812    carbenicillin, 25 μg/ml chloramphenicol as well as a set of 19 amino acids (lacking
813    methionine), each at final concentration of 25 μg/mL. These cultures were inoculated overnight
814    to final $OD_{600}$ of 0.05, and grown at 37 °C with shaking up to of $OD_{600}$ 0.2. At this point, one
815    1-mL aliquot (the pre-induction zero time-point) was transferred to 1.5 mL Eppendorf tubes
816    containing 10 µL of radioisotope – either $^{35}S$ methionine (4.35 µCi, Perkin Elmer), or $^{3}H$
817    uridine (0.65 µCi, Perkin Elmer) or $^{3}H$ thymidine (2 µCi, Perkin Elmer) – and transferred to
818    the heat block at 37 °C. Immediately after, the expression of toxins in the remaining 14 mL
819    culture was induced by addition of L-arabinose (final concentration of 0.2%). Throughout the
820    toxin induction time course, 1-mL aliquots were taken from the 15 mL culture and transferred
821    to 1.5 mL Eppendorf tubes containing 10 µl of radioisotope ($^{35}S$ methionine / $^{3}H$ uridine / $^{3}H$
822    thymidine). The incorporation of radioisotopes was stopped after 8 minutes of incubation at 37
823    °C by adding 200 µL of ice-cold 50% trichloroacetic acid (TCA) to 1 mL cultures. In parallel
824    with taking the time-points for labelling, 1 mL aliquots were taken for $OD_{600}$ measurements.
825    Isotope incorporation was quantified by normalising radioactivity counts (CPM) to $OD_{600}$, with
826    the pre-induction zero time-point set as 100%.
827    All experiments were performed in three biological replicates (i.e. using three independent
828    cultures inoculated from three different colonies).
829
830    **Methods references**
831

832    41.     Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold.
833    *Nature* **596**, 583–589 (2021).
834    42.     Hagberg, A. A., Schult, D. A. & Swart, P. J. Exploring Network Structure,
835    Dynamics, and Function using NetworkX. in *Proceedings of the 7th Python in Science*
836    *Conference* (eds. Varoquaux, G., Vaught, T. & Millman, J.) 11–15 (2008).
837    43.     Sehnal, D. *et al.* Mol\* Viewer: modern web app for 3D visualization and
838    analysis of large biomolecular structures. *Nucleic Acids Res.* **49**, W431–W437 (2021).
839    44.     Durairaj, J., Akdel, M., de Ridder, D. & van Dijk, A. D. J. Geometricus
840    Represents Protein Structures as Shape-mers Derived from Moment Invariants. Preprint

841       at https://doi.org/10.1101/2020.09.07.285569.

842       45.      Flusser, J., Boldys, J. & Zitova, B. Moment forms invariant to rotation and
843       blur in arbitrary number of dimensions. *IEEE Transactions on Pattern Analysis and*
844       *Machine Intelligence* vol. 25 234–246 Preprint at
845       https://doi.org/10.1109/tpami.2003.1177154 (2003).

846       46.      Flusser, J., Suk, T. & Zitová, B. 2D and 3D Image Analysis by Moments.
847       Preprint at https://doi.org/10.1002/9781119039402 (2016).

848       47.      Mamistvalov, A. G. n-dimensional moment invariants and conceptual
849       mathematical theory of recognition n-dimensional solids. *IEEE Transactions on Pattern*
850       *Analysis and Machine Intelligence* vol. 20 819–831 Preprint at
851       https://doi.org/10.1109/34.709598 (1998).

852       48.      Hattne, J. & Lamzin, V. S. A moment invariant for evaluating the chirality of
853       three-dimensional objects. *J. R. Soc. Interface* **8**, 144–151 (2011).

854       49.      Paszke, A. *et al.* PyTorch: An imperative style, high-performance deep
855       learning library. (2019) doi:10.48550/ARXIV.1912.01703.

856       50.      Das, S. *et al.* Functional classification of CATH superfamilies: a domain-
857       based approach for protein function annotation. *Bioinformatics* vol. 32 2889–2889
858       Preprint at https://doi.org/10.1093/bioinformatics/btw473 (2016).

859       51.      Zhang, C., Shine, M., Pyle, A. M. & Zhang, Y. US-align: universal structure
860       alignments of proteins, nucleic acids, and macromolecular complexes. *Nat. Methods* **19**,
861       1109–1115 (2022).

862       52.      AlQuraishi, M. ProteinNet: a standardized data set for machine learning of
863       protein structure. *BMC Bioinformatics* **20**, 311 (2019).

864       53.      Bojanowski, P., Grave, E., Joulin, A. & Mikolov, T. Enriching Word Vectors
865       with Subword Information. *Transactions of the Association for Computational*
866       *Linguistics* vol. 5 135–146 Preprint at https://doi.org/10.1162/tacl_a_00051 (2017).

867       54.      Rehurek, R. & Sojka, P. Gensim--python framework for vector space
868       modelling. *NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech*
869       *Republic*.

870       55.      Holm, L. Using Dali for Protein Structure Comparison. *Methods Mol. Biol.*
871       **2112**, 29–42 (2020).

872       56.      Zhang, Y. & Skolnick, J. TM-align: a protein structure alignment algorithm
873       based on the TM-score. *Nucleic Acids Res.* **33**, 2302–2309 (2005).

874       57.      Mavridis, L. & Ritchie, D. W. 3D-blast: 3D protein structure alignment,
875       comparison, and classification using spherical polar Fourier correlations. *Pac. Symp.*
876       *Biocomput.* 281–292 (2010).

877       58.      Wang, S. & Zheng, W.-M. CLePAPS: fast pair alignment of protein structures
878       based on conformational letters. *J. Bioinform. Comput. Biol.* **6**, 347–366 (2008).

879       59.      Liu, F. T., Ting, K. M. & Zhou, Z.-H. Isolation Forest. *2008 Eighth IEEE*
880       *International Conference on Data Mining* Preprint at
881       https://doi.org/10.1109/icdm.2008.17 (2008).

882       60.      Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn.*
883       *Res.* **12**, 2825–2830 (2011).

884       61.      Gabler, F. *et al.* Protein Sequence Analysis Using the MPI Bioinformatics
885       Toolkit. *Curr. Protoc. Bioinformatics* **72**, e108 (2020).

886       62.      Pereira, J. & Alva, V. How do I get the most out of my protein sequence using
887       bioinformatics tools? *Acta Crystallogr D Struct Biol* **77**, 1116–1126 (2021).

888       63.      Evans, R. *et al.* Protein complex prediction with AlphaFold-Multimer. *bioRxiv*
889       2021.10.04.463034 (2021) doi:10.1101/2021.10.04.463034.

890       64.      Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and

891 high throughput. *Nucleic Acids Research* vol. 32 1792–1797 Preprint at
892 https://doi.org/10.1093/nar/gkh340 (2004).

893 65. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of
894 protein database search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).

895 66. Frickey, T. & Lupas, A. CLANS: a Java application for visualizing protein
896 families based on pairwise similarity. *Bioinformatics* **20**, 3702–3704 (2004).

897 67. Wang, Y., Huang, H., Rudin, C. & Shaposhnik, Y. Understanding how
898 dimension reduction tools work: an empirical approach to deciphering t-SNE, UMAP,
899 TriMAP, and PaCMAP for data visualization. *arXiv preprint arXiv:2012. 04456* (2020).

900 68. Ester, M., Kriegel, H.-P., Sander, J. & Xu, X. A density-based algorithm for
901 discovering clusters in large spatial databases with noise. in *Proceedings of the Second*
902 *International Conference on Knowledge Discovery and Data Mining* 226–231 (AAAI
903 Press, 1996).

904 69. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**,
905 e1002195 (2011).

906 70. Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: lightning-fast
907 iterative protein sequence searching by HMM-HMM alignment. *Nat. Methods* **9**, 173–
908 175 (2011).

909 71. Quan, J. & Tian, J. Circular polymerase extension cloning for high-throughput
910 cloning of complex and combinatorial DNA libraries. *Nat. Protoc.* **6**, 242–251 (2011).

911 72. Guzman, L. M., Belin, D., Carson, M. J. & Beckwith, J. Tight regulation,
912 modulation, and high-level expression by vectors containing the arabinose PBAD
913 promoter. *J. Bacteriol.* **177**, 4121 (1995).

914 73. Jaskólska, M. & Gerdes, K. CRP-dependent positive autoregulation and
915 proteolytic degradation regulate competence activator Sxy of Escherichia coli. *Mol.*
916 *Microbiol.* **95**, 833–845 (2015).

917 74. Neidhardt, F. C., Bloch, P. L. & Smith, D. F. Culture medium for
918 enterobacteria. *J. Bacteriol.* **119**, 736–747 (1974).

935

936 **Author contributions**

937 J.P. and J.D. conceptualised the study. J.P. performed the functional darkness analysis and
938 constructed the sequence-based network. J.D. performed the structure outlier analysis. A.M.W.
939 developed the interactive web resource and J.P., J.D. and G.T. coordinated its development.
940 J.P., J.D., A.B. and A.A. performed the computational analysis of selected examples. G.S.,
941 M.Akdel, J.P., J.D. and A.M.W. developed computational methodologies. T.M., T.B. and M.
942 Abdullah carried out wet-lab experiments. V.H. and T.T. conceptualised, coordinated and
943 supervised wet-lab experiments. T.S., A.B., V.H., T.T., G.T. and J.P. acquired funding. J.P.
944 and J.D. wrote the original draft. All authors contributed, reviewed, edited and approved the
945 manuscript.

948 **Data availability statement**
949 All data used for this study is publicly available in UniProtKB (https://www.uniprot.org/,
950 UniRef version 2022_03), the AlphaFold database (https://alphafold.ebi.ac.uk/, version 4, with
951 specific examples corresponding to UniProt IDs A0A0E3S9F7, A0A3R7AQ40,
952 A0A520JWH3, A0A1W9UY89, A0A7J4P9B0, A0A0F9A5W1, A0A0P9GTS8,
953 AOA418VYX3, A0A2S5M855, A0A2K2VML8, A0A098EYBO, G0TGH8, A0A015IZK3,
954 A0A377W562, A0A494VZL1, A0A0S7BXY3, A0A7X7MB17, YFHO_BACSU,
955 A8JBY2_CHLRE, and A0A3A8FAL8), the CATH database (https://www.cathdb.info/,
956 version 4.2.0), ProteinNet (https://github.com/aqlaboratory/proteinnet, CASP12 dataset),
957 Foldseek benchmark data (https://wwwuser.gwdg.de/~compbiol/foldseek), the Protein Data
958 Bank (https://www.ebi.ac.uk/pdbe/, PDB IDs 5FMT, 5GKH, 8D3P, 6SK0, 2FIM, 1ZXU,
959 6GXC and 7OCI), and NCBI GenBank (https://www.ncbi.nlm.nih.gov/protein/, EntrezIDs
960 WP_213381069.1 and WP_213381068.1).
961 For the laboratory experiments all data generated are included in the manuscript and
962 supplementary materials. All data and metadata generated supporting the large and the
963 individual sequence similarity networks are available at https://zenodo.org/record/8121336
964 (CC-BY 4.0). An interactive version of the large sequence similarity network, queryable by
965 keyword, UniProt ID, connected component ID, community ID, protein sequence, and protein
966 structure, is available at https://uniprot3d.org/atlas/AFDB90v4. The interactive resource allows
967 also for the downloading of the metadata associated with each individual connected component
968 and community, as well as for the results of any search.
969
970 **Code availability statement**
971 All the code to collect and process the annotation data in UniProtKB, UniParc and InterPro,
972 and the pLDDT data from AFDB is available at
973 https://github.com/ProteinUniverseAtlas/dbuilder. Model and training code for shape-mer
974 generation can be found in https://github.com/TurtleTools/geometricus/tree/master/training.
975 All analysis code, including that to process the large sequence similarity network, decompose
976 structures and generate the plots displayed, is available at
977 https://github.com/ProteinUniverseAtlas/AFDB90v4 (Apache).
978

**Additional information statement**

Supplementary Information is available for this paper. Correspondence and requests for materials should be addressed to Joana Pereira (joana.pereira@unibas.ch) or Torsten Schwede (torsten.schwede@unibas.ch). Reprints and permissions information is available at www.nature.com/reprints.

**Extended data**

**Extended data figure 1. Distribution of functional darkness in UniProt and AFDB (version 4).** Functional brightness distribution in (a) UniRef50, (b) UniRef50 clusters with models in AFDB (which excludes long proteins, and those UniRef50 clusters composed solely of UniParc entries and viral proteins), (c) UniRef50 clusters whose best structural representative has an average pLDDT > 70, and (d) UniRef50 clusters whose best structural representative has an average pLDDT > 90. For each set, the percentage of fully dark UniRef50 clusters, and corresponding brightness bin, are highlighted in purple. The bar associated with functionally bright UniRef50 clusters (functional brightness >95%) is marked in white. (e) Percentage of fully dark UniRef50 clusters with proteins annotated as a domain of unknown function (DUF) in each set a-e.

**Extended data figure 2. Structural conservation and structure-based function prediction of TumE.** Structural superposition of five randomly selected members of component 159 (UniProt IDs A0A0E3S9F7, A0A3R7AQ40, A0A520JWH3, A0A1W9UY89, A0A7J4P9B0) with secondary structure elements labelled.

**Extended data figure 3. Testing the toxicity of putative TumA antitoxins.** Antitoxin expression plasmids were cotransformed with empty toxin expression vectors (pBAD33) into *E. coli* BW25113 cells. The bacterial cultures were started from a single colony and grown for five hours in liquid LB media supplemented with appropriate antibiotics. The cultures were normalised to $OD_{600} = 1.0$, serially diluted and spotted on LB agar plates containing appropriate antibiotics and 500 μM IPTG for antitoxin induction and 0.2% arabinose to mimic the conditions in toxin neutralisation assay. The experiment was made in n=3 biologically independent replicates. For source data, see Supplementary figure 2.

**Extended data figure 4. Diversity of the (a) names predicted by ProtNLM and (b) their word composition, as well as the (c) fraction of structural outliers, for all fully dark and fully bright connected components.** Name diversity is calculated as the number of unique protein names within a component by the total number of component proteins. Word diversity is calculated as the number of unique words across all protein names within a component by the total number of words, ignoring the words "protein", "domain", "family", "containing", and "superfamily". Outlier content is calculated as the percentage of UniRef50 clusters with negative structural outlier scores within that component. Fully bright and fully dark distributions were compared using a two-sided Kolmogorov–Smirnov test, resulting in a test

23

1022  statistic of 0.2915 and P-value = $8.8829 \times 10^{-16}$ for (b) and test statistic 0.05859 and P-value =
1023  $5.245 \times 10^{-81}$ for (c).

1024

1025  **Extended data figure 5. The highly semantically diverse prophage-associated connected**
1026  **components 3314 and 6732.** (a) Sequence similarity network of homologs of members of
1027  connected component 3314 and the tubulin-binding domain of TRAF3-interacting protein 1,
1028  as computed with CLANS at an E-value threshold of $1 \times 10^{-5}$. Points represent individual
1029  proteins and grey lines BLASTp matches at an E-value better than $1 \times 10^{-4}$. Individual
1030  subclusters are labelled 1-2 and structural representatives are shown. For subcluster 1, 5
1031  randomly selected structural representatives of component 3314 are superposed (UniProt IDs
1032  A0A0F9A5W1, A0A0P9GTS8, AOA418VYX3, A0A2S5M855, A0A2K2VML8). For
1033  subcluster 2, the tubulin-binding domain of *Chlamydomonas reinhardtii* TRAF3-interacting
1034  protein 1 (PDB ID 5FMT, chain B) is shown. (b) Genomic context conservation of 30
1035  sequences from subcluster 1 with a maximum sequence identity of 30%, as computed with
1036  GCsnap. (c) Structure superposition of component 6732 representative (A0A098EYBO,
1037  purple) and mismatch restriction endonuclease EndoMS (PDB ID 5GKH, chain A, grey). The
1038  grey box indicates the active site pocket with conserved residues labelled. Note that the residue
1039  D165 corresponding to D86 is mutated to alanine in the PDB structure. Structural homologs
1040  were searched both with Foldseek, which resulted in a hit to Cas4 endonuclease PDB ID 8D3P
1041  with TM-score 0.34, and Dali[55] multiple hits to restriction endonucleases, the top-ranking with
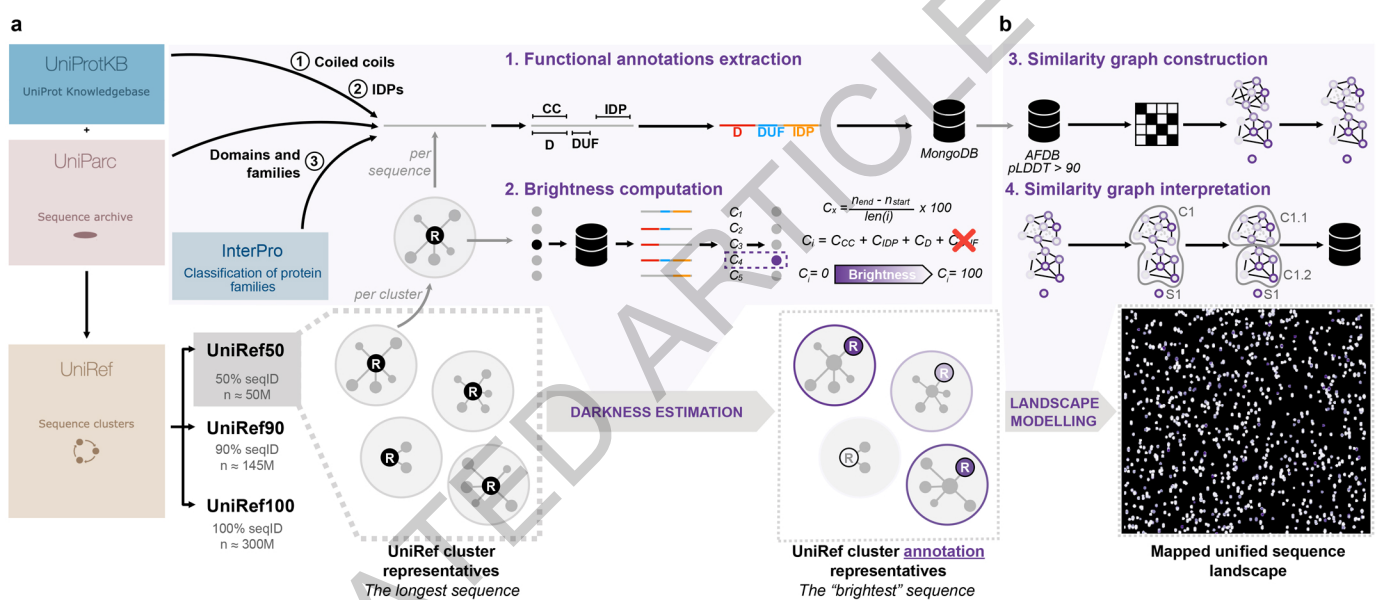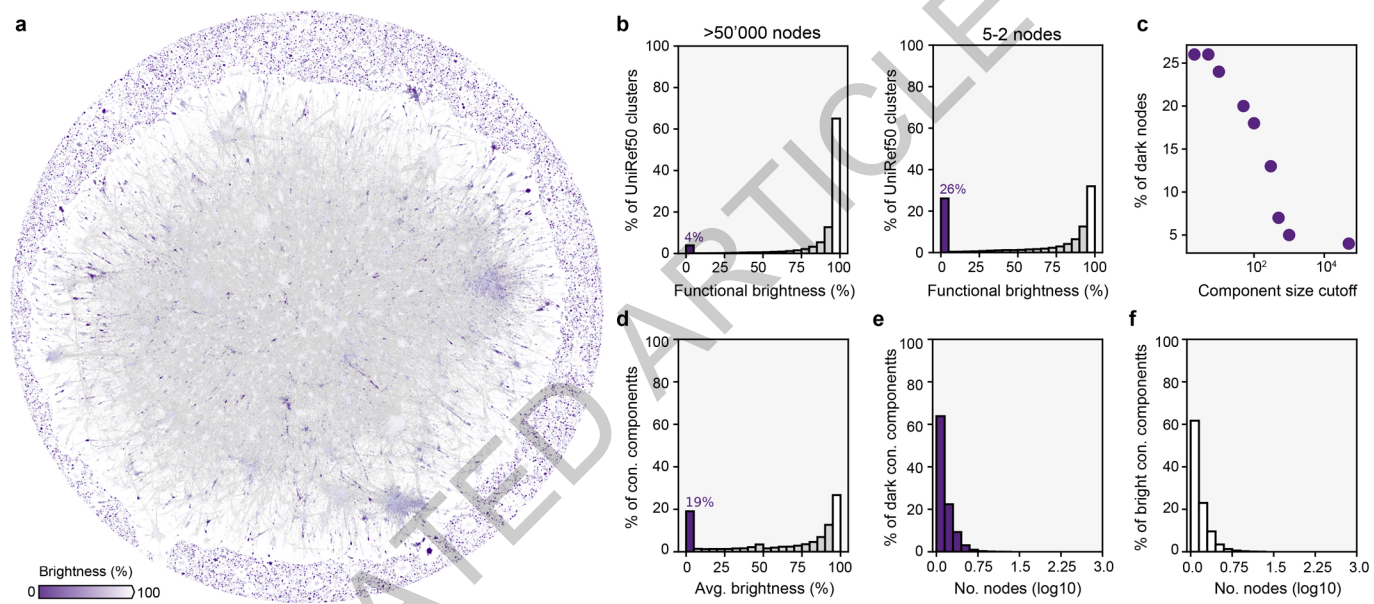1042  a Z-score of 8.2.

1043

1044  **Extended data figure 6. An example of substructure decomposition.** (a) An example
1045  AlphaFold protein model with its 6 most common shape-mers highlighted in different colours.
1046  Spheres mark the shape-mer central residue and backbone atoms within 4Å are coloured. (b-g)
1047  Four random representatives of each selected shape-mer, obtained from CATH proteins with
1048  <20% sequence identity. Spheres depict positions within 8 residues in sequence and 10Å
1049  spatially from the central residue.

1050

1051  **Extended data figure 7. Shape-mer representations combined with FastText can**
1052  **discriminate between protein families.** (a) Cumulative distributions of sensitivity for
1053  homology detection on the SCOPe40 database of single-domain structures. True positives
1054  (TPs) are matches within the same SCOPe family, false positives (FPs) are matches between
1055  different folds. Sensitivity is the area under the ROC curve up to the first FP. Results based on
1056  shape-mer FastText Smith-Waterman alignment are shown in black. (b)  Protein-level
1057  embedding distance measured as the cosine distance of FastText sentence vectors for proteins
1058  within the same SCOPe family (top) and from different SCOPe folds (bottom).

1059

**a**

UniProtKB
UniProt Knowledgebase

+

UniParc
Sequence archive

UniRef
Sequence clusters

UniRef50
50% seqID
n ≈ 50M

UniRef90
90% seqID
n ≈ 145M

UniRef100
100% seqID
n ≈ 300M

① Coiled coils
② IDPs
③ Domains and families

InterPro
Classification of protein families

**1. Functional annotations extraction**

$$C_x = \frac{n_{end} - n_{start}}{len(i)} \times 100$$

$$C_i = C_{CC} + C_{IDP} + C_D \ \times \ C_{DUF}$$

$$C_i = 0 \quad \boxed{\text{Brightness}} \quad C_i = 100$$

**2. Brightness computation**

$C_1$
$C_2$
$C_3$
$C_4$
$C_5$

*per sequence*

*per cluster*

MongoDB

DARKNESS ESTIMATION

**UniRef cluster representatives**
*The longest sequence*

**UniRef cluster annotation representatives**
*The "brightest" sequence*

**b**

**3. Similarity graph construction**

AFDB
pLDDT > 90

**4. Similarity graph interpretation**

C1
C1.1
C1.2
S1
S1

LANDSCAPE MODELLING

**Mapped unified sequence landscape**

**a** Brightness (%) 0 — 100

**b** >50'000 nodes — % of UniRef50 clusters vs Functional brightness (%); 4%

5-2 nodes — % of UniRef50 clusters vs Functional brightness (%); 26%

**c** % of dark nodes vs Component size cutoff

**d** % of con. components vs Avg. brightness (%); 19%

**e** % of dark con. componentts vs No. nodes (log10)

**f** % of bright con. components vs No. nodes (log10)

**a**

C27

C27 *YfhO-like*

AglB

STT3     PglB

YfhO

● All sequences in C27
○ Bright sequences in C27

**b**

| C27 | YfhO | PglB | STT3 |
|---|---|---|---|
| A0A7X7MB17 | YFHO_BACSU | PGLB_CAMJR (6GXC_A) | STT3_YEAST (7OCI_F) |
| | *TM-score: 0.583* | *TM-score: 0.525* | *TM-score: 0.555* |

Membrane

**a**

C159

**2** — DUF6516
HP181

DUF6516 **3**
HP167

**7**

**1a** **1b**
**1c**
DUF6516
RelB-L

**6**

**5** UnkD-DUF6516

**4** DUF6516
HP185

**b**

RelB-L
*TumA antitoxin*

RelB-L
*TumA antitoxin*

DUF6516
*TumE toxin*

DUF6516'
*TumE toxin*

*DNA-binding region*

**c**

DNA-binding
*56%*

Hydrolase activity
acting on ester bonds
*40%*

**d**

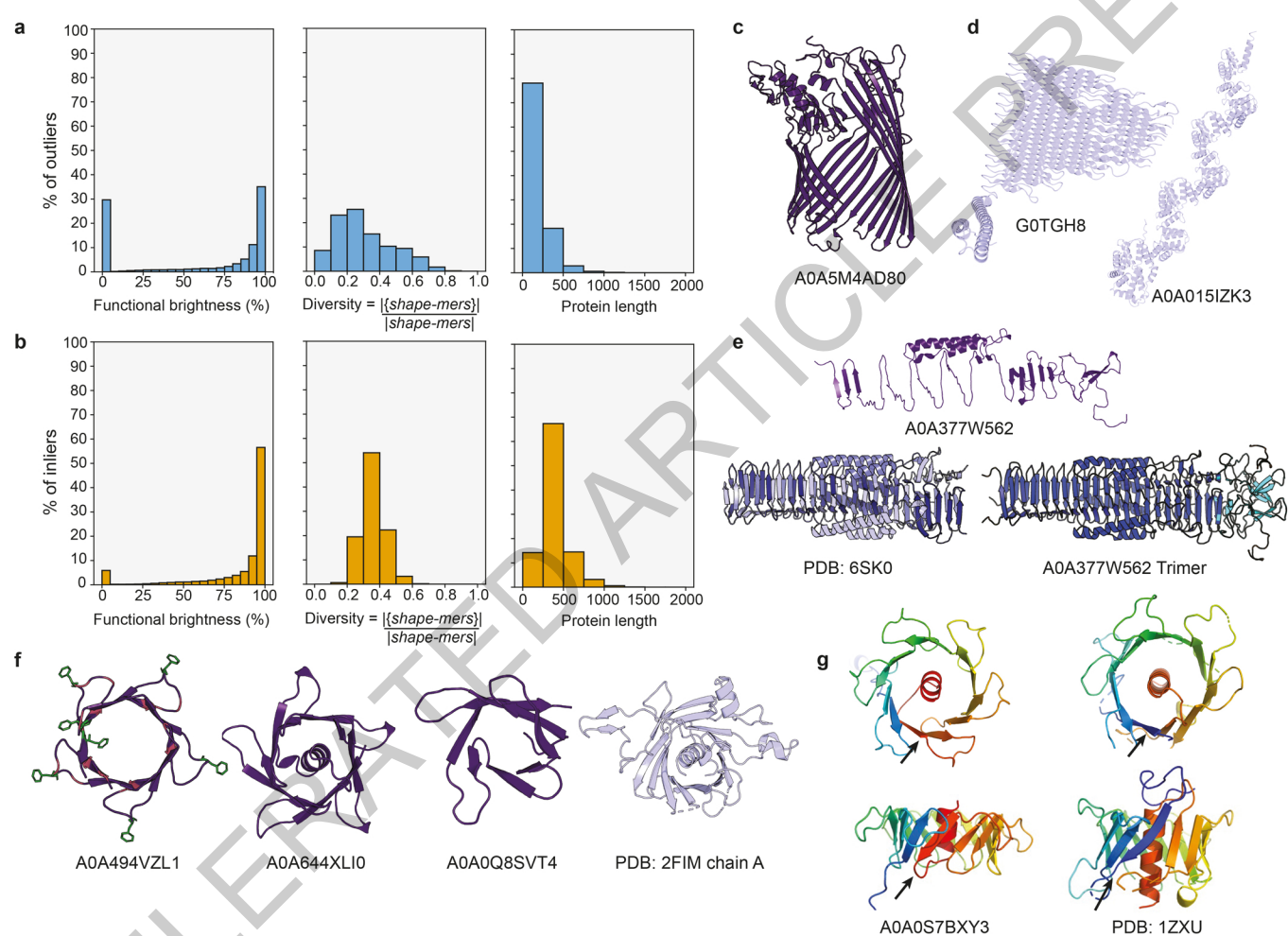| | TumE toxin | TumE toxin + TumA antitoxin |
|---|---|---|
| | $10^{-1}$ $10^{-2}$ $10^{-3}$ $10^{-4}$ $10^{-5}$ $10^{-6}$ $10^{-7}$ | $10^{-1}$ $10^{-2}$ $10^{-3}$ $10^{-4}$ $10^{-5}$ $10^{-6}$ $10^{-7}$ |
| *T. ingrica* | | |
| *T. litoralis* | | |
| *Crenothrix* sp. D3 | | |
| *M. alcaliphilum* | | |
| *E. bacterium* | | |
| *A. tepidum* | | |

**e**

Incorporation (%) vs Time (min)
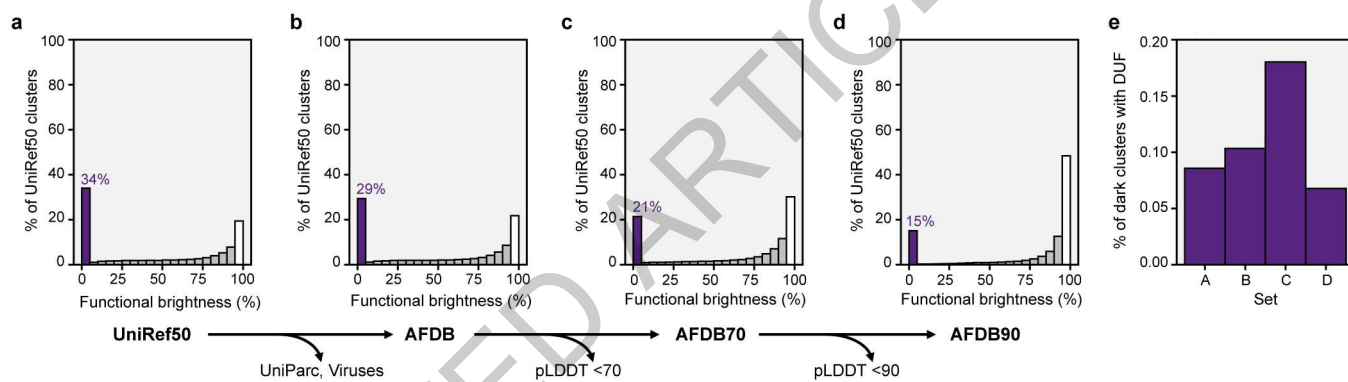
- $^{35}$S methionine
- $^3$H uridine
- $^3$H thymidine
- $OD_{600}$

**Extended Data Fig. 1**

**Extended Data Fig. 2**

no induction        500 µM IPTG

$10^{-1}$ $10^{-2}$ $10^{-3}$ $10^{-4}$ $10^{-5}$ $10^{-6}$ $10^{-7}$     $10^{-1}$ $10^{-2}$ $10^{-3}$ $10^{-4}$ $10^{-5}$ $10^{-6}$ $10^{-7}$

*T. ingrica*

*T. litoralis*

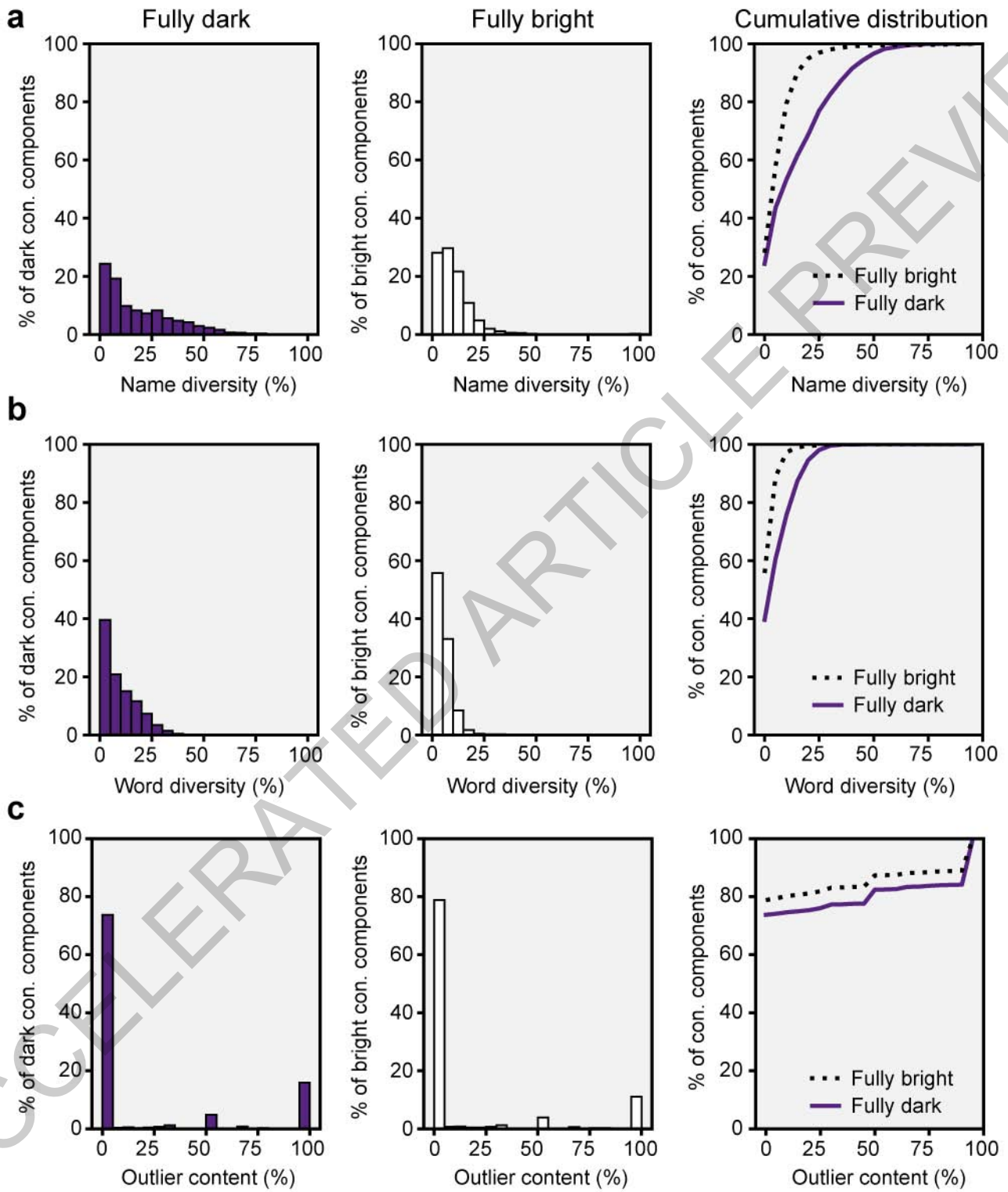*Crenothrix* sp. D3

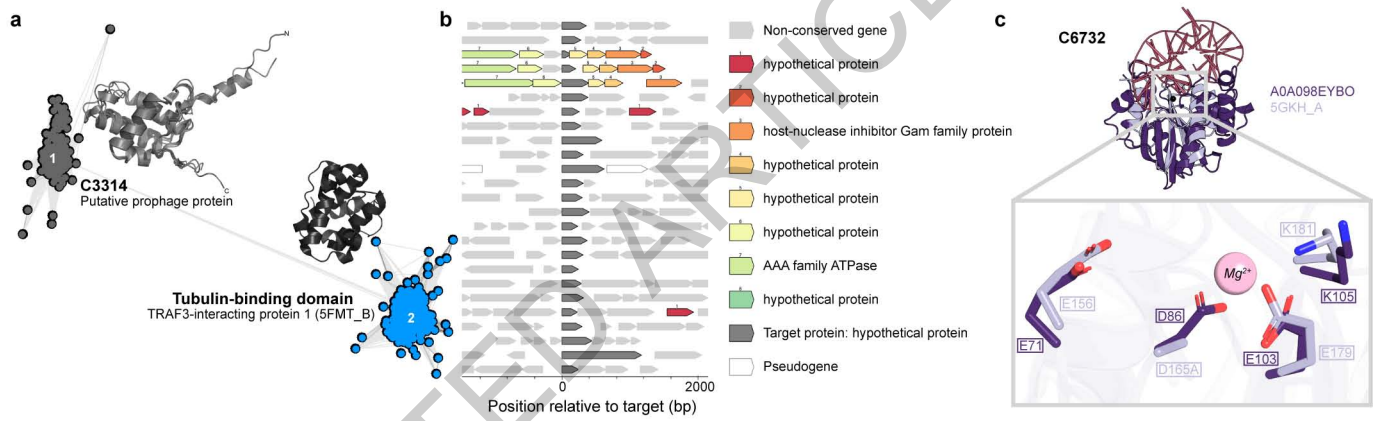*M. alcaliphilum*

*E. bacterium*
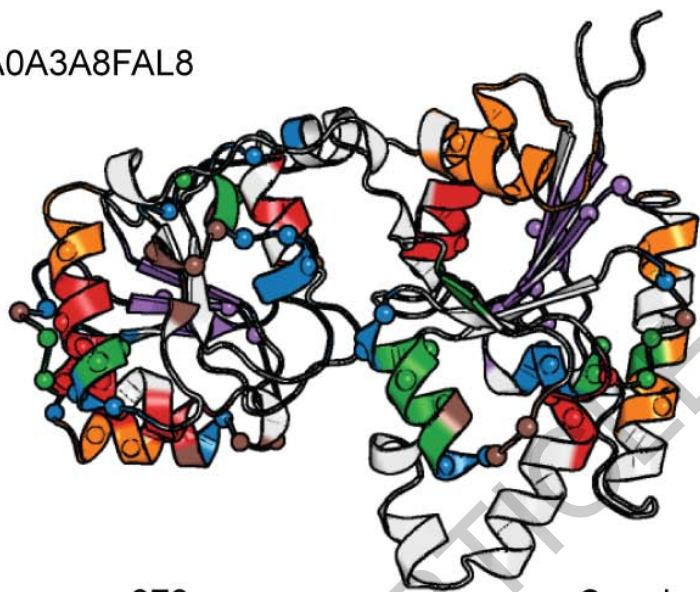
*A. tepidum*

**Extended Data Fig. 3**

**Extended Data Fig. 4**

**Extended Data Fig. 5**

**a** A0A3A8FAL8

**b** shape-mer 878

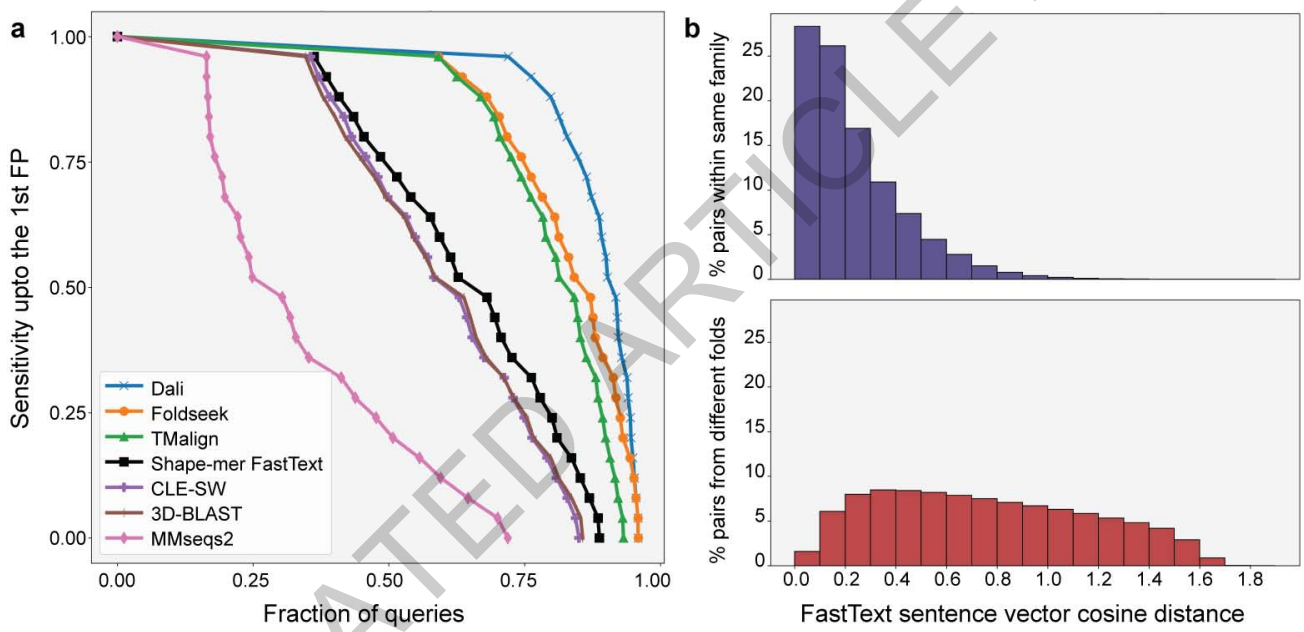**c** shape-mer 621

**d** shape-mer 622

**e** shape-mer 110

**f** shape-mer 662

**g** shape-mer 846

**Extended Data Fig. 6**

**Extended Data Fig. 7**

Corresponding author(s): Joana Pereira
Torsten Schwede

Last updated by author(s): 29.08.2023

# Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our Editorial Policies and the Editorial Policy Checklist.

## Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

| n/a | Confirmed | |
|-----|-----------|---|
| ☐ | ☒ | The exact sample size (*n*) for each experimental group/condition, given as a discrete number and unit of measurement |
| ☐ | ☒ | A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly |
| ☐ | ☒ | The statistical test(s) used AND whether they are one- or two-sided<br>*Only common tests should be described solely by name; describe more complex techniques in the Methods section.* |
| ☒ | ☐ | A description of all covariates tested |
| ☒ | ☐ | A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons |
| ☐ | ☒ | A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| ☐ | ☒ | For null hypothesis testing, the test statistic (e.g. *F*, *t*, *r*) with confidence intervals, effect sizes, degrees of freedom and *P* value noted<br>*Give P values as exact values whenever suitable.* |
| ☒ | ☐ | For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings |
| ☒ | ☐ | For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes |
| ☐ | ☒ | Estimates of effect sizes (e.g. Cohen's *d*, Pearson's *r*), indicating how they were calculated |

*Our web collection on statistics for biologists contains articles on many of the points above.*

## Software and code

Policy information about availability of computer code

| Data collection | Data collection for the annotated network was carried out using custom code available at https://github.com/ProteinUniverseAtlas/dbuilder which uses Python 3.6 and PyMongo v3.11.3.<br>Training data for protein substructure decomposition and outlier detection was created using custom code available at https://github.com/TurtleTools/geometricus/tree/master/training using Python 3.9, cath-tools-genomescan (version 17/12/2019), and ProteinNet (CASP12 dataset) |
|-----------------|---|
| Data analysis | Custom code for data analysis can be found at https://github.com/ProteinUniverseAtlas/AFDB90v4 and uses:<br>Python 3.6, 3.9<br>SciPy (v1.5.4)<br>NetworkX (v2.5.1)<br>ProDy (v2.2.0)<br>Geometricus (v0.5.0)<br>PyTorch (v1.12.0)<br>Gensim (v4.2.0)<br>scikit-learn (v1.1.1)<br>Datashader (v0.12.1)<br><br>In addition, the following tools were used for analyses as described in the Methods:<br>MMseqs (release 13-45111)<br>MUSCLE (v5.1)<br>GCsnap (v1.0.17) |

DeepFRI (v1.0.0)
Foldseek (Version: 7.04e0ec8)
AlphaFold (v2.3.0)
PyMol (open-source v2.5.0)

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio guidelines for submitting code & software for further information.

## Data

Policy information about availability of data

All manuscripts must include a data availability statement. This statement should provide the following information, where applicable:
- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our policy

All data used for this study is publicly available in UniProtKB (https://www.uniprot.org/, UniRef version 2022_03), the AlphaFold database (https://alphafold.ebi.ac.uk/, version 4, with specific examples corresponding to UniProt IDs A0A0E3S9F7, A0A3R7AQ40, A0A520JWH3, A0A1W9UY89, A0A7J4P9B0, A0A0F9A5W1, A0A0P9GTS8, AOA418VYX3, A0A2S5M855, A0A2K2VML8, A0A098EYBO, G0TGH8, A0A015IZK3, A0A377W562, A0A494VZL1, A0A0S7BXY3, A0A7X7MB17, YFHO_BACSU, A8JBY2_CHLRE, and A0A3A8FAL8), the CATH database (https://www.cathdb.info/, version 4.2.0), ProteinNet (https://github.com/aqlaboratory/proteinnet, CASP12 dataset), Foldseek benchmark data (https://wwwuser.gwdg.de/~compbiol/foldseek), the Protein Data Bank (https://www.ebi.ac.uk/pdbe/, PDB IDs 5FMT, 5GKH, 8D3P, 6SK0, 2FIM, 1ZXU, 6GXC and 7OCI), and NCBI GenBank (https://www.ncbi.nlm.nih.gov/protein/, EntrezIDs WP_213381069.1 and WP_213381068.1).
For the laboratory experiments all data generated are included in the manuscript and supplementary materials. All data and metadata generated supporting the large and the individual sequence similarity networks are available at https://zenodo.org/record/8121336 (CC-BY 4.0). An interactive version of the large sequence similarity network, queryable by keyword, UniProt ID, connected component ID, community ID, protein sequence, and protein structure, is available at https://uniprot3d.org/atlas/AFDB90v4. The interactive resource allows also for the downloading of the metadata associated with each individual connected component and community, as well as for the results of any search.

## Research involving human participants, their data, or biological material

Policy information about studies with human participants or human data. See also policy information about sex, gender (identity/presentation), and sexual orientation and race, ethnicity and racism.

| | |
|---|---|
| Reporting on sex and gender | Not applicable. No human participants or human data was used in this study. |
| Reporting on race, ethnicity, or other socially relevant groupings | Not applicable. No human participants or human data was used in this study. |
| Population characteristics | Not applicable. No human participants or human data was used in this study. |
| Recruitment | Not applicable. No human participants or human data was used in this study. |
| Ethics oversight | Not applicable. No human participants or human data was used in this study. |

Note that full information on the approval of the study protocol must also be provided in the manuscript.

# Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences  ☐ Behavioural & social sciences  ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

# Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

| | |
|---|---|
| Sample size | We followed the standard practices in the toxin-antitoxin molecular microbiology field. The sample size of three is regularly used in toxin-antitoxin microbiology field for growth and metabolic labeling experiments. The effects were strong and do not require further statistical analysis. |
| Data exclusions | No data were excluded. |
| Replication | The experiments were repeated in at least three biological independent replicates. All of the attempts were successful and showed the same results. |

| | |
|---|---|
| Randomization | We followed the standard practices in the toxin-antitoxin molecular microbiology field. Randomization of samples is generally not practiced. |
| Blinding | We followed the standard practices in the toxin-antitoxin molecular microbiology field. Blinding is generally not practiced. |

# Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

## Materials & experimental systems

| n/a | Involved in the study |
|---|---|
| ☒ | Antibodies |
| ☒ | Eukaryotic cell lines |
| ☒ | Palaeontology and archaeology |
| ☒ | Animals and other organisms |
| ☒ | Clinical data |
| ☒ | Dual use research of concern |
| ☒ | Plants |

## Methods

| n/a | Involved in the study |
|---|---|
| ☒ | ChIP-seq |
| ☒ | Flow cytometry |
| ☒ | MRI-based neuroimaging |