

1 Energetic Profile-Based Protein Comparison: A New and Fast Approach for
2 Structural and Evolutionary analysis

3 Peyman Choopanian^{1, 2}, Jaan-Olle Andressoo^{1, 2, 3*}, and Mehdi Mirzaie^{1, 2*}

4 ¹Translational Neuroscience, Department of Pharmacology, Faculty of Medicine and Helsinki Institute of
5 Life Science, 00014 University of Helsinki, Finland

6 ²Department of Pharmacology, Faculty of Medicine, 00014 University of Helsinki, Finland

7 ³Division of Neurogeriatrics, Department of Neurobiology, Care Sciences and Society (NVS), 17177
8 Karolinska Institutet, Sweden

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30 **Abstract**

31 In structural bioinformatics, the efficiency of predicting protein similarity, function, and evolutionary
32 relationships is crucial. Our innovative approach leverages protein energy profiles derived from a
33 knowledge-based potential, deviating from traditional methods relying on structural alignment or atomic
34 distances. This method assigns unique energy profiles to individual proteins, facilitating rapid comparative
35 analysis for both structural similarities and evolutionary relationships across various hierarchical levels. Our
36 study demonstrates that energy profiles contain substantial information about protein structure at class,
37 fold, superfamily, and family levels. Notably, these profiles accurately distinguish proteins across species,
38 illustrated by the classification of coronavirus spike glycoproteins and bacteriocin proteins. Introducing a
39 novel separation measure based on energy profile similarity, our method shows significant correlation with
40 a network-based approach, emphasizing the potential of energy profiles as efficient predictors for drug
41 combinations with faster computational requirements. Our key insight is that the sequence-based energy
42 profile strongly correlates with structure-derived energy, enabling rapid and efficient protein comparisons
43 based solely on sequences.

44

45

46 **Keywords:** Energy-based annotation, Structural dissimilarity, Evolutionary relationships, Profile of
47 energy, Knowledge-based potential.

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66 **Introduction**

67 A thorough understanding of protein function holds paramount importance within the domains of biology,
68 medicine, and pharmacy. While experimental methods exhibit high accuracy in protein function
69 associations, their inherent limitations, such as being time-intensive and expensive, have instigated the
70 exploration of computational alternatives. The evaluation of protein similarity by comparing two proteins
71 has consistently emerged as a key methodology. This assessment plays a pivotal role in uncovering
72 insights into the functions and evolutionary relationships of proteins. Advances in high-throughput
73 technologies have led to the establishment of extensive repositories containing protein sequences, a
74 substantial proportion of which, however, lack annotations¹. The significant advancements in omics data
75 and the evolution of machine learning techniques have propelled progress in protein research, transitioning
76 from traditional methods like PSI-BLAST² to more sophisticated approaches³. In the realm of machine
77 learning research, a crucial step is encoding data as input. Although there is no universal approach,
78 encoding amino acid sequences or structural features has been widely adopted for various protein function
79 predictions, including drug-protein interactions⁴, anti-hypertensive peptides⁵, and RNA-protein
80 interactions⁶. Despite the diversity in methodologies, the underlying commonality revolves around
81 determining protein similarity either through sequence alignment or structural comparison.

82 Protein structure is a fundamental feature affecting function, and activity. The energy of a protein structure
83 plays a key role in determining its structure. Knowledge-based potentials, categorized as statistical energy
84 functions, are derived from databases of known protein structures that empirically capture the most
85 probable state of a protein and describe microstates of interactions within protein structure^{7, 8, 9}. It is
86 generally assuming that a native protein structure is confined to a state with the minimum total energy, and
87 the more similar a structure is to the native state, the closer its total energy is to the native state. However,
88 we take a step beyond this assumption and suggest a hypothesis that two similar proteins possess
89 analogous energy profiles. To evaluate this hypothesis, we assigned an energetic feature vector to each
90 protein, with each entry representing the summation of energies for a specific pair of amino acids. With 20
91 amino acids in proteins, this resulted in 210 pairwise interaction types. This is the first study to assign an
92 energetic feature vector to each protein for comparative analysis. This 210-dimensional vector represents
93 the intricate energy landscape inherent to the structural diversity of proteins. This vector of energies serves
94 as the cornerstone of our analytical approach and provides a robust foundation for further investigative
95 pursuits. Given the current issue of experimentally determining the three-dimensional structures of proteins,
96 estimating energy based on sequence emerges as a crucial consideration. Dostari et al.¹⁰ introduced a
97 method to estimate energy based on amino acid composition. In our study, we drew inspiration from their
98 approach to extract the energy profile based on protein sequence.

99 The stratification of proteins into distinct folds, superfamilies, and families, guided by evolutionary
100 consanguinity or shared structural and functional attributes, is crucial for precise function prediction.
101 Databases like CATH (Class, Architecture, Topology, Homologous superfamily)¹¹ and SCOP (Structural
102 Classification Of Proteins)¹² categorize proteins into hierarchical groups based on structural feature, from

103 broad classifications like folds and classes to finer details such as superfamilies and families. To assess
104 profile of energies at various levels, we utilized the ASTRAL40 (95) datasets from SCOPe as a benchmark
105 dataset¹³. Comparing energy and distance between profiles estimated from both sequence and structure
106 revealed a high correlation on protein domains from both ASTRAL40 and ASTRAL95 datasets. UMAP
107 projections provided additional evidence that the profile of energy encapsulates structural information at
108 fold, superfamily, and family levels, as observed through random selections. Our method demonstrated
109 superior performance in terms of both accuracy and computational efficiency compared to currently
110 available tools.

111 In the realm of structural biology and evolutionary analysis, three-dimensional protein structure
112 classification and the alignment of multiple sequences stand as formidable tools for uncovering structural
113 similarities and deducing phylogenetic relationships. We also evaluated our method to reconstruct
114 evolutionary relationships among proteins from the ferritin-like superfamily that are beyond the "twilight
115 zone"¹⁴ —a sequence similarity range (typically 20-35% identity) that complicates the differentiation
116 between true homologs and random matches due to insufficient sequence conservation. Our findings
117 strongly suggest that a substantial and valuable evolutionary signal is preserved within the profile of energy,
118 serving as a representative indicator of protein structure. To assess the discriminatory capacity of energy
119 profiles in discerning proteins across various species, we chose the spike glycoproteins from three
120 coronavirus species¹⁵. Our findings indicate that both sequence-level and structural-level energy profiles
121 successfully cluster proteins from distinct species. In a separate analysis, we computed the sequence-
122 based energy profile for a diverse set of bacterial protein families known as bacteriocins. The identification
123 and understanding of these peptides are crucial due to their ecological significance, but their diverse
124 sequences and structures present challenges for conventional identification methods. The BAGEL data set
125 includes 689 proteins¹⁶, each with a length greater than 30 amino acids, providing a comprehensive and
126 challenging benchmark for evaluating peptide identification techniques. Our findings highlight that the
127 energy profile can categorize these proteins based on BAGEL annotation, demonstrating the effectiveness
128 of our method in handling the complexity and diversity inherent in bacteriocin sequences.

129 The identification of effective drug combinations, essential for treating complex diseases, face challenges
130 due to the combinatorial explosion of potential drug pairs. Cheng et al. introduced a network-based
131 methodology leveraging the human protein-protein interactome to discover clinically effective drug
132 combinations, demonstrating that topological relationships between drug-target modules, as indicated by a
133 separation measure, reflect both biological and pharmacological relationships¹⁷. In our study, we introduce
134 a separation measure based on the similarity between profile of energies of protein targets, revealing a
135 significant correlation with the separation measure derived from the protein-protein interaction network.
136 This suggests that the profile of energy holds promise as a reliable predictor for drug combinations,
137 requiring only protein sequences and offering quicker computation compared to network-based
138 approaches. Therefore, this study offers a means to characterize and compare proteins using profile of
139 energies, enabling predictions of their structural and functional properties.

140

141
142 **Results**
143 Knowledge-based potentials are derived from databases of known protein structures. Various potential
144 functions, such as distance-dependent, dihedral angles, and accessible surface energies leverage
145 information from known protein structures to estimate energies of pairwise interactions ⁸. In our
146 investigation, we employed the potential function, where atom contacts were identified using the tessellation
147 method, as outlined in the method section ¹⁸. For each protein structure, contacts between atoms were
148 determined through the tessellation method, and the energy for each pair of amino acid types was estimated
149 using equation (2) in the Method section. With 20 amino acids present in proteins, this process resulted in
150 210 pairwise interaction types. This 210-dimensional vector represents the energy landscape inherent in
151 the structural diversity of proteins. For each pair of proteins, the Manhattan distance between the profiles
152 of energies is considered a measure of dissimilarity between them. Given the current issue of
153 experimentally determining the three-dimensional structures of proteins, estimating energy based on
154 sequence pertains as a crucial subject. We utilized equation (5) to construct the profile of energy based on
155 sequence.

156 **Correlation between Energy estimated based on structure and Sequence.**
157 To examine the profile of energy at various levels of SCOP, we employed the ASTRAL40 (95) database
158 (version 2.08) from SCOPe as a benchmark dataset, comprising domains with no more than 40% (95%)
159 sequence similarity, as determined by BLAST identity, and filtered for E-value similarity scores ¹³. This
160 dataset offers a comprehensive description of structural and evolutionary relationships among proteins from
161 the Protein Data Bank. At first, we calculated energies for protein domains in the ASTRAL40 and
162 ASTRAL95 datasets using both structure- and sequence-based methods. Fig. 1A depicts the relationship
163 between the total energy derived from the structure (on the y-axis) and from the sequence (on the x-axis),
164 with the ASTRAL40 on the left and ASTRAL95 on the right side of the figure. The observed high correlation
165 coefficient suggests that sequence-based energy estimation serves as a reliable approximation and can be
166 effectively used in scenarios where the protein structure is unidentified.
167 For every pair of domains within the ASTRAL40 (ASTRAL95) datasets, the distances between their profile
168 of energy were computed utilizing both structural and sequence-based energy estimation. In Fig. 1B, the
169 x-axis denotes the distance between Compositional Profile of Energies (CPE), while the y-axis represents
170 the distance between Structural Profile of Energies (SPE)(for more details see the method section). The
171 figure reveals a strong correlation between the distances estimated through structural and sequence-based
172 approaches. Hence, the energy estimation based on sequence data is deemed sufficiently reliable.
173 The stability, mutational robustness, and design adaptability of α -helices relative to β -strands in natural
174 proteins have been widely acknowledged in scientific literature. To investigate this phenomenon, Fig. 2
175 presents the distribution of total energy within protein domains from the ASTRAL40 and ASTRAL95
176 datasets, categorized into four structural scope classes: all-alpha, all-beta, alpha + beta, and alpha/beta.
177 Total energies, normalized by protein length, are analyzed to discern patterns across these structural

178 classes. The figure highlights significant differences in total energy among domains with different structural
179 compositions, suggesting diverse energetic landscapes associated with distinct protein structures. This
180 observation is consistent with similar trends observed in energy estimations derived from sequence
181 information (Fig. 2B).

182

183 **3.2 Unveiling the Energy Patterns Across SCOP Hierarchy**

184 We visualized energy profiles derived from sequence and structure for domains within the all-alpha and all-
185 beta classes. As shown in Fig. 3, UMAP embeddings effectively capture structural characteristics
186 distinguishing all-alpha and all-beta domains. This visualization reveals distinct energy patterns between
187 these classes, a consistency also found in sequence-based analyses. To explore structural information at
188 lower hierarchical levels of SCOP, two folds (a.100 and a.104) from the all-alpha class, two superfamilies
189 (a.29.2 and a.29.3) from fold a.29, and two families (a.25.1.0 and a.25.1.2) from superfamily a.25.1 were
190 randomly selected. Fig. 4 displays two figures per panel, with the left figure illustrating CPE profiles and the
191 right figure showcasing SPE profiles. UMAP plots in Fig. 4 demonstrate that protein domains within the
192 same fold, superfamily, or family share similar energy patterns and cluster together.

193 To delve deeper into differences in distances among protein domains within the same class, we calculated
194 pairwise distances for domains within the all-alpha class from the ASTRAL95 dataset. Subsequently, these
195 distances were compared with distances from domains across different classes. As shown in Fig. 5A-B,
196 intraclass distances in purple are significantly lower than interclass distances in yellow. Similar results were
197 obtained when calculating pairwise distances from domains within fold a.29 and comparing them with
198 distances from domains in different folds within the all-alpha class. Likewise, distances between energy
199 patterns of domains within the same superfamily a.29.3 are significantly less than distances between
200 energy patterns of domains within fold a.29 that belong to different superfamilies (Fig. 5C-D). Consequently,
201 it can be inferred that energy patterns of domains belonging to the same superfamily/fold/class exhibit
202 higher similarity than those from different superfamilies/folds/classes.

203 It is commonly assumed that proteins sharing similar structures also exhibit similar functions. Several
204 measurements have been developed to assess protein structure similarity, each offering unique insights.
205 Root Mean Square Deviation (RMSD)¹⁹ quantifies average spatial variance between corresponding atoms
206 or components within superimposed proteins, providing a fundamental measure of structural deviation. The
207 TM-score (Template Modeling score)²⁰ evaluates similarity by considering both residue-level alignment
208 and overall topology, offering a nuanced assessment of structural resemblance. TM-Vec²¹, a recent
209 advancement, employs deep learning techniques trained on diverse protein structures to enhance accuracy
210 and efficiency in similarity assessment. On the alignment front, GR-align²² stands out, likely utilizing
211 geometric reasoning for accurate structural alignment. Its robustness to structural variations makes it
212 invaluable in structural bioinformatics. Additionally, the Hausdorff distance²³ provides a measure of
213 dissimilarity between sets of points, offering further insight into structural comparisons. Here, we employed
214 a benchmark dataset sourced from the CATH v4.2.0 database, comprising 251 protein domains from two
215 distinct protein families: the C-terminal domain in the DNA helicase RuvA subunit (representing the Alpha

216 class, characterized by Orthogonal Bundle Architecture, Helicase, and Ruva Protein fold, with CATH Code:
217 1.10.8.10), and the Homing endonucleases (belonging to the Alpha and Beta class, featuring Roll
218 Architecture, and Endonuclease I-crel fold, with CATH Code: 3.10.28.10). The protein domains varied in
219 the number of residues, ranging from 44 to 854, with an average of 211.

220 We used the 1-NN classification method to categorize proteins based on GR-Align, RMSD, TM-score, Yau-
221 Hausdorff distance, TM-Vec, and the distance between energy profiles as a measure of protein dissimilarity.
222 As shown in Table 1 and Fig. 6, our method achieves a computation time faster than TM-Vec. Moreover,
223 our method significantly outperforms GR-Align, RMSD, TM-Score, and YH in terms of accuracy and
224 efficiency, highlighting a substantial advantage. Our method eliminates the need for superimposing protein
225 structures or conducting structural alignments; instead, we calculate energy profiles and measure the
226 distance between them. Table 1 details result and processing times, demonstrating the efficient
227 implementation of CPE calculation and the 1-NN algorithm, completed in about 3 minutes on a system with
228 a 2.4 GHz processor and 4GB RAM. Our methodology achieved a remarkable classification accuracy of
229 97% in distinguishing between two protein families.

230 To assess the profile of energy in protein superfamily classification, we investigated five distinct SCOP
231 superfamilies: winged helix (a.4.5), PH domain-like (b.55.1), NTF-like (d.17.4), Ubiquitin-like (d.15.1), and
232 Immunoglobulins (b.1.1) ²⁴. Our classification strategy incorporated energetic profiles CPE as features,
233 employing 1-nearest neighbor (1-NN) and Random Forest (RF) classifiers as our models. To ensure the
234 robustness and generalization of our models, we subjected RF to rigorous 10-fold cross-validation. The
235 results, summarized in Table 2, include metrics for accuracy and F1-score, demonstrating the effectiveness
236 of our model. Both classifiers show performance levels close to 100%, as illustrated in Table 2. We
237 compared the CPE method with TM-Vec. As depicted in Table 2, our results (CPE) are not only comparable
238 to TM-Vec in terms of accuracy but also demonstrate a faster performance.

239

240 **3.4 Phylogeny Inference of the Ferritin-Like Superfamily**

241 In conjunction with the organizational frameworks provided by SCOP, CATH, and Pfam for the protein
242 universe, it is important to note their limitations, as they may present conflicting classifications and lack the
243 ability to elucidate evolutionary relationships between individual superfamilies across long evolutionary
244 distances. Lundin et al. conducted a comprehensive analysis of protein structures within the functionally
245 diverse ferritin-like superfamily. They employed an evolutionary network construction approach to unveil
246 relationships among proteins beyond the "twilight zone", where sequence similarity alone fails to facilitate
247 meaningful evolutionary analysis. Building on this context, our study leverages profiles of energies to
248 reconstruct a phylogenetic network. Our findings strongly suggest that a substantial and valuable
249 evolutionary signal is preserved within the profile of energy, serving as a representative indicator of protein
250 structure. Lundin et al. ¹⁴ investigated how ferritin-like proteins are classified across Pfam, SCOP, and
251 CATH. Notably, this superfamily encompasses a diverse range of proteins, including iron-storing ferritins,
252 methane monooxygenases, the small subunit of Ribonucleotide reductase-like (RNR R2), rubrerythrins,
253 bacterioferritins, Dps (DNA binding protein from starved cells that protects against oxidative DNA damage),

254 and Dps-like proteins. As discussed by Lundin et al.¹⁴ at the superfamily level, the classification of the
255 “ferritin-like” superfamily appears consistent across these databases but does differ in the amount of
256 information provided regarding the relationships and functions of superfamily constituents. So, although the
257 classification in all three databases is hierarchical, they do not encompass all level of functional and
258 evolutionary information. The low sequence similarities across this superfamily make it feasible to construct
259 sequence-based phylogenies only for specific subsets. Consequently, addressing this challenge requires
260 efforts to integrate structural information with sequence-based phylogenies. Malik et al.²⁵, and Puente-
261 Lelievre et al.²⁶ delved into the evolutionary relationships of this superfamily by creating a phylogenetic
262 network.

263 They employed the distance-based NeighborNet network method²⁷, utilizing distances calculated through
264 structure-based alignment methods. . Fig. 7A depicts the schematic tree built by Malik et al.²⁵, and Lelievre
265 et al.²⁶. We employed the same protein structures within this superfamily as utilized by Malik et al.²⁵ and
266 Lelievre et al.²⁶ to reconstruct the phylogeny of based on profile of energies. The dataset specifically
267 focuses on the SCOP superfamily, Ferritin-like (a.25.1) encompassing two manually curated protein
268 families: Ferritin (a.25.1.1) and Ribonucleotide Reductase-like [RNR] (a.25.1.2). The “Ferritin” family
269 contains ferritins, bacterioferritins, and Dodecameric ferritin homolog (Dps) proteins and the
270 “Ribonucleotide Reductase-like” family contains the activating subunit of class I ribonucleotide reductase
271 (RNR R2), BMM, and Fatty acids¹⁴. Following this, we computed SPE for each protein and calculated all
272 pairwise distances between SPEs. The phylogenetic tree, constructed using the phangorn package²⁸, was
273 visualized through the SplitTree software²⁹ and is presented in Fig. 7B.

274 Our results suggest that the energetic phylogenies within the ferritin-like superfamily unveil significant
275 relationships among its members, aligning with known evolutionary relationships and functional roles. In
276 line with prior investigations, a key observation is that the resulting phylogenetic tree exhibits two primary
277 branches, corresponding to two families a.25.1.1 and a.25.1.2. Thus, our methodology accurately bifurcates
278 this superfamily into two families. Delving into specifics, the family a.25.1.1 (depicted by orange color
279 triangles) further divides into four subgroups: “ferritins”, “Dps”, “Rubrerythrin”, and “Bacterioferritins”
280 indicated by distinct colors in Fig. 7B. On the other hand, the second branch related to the a.25.1.2 family
281 (dark blue triangles), despite SCOP and CATH assigning these proteins to a unified RNR-like family,
282 reveals three distinct families according to Pfam—Phenol_Hydrox (PF02332), Ribonuc_red_sm (PF00268),
283 and Fatty acid desaturase (PF03405). Our results consistently support this more detailed sequence-based
284 classification, as well as the further subdivision of the BMMs into BMMA and BMMb. The protein groupings
285 presented by Lelievre et al. in Fig. 7A are color-coded, corresponding to the colors used in Fig. 7B, C. Our
286 approach successfully reconstructed the phylogenetic tree using the energy profile. However, as shown in
287 Fig. 7C, the phylogenetic tree generated by the TM-Vec representation and cosine similarity could delineate
288 two distinct branches corresponding to two protein families but failed to predict the evolutionary
289 relationships within each protein family as proposed by Lelievre et al. The dashed line in Fig. 7B-C

290 demonstrates that both the energy model and the vector model effectively distinguish between the Ferritin
291 and Ribonucleotide reductase-like families.

292 In Fig. 7B, the energy profile model accurately orders the divergence of proteins within the Ferritin family,
293 following Lelievre et al.'s order of rubrerythrins, Ferritins, and then Dps and bacterioferritins. Conversely,
294 for the Ribonucleotide reductase-like family, the energy profile model reconstructs the proposed
295 evolutionary order of Fatty acids, RNR R2, and then BMM. In contrast, the phylogenetic tree reconstructed
296 using the TM-Vec model, while capable of differentiating the two families, does not align with the
297 evolutionary order suggested by Lelievre et al. For example, within the Ferritin family, Lelievre et al.'s model
298 posits that Dps diverged later, whereas the TM-Vec model indicates it as the first group to separate.
299 Similarly, for the Ribonucleotide reductase-like family, the TM-Vec model places the BMM proteins as the
300 earliest branch in the phylogenetic tree, whereas Lelievre et al.'s model suggests they were among the last
301 to diverge.

302

303 **3.5 Clustering of the SARS-CoV-2, SARS-CoV and 2012 MERS-CoV proteins**

304

305 Over the past two decades, Coronaviruses (CoVs) have been associated with various outbreaks, including
306 the 2002–2003 SARS-CoV outbreak, the 2012 MERS-CoV incident, and the recent COVID-19 pandemic
307 initiated by SARS-CoV-2 in late 2019. Since February 2020, a considerable number of SARS-CoV-2 protein
308 structures have been recorded in the Protein Data Bank (PDB). One pivotal viral protein, the spike
309 glycoprotein, has garnered significant attention. As a transmembrane glycoprotein, it plays a central role in
310 viral infection by facilitating host receptor binding and stands as the primary target for neutralizing antibodies
311 and vaccine design. To thoroughly investigate the structural landscape of these spike glycoproteins and
312 gain insights into their evolutionary connections, we utilized the CoV3D database
313 (<https://cov3d.ibbr.umd.edu>), a comprehensive repository containing diverse coronavirus protein structures
314 and their complex interactions with antibodies, receptors, and small molecules¹⁵.

315 From the CoV3D database, we curated a dataset comprising 143 spike glycoprotein structures
316 distinguished by the presence of the closed receptor binding domain (RBD) within their structure. This
317 dataset encompasses 80 chains from SARS-CoV-2, 31 chains from SARS-CoV, and 32 chains from MERS-
318 CoV. To scrutinize the structural variations and relationships among these spike glycoproteins, we
319 generated a 210-dimensional profile of energies at both sequence and structure levels. By calculating
320 Manhattan distances between all pairs of energetic profiles, we successfully categorized the spike
321 glycoprotein structures into three distinct clusters through unsupervised clustering based on these
322 distances. These clusters correspond to the SARS-CoV, MERS-CoV, and SARS-CoV-2 viruses, offering a
323 visually informative representation of the structural and evolutionary relationships within this protein family.
324 As depicted in Fig. 8A-B, the lineage of SARS-CoV and SARS-CoV-2 is clearly distinguished from that of
325 MERS-CoV, which belongs to a distinct subgenus within the family of coronaviruses.

326 As shown in Fig. 8A-D, our methods, CPE and SPE, demonstrate superior results in clustering proteins
327 from different viruses. Based on the evolutionary history of these viruses, it is known that MERS-CoV

328 emerged first, followed by SARS-CoV, and then SARS-CoV-2. CPE and SPE accurately reveal this
329 evolutionary pattern. However, Fig. 8C and 8D reveal that RMSD and TM-Vec are not effective in accurately
330 reconstructing this pattern for certain proteins. As highlighted by the orange circles in Fig. 8C, some SARS-
331 CoV-2 proteins are positioned far from their respective groups. Similarly, SARS-CoV proteins in the tree
332 generated by the TM-Vec method are also mis grouped, as indicated by the orange circle. Additionally, this
333 method failed to correctly group proteins from MERS-CoV. In addition to clustering performance metrics,
334 our methods, CPE and SPE, are significantly more efficient in time of computation, with CPE taking only
335 0.9 seconds and SPE requiring just 3 minutes. In comparison, TM-Vec takes 89 seconds, RMSD takes 70
336 minutes, and TM-score requires a substantial 9.7 hours. This underscores the computational advantage of
337 our methods, making them both effective and efficient for clustering analyses. To compare the clustering
338 results of our methods with those obtained from RMSD, TM-score and TM-Vec, we used the Adjusted Rand
339 Index (ARI). The ARI measures the similarity between two clusterings, accounting for chance agreement,
340 with values ranging from -1 to 1, where 1 indicates perfect agreement. Fig. 8E and Table S1 show that with
341 a cut tree of 4, our sequence-based methods achieved the highest clustering performance, with an ARI of
342 0.95. TM-Vec's optimal result was at a cut tree of 5, achieving an ARI of 0.87. For structure-based methods,
343 our SPE method achieved a perfect ARI of 1 with a cut tree of 3. RMSD's best performance was at a cut
344 tree of 6 with an ARI of 0.73, while TM-score performed best at a cut tree of 4, with an ARI of 0.56. These
345 results, summarized in Table 3, underscore the effectiveness of our proposed methods in both sequence
346 and structure-based clustering analyses. Fig. 8F highlights that our methods are significantly faster than
347 others. The list of PDBIDs of the spike proteins used in this analysis is summarized in the table S2 and the
348 result of TM-score is presented in **Fig. S1**.

349

350 **3.6. Clustering of Bacteriocins**

351 Bacteriocins are peptides produced by bacteria that act as strong antibacterial agents against other,
352 typically closely related microbial species. We analyzed the bacteriocins family available in the BAGEL
353 database, those with a length larger than 30 amino acids, including a total of 689 proteins¹⁶. Detecting and
354 understanding these peptides is crucial due to their ecological importance, but their diverse sequences and
355 structures make them challenging to identify using traditional methods. To address this issue, the BAGEL
356 tool was developed in 2006, specifically designed for identifying Ribosomally synthesized and post-
357 translationally modified peptides (RiPP) and bacteriocin biosynthetic gene clusters (BGCs). BAGEL
358 categorizes bacteriocins based on size and stability into RiPPs (also defined as class I bacteriocins by
359 BAGEL), class II bacteriocins (small heat stable proteins < 10 kDa) and class III bacteriocins (large heat-
360 labile proteins > 10 kDa). As shown in Fig. 9A, our analysis revealed that profile of energy (CPE) can clearly
361 partition bacteriocins according to BAGEL annotation. Hamamsy et al.³⁰ leveraged the deep protein
362 language models to develop the TM-Vec model, which is trained on pairs of protein sequences and their
363 TM-scores. We compared CPE distances to the TM-scores of protein structures predicted by AlphaFold2
364³¹, OmegaFold³², and ESMFold³³, as well as the TM-Vec predicted by the TM-Vec model. As demonstrated
365 in Fig. 9B, the TM-score of proteins predicted by AlphaFold2, OmegaFold, and ESMFold from the same

366 class is similar to proteins from different classes. TM-Vec is effective at distinguishing between bacteriocins
367 from the same class and proteins from different classes. Although there is some overlap between TM-Vec
368 values from proteins from the same class and other classes. Our method also effectively distinguishes
369 between proteins from the same class and those from other classes in bacteriocin dataset.

370

371 **3.7. Effective Drug Combination suggestion using Energetic Signatures**

372 The identification and validation of effective drug combinations are crucial in the treatment of various
373 complex diseases, aiming to enhance therapeutic efficacy while minimizing toxicity. However, this task is
374 hindered by a combinatorial explosion resulting from the multitude of potential drug pairs. Cheng et al.
375 introduced a network-based methodology to pinpoint clinically effective drug combinations tailored to
376 specific diseases¹⁷. This approach involved assessing the network-based relationships among drug targets
377 and disease proteins within the human protein-protein interactome. By quantifying these relationships, they
378 identified clusters of drugs that exhibited correlations with therapeutic effects. The drugs within these
379 clusters targeted the same disease module but belonged to separate neighborhoods. This innovative
380 network methodology presented by Cheng et al. provides a generic and powerful means to discover
381 effective combination therapies during drug development. Disease proteins were observed to form localized
382 neighborhoods, referred to as disease modules, rather than being randomly distributed throughout the
383 interactome. To characterize the mutual relationship between two drugs and a disease module, they
384 employed the following network-based proximity measure:

$$385 \quad s_{AB} = \langle d_{AB} \rangle - \frac{\langle d_{AA} \rangle + \langle d_{BB} \rangle}{2} \quad (5)$$

386 This measure assessed the network proximity of drug-target modules A and B by comparing the mean
387 shortest distance within the interactome between the targets of each drug ($\langle d_{AA} \rangle$ and $\langle d_{BB} \rangle$) to the
388 mean shortest distance between A-B target pairs $\langle d_{AB} \rangle$. When $s_{AB} < 0$, the targets of the two drugs are
389 in the same network neighborhood; when $s_{AB} > 0$, the two targets are topologically separated.

390 The authors demonstrated that the topological relationship between two drug-target modules, as indicated
391 by s_{AB} , reflects both biological and pharmacological relationships. They also showed that the network
392 proximity (s_{AB}) of drug-drug pairs in the human interactome correlates with chemical, biological, functional,
393 and clinical similarities. This led them to conclude that each drug-target module possesses a well-defined
394 network-based footprint. If the footprints of two drug-target modules are topologically separated, the drugs
395 are considered pharmacologically distinct. Conversely, if the footprints overlap, the magnitude of the
396 overlap indicates the strength of their pharmacological relationship. A closer network proximity of targets in
397 a drug pair suggests higher similarities in their chemical, biological, functional, and clinical profiles.

398 Here, we used the following separation measure, denoted by E_{AB} , based on similarity between profiles of
399 energies of protein targets:

$$400 \quad E_{AB} = \langle d_{AB} \rangle - \frac{\langle d_{AA} \rangle + \langle d_{BB} \rangle}{2} \quad (6)$$

401 where

402

$$\langle d_{AB} \rangle = \frac{1}{\|A\| + \|B\|} \left(\sum_{a \in A} \min_{b \in B} d(a, b) + \sum_{b \in B} \min_{a \in A} d(a, b) \right)$$

403 and $d(a, b)$ represents the Manhattan distance between the energy profiles of proteins a and b.
 404 Fig. 10 depicts the correlation between s_{AB} values, as computed by Cheng et al.¹⁷, for a set of 65
 405 antihypertensive drugs exhibiting complementary exposure to the hypertension disease module, and the
 406 corresponding E_{AB} . The results demonstrate a strong correlation between s_{AB} and E_{AB} , suggesting that the
 407 energy profile holds promise for predicting drug combinations. It is important that our approach only requires
 408 protein sequences and is significantly faster than computing the shortest path in a protein-protein interaction
 409 network.

410

411 Discussion

412

413 The continuous growth of protein databases highlights the importance of understanding their functional
 414 characteristics. It's widely recognized that proteins with similar structures often perform similar functions.
 415 Additionally, there's a common belief that proteins with similar structures also share similar energy levels.
 416 Therefore, our study aims to pioneer a new approach by directly linking protein energy landscapes to their
 417 functional attributes. By investigating this relationship, we seek to uncover new insights into how protein
 418 structure, energetics, and biological activity are interconnected. Knowledge-based potentials are energy
 419 functions derived from known protein structures. In our study, we used the DBNI potential function¹⁸ to
 420 calculate the energy between pairs of amino acids, generating energy profiles based on both sequence and
 421 three-dimensional structure. A significant achievement of our study is the high correlation observed
 422 between energy estimates derived from sequence and those from structural data, allowing for the derivation
 423 of energy profiles based solely on sequence information, which enables fast and accurate computational
 424 analysis. However, it's worth noting that the reliance on knowledge-based potentials is dependent on known
 425 protein structures, potentially limiting the generalizability of results to proteins with varied structural
 426 characteristics or those are underrepresented in existing databases. Furthermore, despite the promising
 427 correlation between energy estimates derived from sequence and structural data, it is possible that there
 428 are complexities in accurately capturing the entirety of protein energetics solely from sequence information,
 429 which could affect the reliability of the resulting energy profiles. To address these issues, one possible
 430 option is to adjust the energy profile, such as through reweighting, to specific applications, such as protein
 431 remote homology detection or drug-target affinity prediction.

432 We employed Uniform Manifold Approximation and Projection (UMAP) to visualize energy profiles at both
 433 sequence and structural levels derived from protein domains within the ASTRAL database, revealing their
 434 capacity to distinguish proteins across various hierarchical levels, including class, fold, superfamily, and
 435 family. Notably, the Manhattan distance between energy profiles serves as a measure of dissimilarity,
 436 eliminating the necessity for structural or sequence alignment in protein comparison and resulting in
 437 significantly faster computational analyses, as demonstrated in Table 2. The comparison table highlights

438 notable differences in both accuracy and computational efficiency among the methods evaluated. The
439 profile of energy (CPE) method demonstrates a remarkable accuracy of 97%, significantly surpassing other
440 methods such as GR-Align, RMSD, and TM-Score, which range from 59.2% to 81.5%. This indicates that
441 the CPE method excels in accurately distinguishing between protein structures at different superfamilies,
442 showcasing its superiority in capturing structural dissimilarities effectively. In terms of computational
443 efficiency, the CPE method stands out as the most time-efficient, requiring a mere 3 minutes for processing.
444 In contrast, traditional methods like RMSD and TM-Score demand significantly longer computational times,
445 ranging from 1 hour to over 9 hours. For instance, the CPE method is approximately 20 times faster than
446 RMSD and 180 times faster than TM-Score. This stark difference underscores the efficiency of the CPE
447 method, particularly in time-sensitive scenarios or large-scale protein structure comparison tasks.

448 Our method's efficacy was further assessed by comparing its results with structural dissimilarity metrics
449 such as RMSD, TM-Score, and GR-align in classifying proteins across five distinct SCOP superfamilies,
450 showcasing its superior accuracy and computational efficiency. Particularly challenging is elucidating
451 evolutionary relationships among superfamilies beyond the "twilight zone," where sequence similarity alone
452 proves inadequate for meaningful analysis. To address this, we examined energy profiles to reconstruct a
453 phylogenetic network of the Ferritin-like superfamily, incorporating proteins from the twilight zone. Our
454 analysis, consistent with previous studies by Lundin et al¹⁴ and Malik et al.²⁵, unveiled substantial and
455 valuable evolutionary signal preserved within energy profiles, indicating their potential as representative
456 indicators of protein structure. Moreover, we examined the structural attributes of spike glycoproteins
457 among three coronaviruses—SARS-CoV, MERS-CoV, and SARS-CoV-2—using a 210-dimensional
458 energy profile combined with Manhattan distances. This study successfully grouped these proteins into
459 specific clusters corresponding to each virus, offering insights into their structural and evolutionary
460 relationships. Additionally, our inquiry extended to 689 proteins within the bacteriocins family,
461 encompassing various sizes and stability levels sourced from the BAGEL database. By employing the
462 energy profile (CPE), we effectively distinguished bacteriocins according to BAGEL classifications,
463 showcasing the usefulness of this method in protein classification, particularly in scenarios where proteins
464 exhibit differing stabilities. Comparative analysis involving TM-scores from a range of prediction models
465 emphasized the effectiveness of our approach in differentiating proteins within and across classes, thereby
466 providing valuable insights into bacteriocins. In summary, our findings underscore the valuable insights
467 offered by energy profiles across structural, functional, and evolutionary scales.

468 One of the significant applications of assessing protein similarity lies in quantifying the proximity between
469 two drugs based on their protein targets. When the protein targets of two drugs exhibit similarity, it is
470 reasonable to anticipate similarities in the drugs themselves. Our method, capable of quantifying the
471 dissimilarity between two proteins, potentially encodes functional information that can be leveraged to
472 gauge the similarity between two drugs according to their protein targets. Comparative analysis with a study
473 conducted by Cheng et al.¹⁷ demonstrates a notable correlation between our results, derived solely from
474 protein sequence data, and theirs, obtained using protein-protein interaction data. It is worth reiterating that

475 our method boasts remarkable speed compared to conventional approaches. By providing a rapid yet
476 effective means of assessing protein similarity, our method offers promising implications for drug discovery
477 and development, facilitating the identification of potential drug candidates with similar protein targets. This
478 underscores the significance of leveraging computational methods to expedite drug discovery processes
479 while maintaining robustness and accuracy. In conclusion, our research introduces the energy profile as an
480 innovative feature set containing significant functional insights that can be utilized to represent proteins
481 within machine learning methodologies for predicting protein function, drug-target interactions, and drug
482 combination outcomes.

483 In our investigation, we examined the energy profile surrounding protein drug targets and demonstrated a
484 strong correlation between our scoring system and that derived from protein-protein interaction networks.
485 It's important to acknowledge that while a more sophisticated computational approach and experimental
486 validation are crucial in drug combination study, these aspects fall beyond the purview of our manuscript.
487 Moreover, while our method bears significant implications for drug discovery and development, its efficacy
488 might be limited by the availability and quality of protein sequence and structural data, as well as the
489 inherent complexity of drug-target interactions. Therefore, it is imperative for independent research
490 endeavors to address this crucial aspect and offer comprehensive insights into the practical application of
491 our approach in real-world therapeutic contexts.

492

493 **Methods**

494 A non-redundant structural dataset of 6944 protein chains was culled by PISCES from PDB with pairwise
495 sequence identity < 50%, resolution < 1.6 Å, R-factor < 0.25, protein length > 40 and < 1000 residues. These
496 proteins were applied to train and calculate the knowledge-based potential function ³⁴.

497

498 **Pairwise distance-dependent knowledge-based potential.**

499 Knowledge-based potentials are derived from databases of known protein structures. Various potential
500 functions, such as distance-dependent, dihedral angles, and accessible surface energies leverage
501 information from known protein structures to estimate energies of pairwise interactions ⁸. In our study, we
502 employed the potential function in which the contact between atoms identified using the tessellation method
503 ¹⁸. To obtain nearest neighbors in the protein structure, all amino acids in a protein chain were represented
504 by heavy atoms and a Delaunay tessellation of the resulting point set was computed using Qhull ³⁵. Two
505 atoms were defined to be in contact if they are two vertices of an edge in a simplex; and therefore, they are
506 not shielded from contact by other atoms. The distance between any two atoms was divided into 30 distance
507 shells, starting at 0.75 Å, with a distance shell equal to 0.5 Å in width to extract the knowledge-based
508 potential. As a result, atoms separated by less than the 6 angstroms interact. There is no direct interaction
509 between two atoms when there is a third atom between two close atoms. All pairwise occurrences outside
510 of this range were excluded. All atom pairs in a structure were considered except those that belonged to
511 the same residue. The study considered 167 atom types by treating all nonhydrogen atoms as having
512 different atom types when they are in different amino acid residues.

513 The energy between the two atoms i and j at distance d, is calculated as follows:

514

$$\Delta E^{ij}(d) = RT \left[\ln(1 + M_{ij}\sigma) - \ln\left(1 + M_{ij}\sigma\left(\frac{f_{ij}(d)}{f_{xx}(d)}\right)\right) \right] \quad (1)$$

515 where RT is constant and equal to 0.582 kcal/mole. M_{ij} is the number of observations for atomic
516 pair i and j , $f_{ij}(d)$ is the relative frequency of occurrence for i and j in distance class d , $f_{xx}(d)$ is the relative
517 frequency of occurrence for all atomic pairs in distance shell d , and σ is the weight given to each
518 observation. As discussed by Sippl⁷, it was assumed that $\sigma = 0.02$.

519 The potential energy associated with the interaction of residues A and B denoted by $\Delta E(A, B)$ is estimated
520 by summing the pairwise potentials between the atoms of each of these residues as follows:

521

$$\Delta E(A, B) = \sum_{i \in A, j \in B} \Delta E^{ij}(d) \quad (2)$$

522 which the sum is over all pairs of atoms in contact with the Delaunay triangulation method.

523 Given that there are 210 unique amino acid-amino acid interaction types among these 20 amino acids, the
524 total number of unique $\Delta E(A, B)$ values are 210. As a result, we create a 210-dimensional vector to
525 represent distance-dependent energy interactions between residues, with each dimension representing the
526 energy interaction between specific pairs of amino acid types. We call this 210-dimensional vector as the
527 **Structural Profile of Energy (SPE)** of a protein structure.

528

529 The pairwise energy content estimated from amino acid composition.

530 The knowledge-based potential function discussed in the previous section relies on having the three-
531 dimensional structure of a protein. Nevertheless, it's worth noting that the three-dimensional structures of
532 numerous proteins have not yet been determined experimentally. Interactions between pairs of atoms are
533 identified using Delaunay tessellation. For each protein S in the training set, e_i^S denotes the energy of
534 interactions between all amino acid residues of type i and all other amino acids in the protein structure S .
535 We estimated the value of \widehat{e}_i^S using the following expression:

536

$$\widehat{e}_i^S = N_i^S \sum_{j=1}^{20} P_{ij} n_j^S \quad (3)$$

537 where, N_i^S represents the frequency of amino acid type i in the structure S , and $n_j^S = \frac{N_j^S}{L}$, is calculated as
538 the ratio of N_j^S to the total number of amino acids in S , denoted by L , and P is the energy predictor matrix,
539 delineating the dependence of amino acid i 's energy on the j th element within the amino acid
540 composition. The parameters of the corresponding row of the matrix P are obtained by minimizing the
541 function

542

$$Z_i = \sum_S (e_i^S - \widehat{e}_i^S)^2 \quad (4)$$

543 Letting $\frac{\partial Z_i}{\partial P_{ij}} = 0$ for all P_{ij} leads to a set of linear equations which are solved for each amino acid type by
544 using Symbolic Math Toolbox in MATLAB.

545 For each pair of amino acid types i and j , we used the following equation to estimate the energy E_{ij} based
546 on amino acid sequence composition:

547

$$E_{ij} = n_i \sum_j^{20} P_{ij} n_j \quad (5)$$

548 where P is the energy predictor matrix estimated using equation 4. As a result, we create a 210-dimensional
549 vector to represent energy between amino acid types using amino acid composition. We call this 210-
550 dimensional vector as the **Compositional Profile of Energy (CPE)** of a protein sequence. The profile of
551 energies is normalized based on protein length.

552

553 Analysis Tools and Packages

554

555 All computational analyses were conducted using the versatile R programming language (www.r-project.org), with the utilization of various specialized packages tailored for specific tasks. Below is an
556 overview of the packages and tools employed throughout our analysis:

558 The BIO3D software was used to read and analyze PDB files ³⁶. The “geometry” package was used to
559 implement the Quickhull algorithm to find direct contacts and nearest neighbors of atoms in pdb files using
560 the Delaunay tessellation method (<https://cran.r-project.org/web/packages/geometry/index.html>). The kNN,
561 and RF classification algorithms were implemented using the “random Forest”, and the “caret” package ³⁷,
562 ³⁸. Figures were generated using the ggplot2 package ³⁹. TM-Vec representations were generated by
563 configuring `tm_vec_model_cpnt` to `tm_vec_cath_model`. The functions can be accessed at
564 (<https://github.com/tymor22/tm-vec/tree/master>). TM-scores were calculated using `tm_align` from the
565 tmtools 0.1.1 module.

566

567 **Data and code availability.** The data that support the findings of this study and all code for data analysis
568 are openly available at:

569 <https://github.com/mirzaie-mehdi/ProteinEnergyProfileSimilarity>

570

571 **Funding:** P.Ch. and M.M. were supported from the grants of J.-O.A - Academy of Finland (grants no.
572 297727 and 350678), Sigrid Juselius Foundation, ERA-NET NEURON grant nr 352077, Helsinki Institute
573 of Life Science Research Fellow, and by European Research Council (ERC, grant no. 724922).

574

575

576

577

578

579

580

581 **References**

- 582
- 583 1. Sayers EW, *et al.* Database resources of the national center for biotechnology information.
584 *Nucleic acids research* **49**, D10 (2021).
- 585
- 586 2. Altschul SF, *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database search
587 programs. *Nucleic acids research* **25**, 3389-3402 (1997).
- 588
- 589 3. Kilinc M, Jia K, Jernigan RL. Improved global protein homolog detection with major gains in
590 function identification. *Proceedings of the National Academy of Sciences* **120**, e2211823120
591 (2023).
- 592
- 593 4. Quan Y, Xiong Z-K, Zhang K-X, Zhang Q-Y, Zhang W, Zhang H-Y. Evolution-strengthened
594 knowledge graph enables predicting the targetability and druggability of genes. *PNAS nexus* **2**,
595 pgad147 (2023).
- 596
- 597 5. Du Z, Ding X, Hsu W, Munir A, Xu Y, Li Y. pLM4ACE: A protein language model based predictor
598 for antihypertensive peptide screening. *Food Chemistry* **431**, 137162 (2024).
- 599
- 600 6. Wang Y, *et al.* RNAincoder: a deep learning-based encoder for RNA and RNA-associated
601 interaction. *Nucleic Acids Research*, gkad404 (2023).
- 602
- 603 7. Sippl MJ. Boltzmann's principle, knowledge-based mean fields and protein folding. An approach
604 to the computational determination of protein structures. *Journal of computer-aided molecular
605 design* **7**, 473--501 (1993).
- 606
- 607 8. Mirzaie M, Sadeghi M. Knowledge-based potentials in protein fold recognition. *Archives of
608 Advances in Biosciences* **1**, (2010).
- 609
- 610 9. Mirzaie M, Eslahchi C, Pezeshk H, Sadeghi M. A distance-dependent atomic knowledge-based
611 potential and force for discrimination of native structures from decoys. *Proteins: Structure,
612 Function, and Bioinformatics* **77**, 454-463 (2009).
- 613
- 614 10. Dosztanyi Z, Csizmok V, Tompa P, Simon I. The pairwise energy content estimated from amino
615 acid composition discriminates between folded and intrinsically unstructured proteins. *Journal
616 of molecular biology* **347**, 827--839 (2005).
- 617
- 618 11. Sillitoe I, *et al.* CATH: increased structural coverage of functional space. *Nucleic acids research*
619 **49**, D266-D273 (2021).
- 620

- 621 12. Lo Conte L, Ailey B, Hubbard TJ, Brenner SE, Murzin AG, Chothia C. SCOP: a structural
622 classification of proteins database. *Nucleic acids research* **28**, 257-259 (2000).
- 623
- 624 13. Fox NK, Brenner SE, Chandonia J-M. SCOPE: Structural Classification of Proteins—extended,
625 integrating SCOP and ASTRAL data and classification of new structures. *Nucleic acids research*
626 **42**, D304-D309 (2014).
- 627
- 628 14. Lundin D, Poole AM, Sjberg B-M, Hgbom M. Use of structural phylogenetic networks for
629 classification of the ferritin-like superfamily. *Journal of Biological Chemistry* **287**, 20565--20575
630 (2012).
- 631
- 632 15. Gowthaman R, Guest JD, Yin R, Adolf-Bryfogle J, Schief WR, Pierce BG. CoV3D: a database of
633 high resolution coronavirus protein structures. *Nucleic acids research* **49**, D282-D287 (2021).
- 634
- 635 16. van Heel AJ, de Jong A, Montalbán-López M, Kok J, Kuipers OP. BAGEL3: automated
636 identification of genes encoding bacteriocins and (non-)bactericidal posttranslationally modified
637 peptides. *Nucleic Acids Research* **41**, W448-W453 (2013).
- 638
- 639 17. Cheng F, Kovács IA, Barabási A-L. Network-based prediction of drug combinations. *Nature
640 communications* **10**, 1197 (2019).
- 641
- 642 18. Mirzaie M, Sadeghi M. Delaunay-based nonlocal interactions are sufficient and accurate in
643 protein fold recognition. *Proteins: Structure, Function, and Bioinformatics* **82**, 415-423 (2014).
- 644
- 645 19. Maiorov VN, Crippen GM. Significance of root-mean-square deviation in comparing three-
646 dimensional structures of globular proteins. *Journal of molecular biology* **235**, 625--634 (1994).
- 647
- 648 20. Zhang Y, Skolnick J. TM-align: a protein structure alignment algorithm based on the TM-score.
649 *Nucleic acids research* **33**, 2302-2309 (2005).
- 650
- 651 21. Hamamsy T, et al. Protein remote homology detection and structural alignment using deep
652 learning. *Nat Biotechnol*, (2023).
- 653
- 654 22. Malod-Dognin N, Pržulj N. GR-Align: fast and flexible alignment of protein 3D structures using
655 graphlet degree similarity. *Bioinformatics* **30**, 1259-1265 (2014).
- 656
- 657 23. Tian K, Zhao X, Zhang Y, Yau S. Comparing protein structures and inferring functions with a novel
658 three-dimensional Yau–Hausdorff method. *Journal of Biomolecular Structure and Dynamics*,
659 (2018).
- 660

- 661 24. Wintjens RT, Roonan MJ, Wodak SJ. Automatic classification and analysis of $\alpha\alpha$ -turn motifs in
662 proteins. *Journal of molecular biology* **255**, 235-253 (1996).
- 663
- 664 25. Malik AJ, Poole AM, Allison JR. Structural phylogenetics with confidence. *Molecular Biology and*
665 *Evolution* **37**, 2711--2726 (2020).
- 666
- 667 26. Puente-Lelievre C, *et al.* Tertiary-interaction characters enable fast, model-based structural
668 phylogenetics beyond the twilight zone. *bioRxiv*, 2023.2012.2012.571181 (2023).
- 669
- 670 27. Bryant D, Moulton V. Neighbor-net: an agglomerative method for the construction of
671 phylogenetic networks. *Molecular biology and evolution* **21**, 255-265 (2004).
- 672
- 673 28. Schliep KP. phangorn: phylogenetic analysis in R. *Bioinformatics* **27**, 592-593 (2011).
- 674
- 675 29. Huson DH, Bryant D. Application of Phylogenetic Networks in Evolutionary Studies. *Molecular*
676 *Biology and Evolution* **23**, 254-267 (2005).
- 677
- 678 30. Hamamsy T, *et al.* Protein remote homology detection and structural alignment using deep
679 learning. *Nature biotechnology*, 1-11 (2023).
- 680
- 681 31. Jumper J, *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583-
682 589 (2021).
- 683
- 684 32. Wu R, *et al.* High-resolution de novo structure prediction from primary sequence. *bioRxiv* 2022.
685 *Google Scholar*.
- 686
- 687 33. Lin Z, *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language
688 model. *Science* **379**, 1123-1130 (2023).
- 689
- 690 34. Wang G, Dunbrack Jr RL. PISCES: a protein sequence culling server. *Bioinformatics* **19**, 1589-1591
691 (2003).
- 692
- 693 35. Barber CB, Dobkin DP, Huhdanpaa H. The quickhull algorithm for convex hulls. *ACM Transactions*
694 *on Mathematical Software (TOMS)* **22**, 469-483 (1996).
- 695
- 696 36. Grant BJ, Rodrigues AP, ElSawy KM, McCammon JA, Caves LS. Bio3d: an R package for the
697 comparative analysis of protein structures. *Bioinformatics* **22**, 2695-2696 (2006).
- 698
- 699 37. Breiman L. Manual on setting up, using, and understanding random forests v3. 1. *Statistics*
700 *Department University of California Berkeley, CA, USA* **1**, 3-42 (2002).

- 701
702 38. Kuhn M, *et al.* Package ‘caret’. *The R Journal* **223**, (2020).
- 703
704 39. Wickham H. *ggplot2. Wiley interdisciplinary reviews: computational statistics* **3**, 180-185 (2011).
- 705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727

728 **Acknowledgments:** The authors would like to thank Vilma Iivanainen, Elina Nagaeva, and Sakari Hietanen
729 for reading the manuscript and providing valuable feedback.

730

731 **Author Contributions:** M.M. designed and supervised the research; M.M. and P.Ch. analyzed data; M.M.,
732 P.Ch., and J.-O.A. wrote the paper. All authors have read and agreed to the published version of the
733 manuscript.

734 **Competing Interest Statement:** The authors declare no conflict of interest.

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753 **Figure Legends and Tables**

754

755

756 **Fig. 1 | Sequence-Structure relationship.** A) The correlation between total energy estimates derived from
757 protein structure and sequence for protein domains within ASTRAL40(left) and ASTRAL95(right) data sets.
758 B) The correlation between the distances of profile of energy estimated from sequence (CPE) and structure
759 (SPE) for all pairs of domains in ASTRAL40(left) and ASTRAL95(right).

760

761 **Fig. 2 | Energy Distribution in Protein Domain Structural Classes.** The distribution of normalized total
762 energy in protein domains from ASTRAL40 and ASTRAL95 datasets based on protein structure (A) and
763 sequence (B) across various structural scope classes. In the ASTRAL40 dataset, there are 2644, 3059,
764 4463, and 3653 protein domains in the all-alpha, all-beta, alpha+beta, and alpha/beta classes, respectively.
765 Similarly, in the ASTRAL95 dataset, there are 3443, 10164, 9344, and 7474 protein domains in the all-
766 alpha, all-beta, alpha+beta, and alpha/beta classes, respectively.

767

768 **Fig. 3 | UMAP Visualization of Energy Profiles in All-Alpha and All-Beta Domains from ASTRAL40
769 and ASTRAL95 Datasets.** UMAP projection of SPE and CPE shows the separation of the all-alpha (green
770 point) and all-beta (pink point) proteins selected from the ASTRAL40 and ASTRAL95 datasets. A) CPE of
771 ASTRAL40, B) SPE of ASTRAL40, C) CPE of ASTRAL95, and D) SPE of ASTRAL95. Dots represent two
772 dimensional UMAP projection of SPE(CPE) for individual sequences. UMAP plots were generated by
773 parameters n_neighbors = 20 and min_dist = 0.1.

774

775 **Fig. 4 | UMAP Visualization of Energy Profiles.** The UMAP projection of Structural Energy Profiles (SPE)
776 and Compositional Energy Profiles (CPE) of protein domains from ASTRAL40 and ASTRAL95 represents
777 the structural information embedded in energy profiles across hierarchical levels of SCOP; each panel
778 includes two figures, one generated by CPE (left panel) and the other by SPE (right panel), revealing that
779 protein domains sharing the same A) fold, B) superfamily, and C) family exhibit comparable energy profile
780 patterns. The folds a.100 and a.104, superfamilies a.29.2 and a.29.3, as well as families a.25.1.0 and
781 a.25.1.2, are randomly selected for analysis, and the UMAP plots were generated using parameters
782 n_neighbors = 20 and min_dist = 0.1.

783

784 **Fig. 5 | Comparative Boxplots of Pairwise Distances among Energy Profiles in ASTRAL40 and
785 ASTRAL95.** Comparative Boxplots of Pairwise Distances among Energy Profiles in ASTRAL40 and
786 ASTRAL95, depicting A-B) intraclass distances within the all-alpha class (in purple) versus interclass
787 distances (in yellow), C-D) intraclass distances within the a.29 fold (in purple) versus distances from protein
788 domains in different folds within the all-alpha class (in yellow), and E-F) intraclass distances within the
789 a.29.3 superfamily (in purple) versus distances from protein domains in different superfamilies within the
790 fold a.29 (in yellow). Each panel presents two figures, one generated using Compositional Energy Profiles
791 (CPE, left panel) and the other using Structural Energy Profiles (SPE, right panel).

792

793 **Fig. 6 | Performance and Computational Efficiency of Protein Dissimilarity Measures.** A) Time versus
794 accuracy for the 1-NN classifier using GR-Align, RMSD, TM-score, Yau-Hausdorff distance, TM-Vec, and
795 the distance between energy profiles (CPE) as measures of protein dissimilarity. B) Running times of the
796 evaluated methods, scaled to 12 hours, with an inset zooming in on the region indicated by the dashed
797 circle. The entire circle represents 130 seconds. Each method is represented by different colors as indicated
798 in the figure legend.

799

800

801 **Fig. 7 | Phylogenetic network reconstruction for the ferritin-like superfamily.** A) Schematic view of the
802 relationships between the major ferritin-like protein families. B) The network recontacted by SPE
803 demonstrates the distinct separation (red dotted line) of two SCOP families: ferritins (SCOP id a.25.1.1),
804 which includes Bacteri, ferritins, Dps, and rubrerythrin subgroups, and the Ribonucleotide Reductase-like
805 family (SCOP id a.25.1.2), which includes BMM_alpha, BMM_beta, Fatty_acid, and RNRR2 subgroups.
806 Smaller groups are clearly distinguished. C) The network recontacted by TM-Vec representation and cosine
807 similarity.

808

809 **Fig. 8 | Clustering Analysis of Spike Glycoprotein Structures from SARS-CoV, SARS-CoV-2, and**
810 **MERS-CoV.** The dendrograms illustrate the clustering of spike glycoprotein structures from three viruses
811 SARS-CoV, SARS-CoV-2, and MERS-CoV. The clustering is based on pairwise distances of energy
812 profiles derived from A) protein structure, B) protein sequence C) RMSD D) TM-Vec. Each leaf on the
813 dendrogram is labeled with the PDB-IDs of the corresponding chains, and the leaves are color-coded to
814 represent the host virus of the spike glycoprotein structure. E) The ARI values of CPE, SPE, TM-Vec,
815 RMSD, and TM-Score scores, and F) The running time.

816

817 **Fig. 9 | UMAP Visualization and Comparison of Embeddings for Bacteriocins.** Visualization of profile
818 of energies embeddings using UMAP for 689 peptides across three classes of bacteriocins. A) The UMAP
819 projection of Compositional Energy Profiles (CPE) on bacteriocins at different classes. B) Comparison of
820 CPE distances (CPE_dis) with the TM-scores produced by running TM-align on structures predicted by
821 AlphaFold2, OmegaFold and ESMFold, and TM-Vec for 238,000 pairs of bacteriocins. CPE_dis is
822 normalized by min-max normalization

823

824

825 **Fig. 10 | Correlation between Protein-Protein Interaction Network Distances and Profile of Energies**
826 **Distances.** The correlation between separation distances estimated by protein-protein interaction network
827 (X-axis) and the distance between profiles of energies (Y-axis).

828

829 **Table 1** | The accuracy and computation time for 1-NN classifier based on GR-Align, RMSD, TM-score,
830 Yau-Hausdorff distance, TM-Vec, and the distance between profiles of energy (CPE) as a measure of
831 protein dissimilarity.

832

Method	Accuracy	Time
GR-Align	62.3%	2 min
RMSD	59.2%	1 h
TM-Score	61.5%	9 h 20 min
YH (10 Rotation)	70.8%	10 min
YH (2500 Rotation)	81.5%	4h 10 min
TM-Vec	100%	67 sec
CPE	100%	1 sec

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853 **Table 2 |** Total accuracy and F1 measure for each of the five superfamilies by 1-NN and the results of 10-
 854 Fold cross validation with random forest (RF).

855

Method	Time	Accuracy	F1 Measure				
			wigend_he lix	PH.domain- like	NTF-like	Ubiquitin-like	Immunoglobulin s
CPE (1NN)	103 Sec	0.98	0.98	0.96	0.99	0.99	0.99
CPE (RF)	103 Sec	0.99	0.97	0.97	0.99	0.99	0.99
TM_Vec (1NN)	955 Sec	0.99	0.99	1	1	0.99	0.99

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

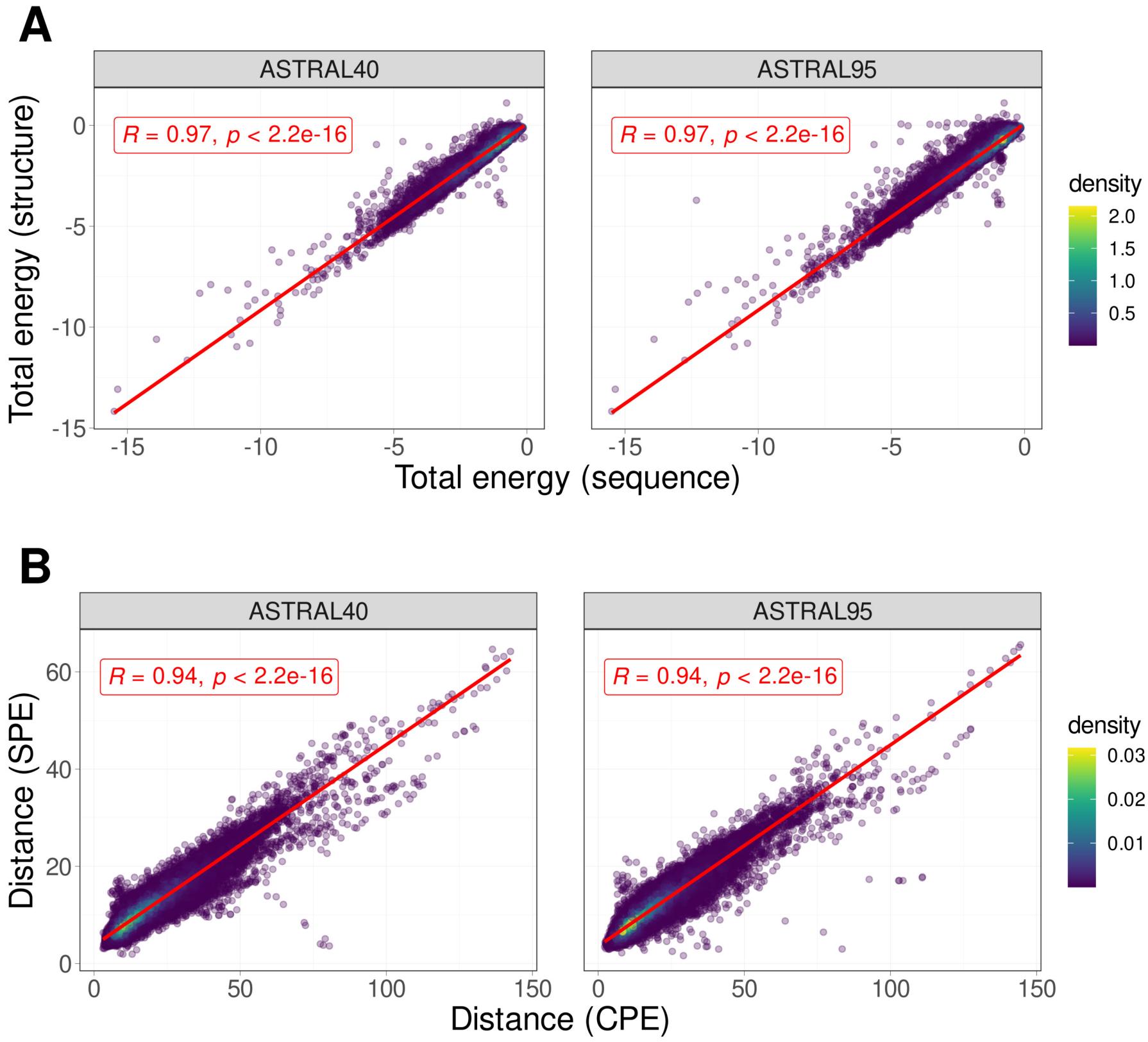
877

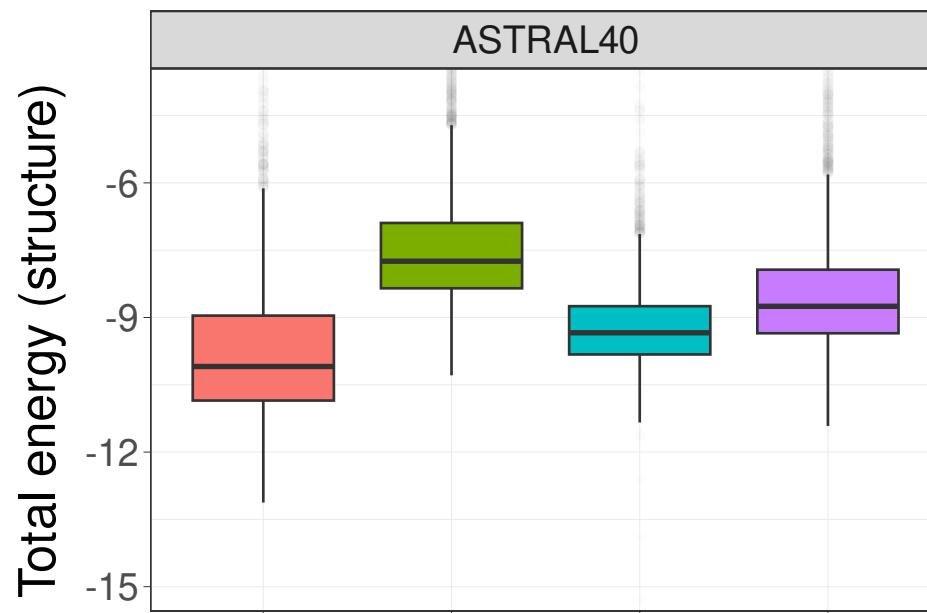
878 **Table 3** | Comparison of clustering results using Adjusted Rand Index (ARI)

Method	Type	Cut Tree	ARI
CPE	Sequence-Based	4	0.95
TM-Vec	Sequence-Based	5	0.87
SPE	Structure-Based	3	1.00
RMSD	Structure-Based	6	0.73
TM-Score	Structure-Based	4	0.56

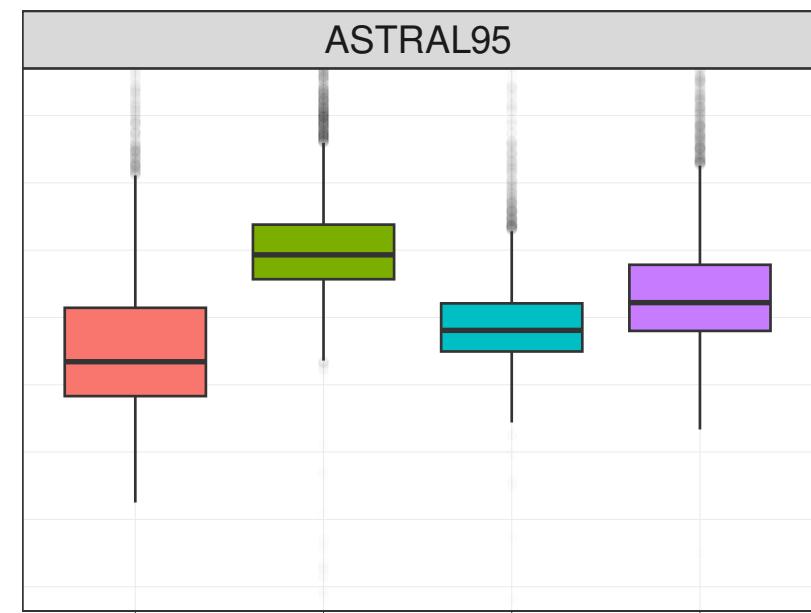
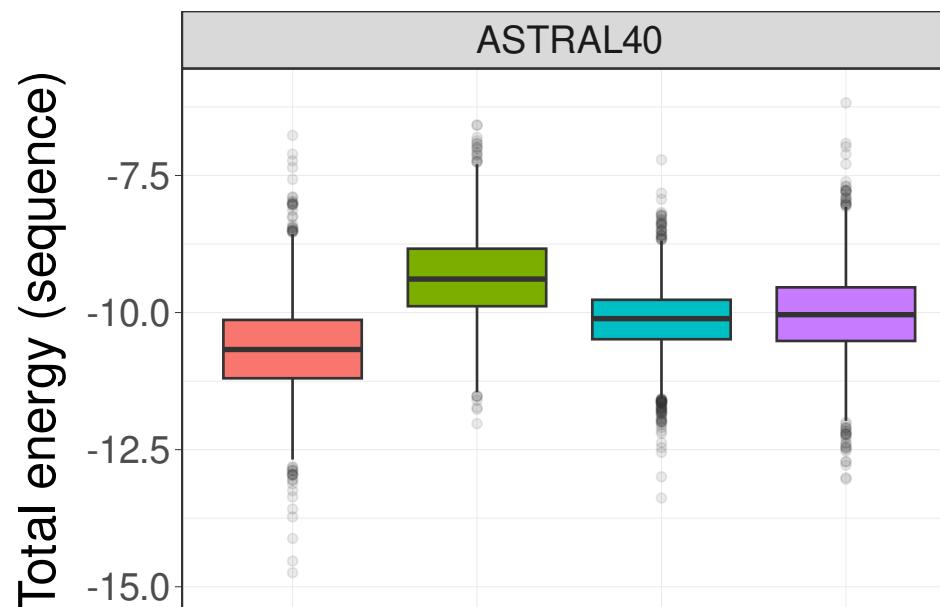
879

880

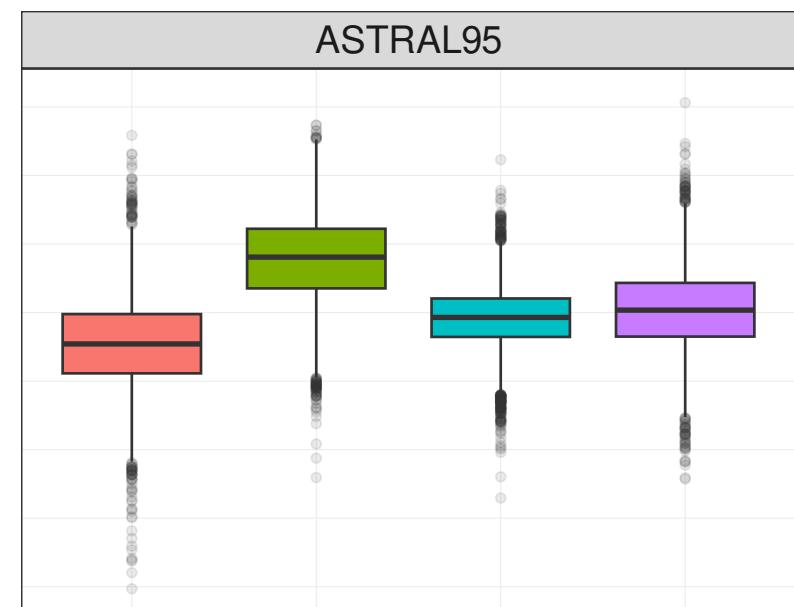


A

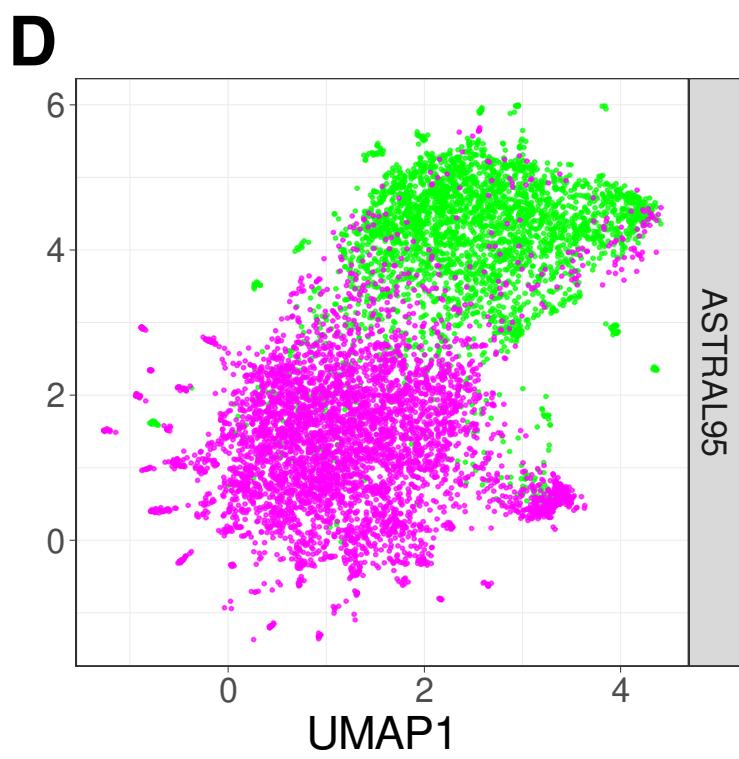
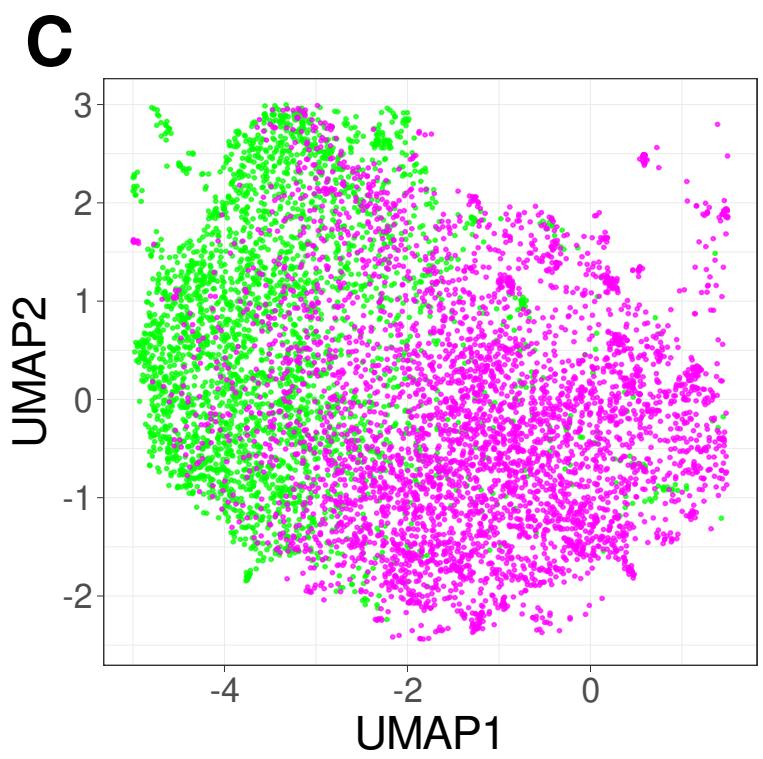
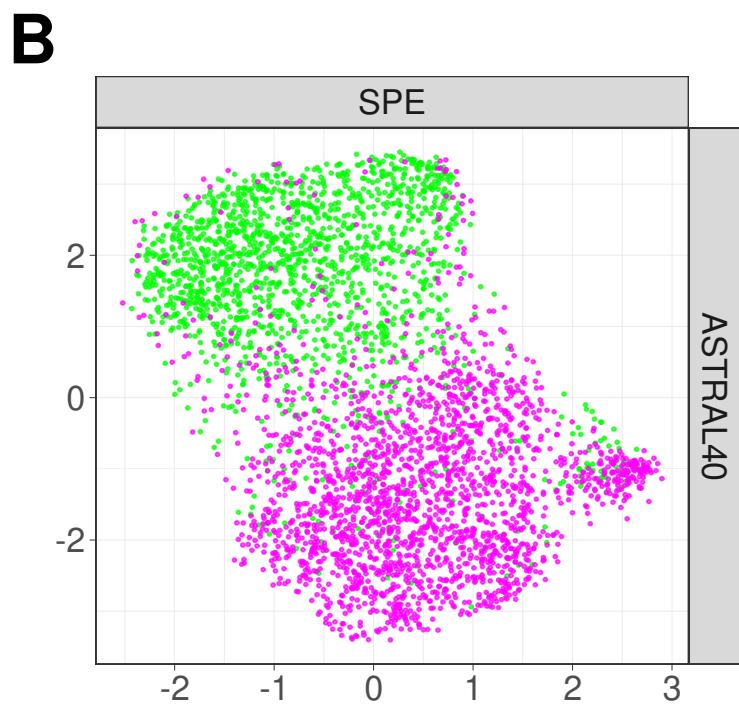
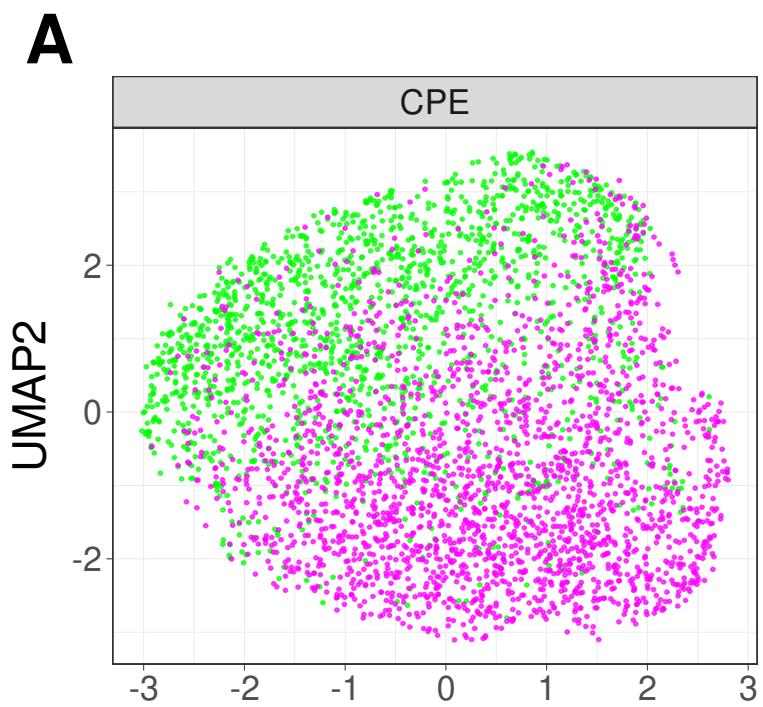
ASTRAL95

**B**

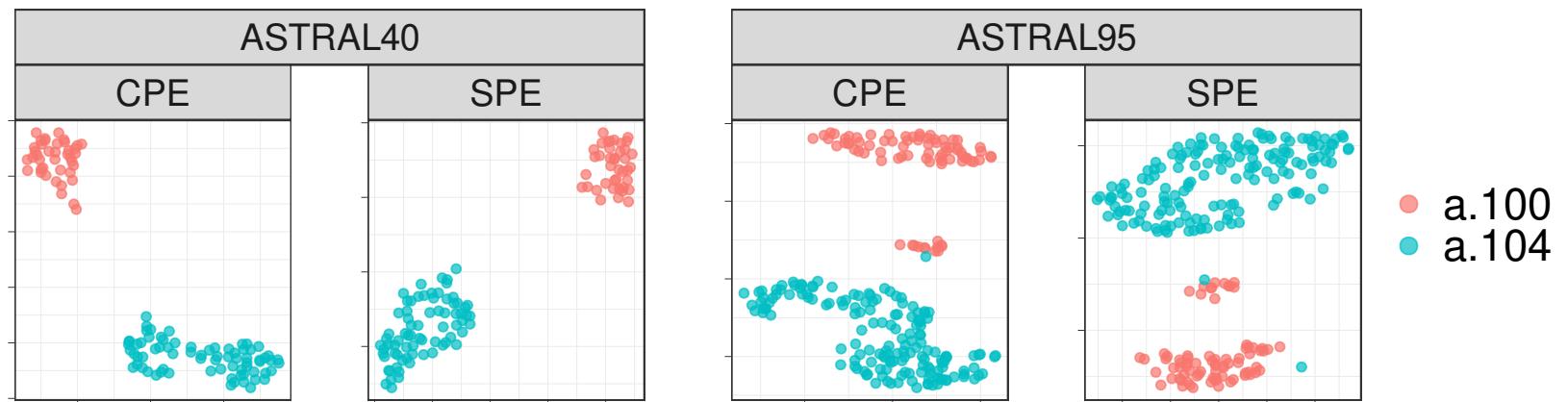
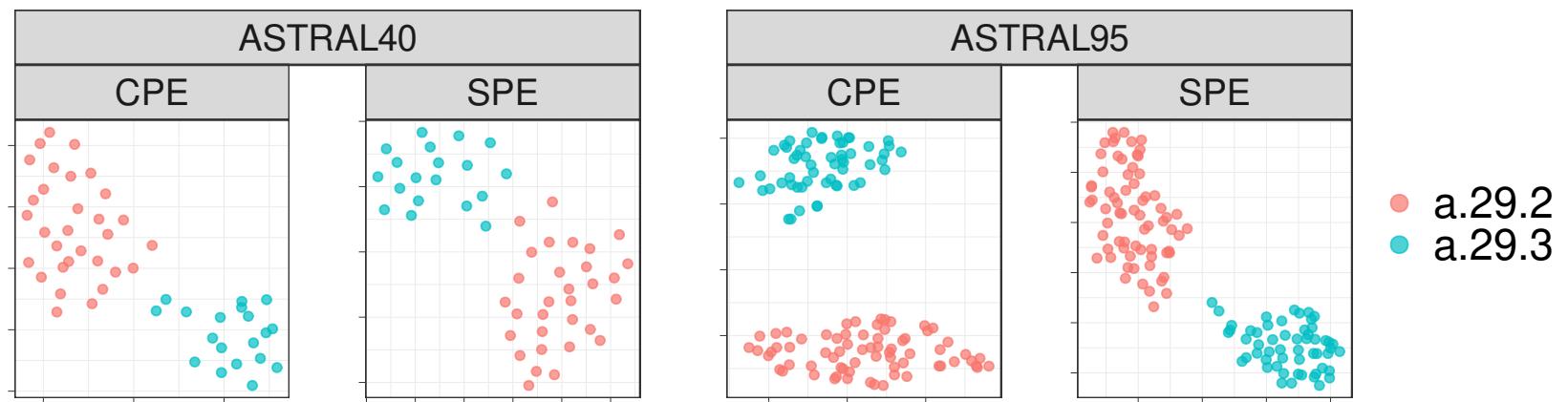
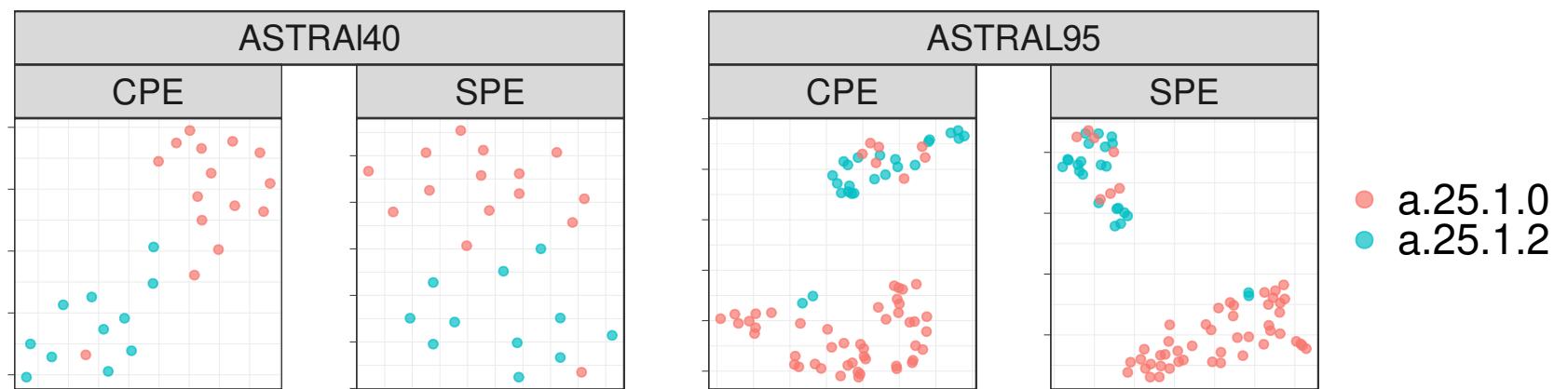
ASTRAL95

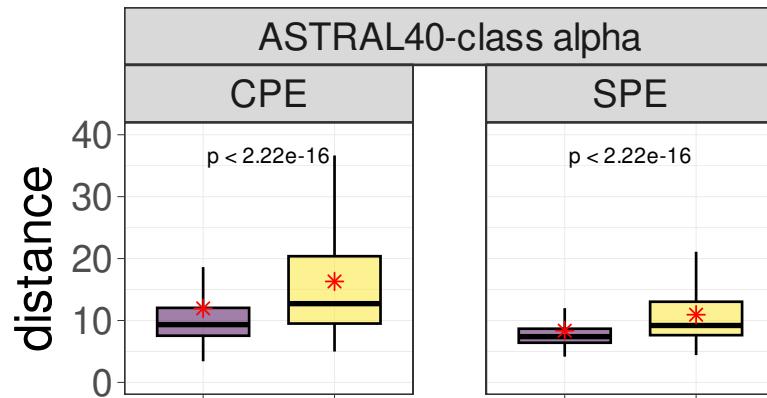
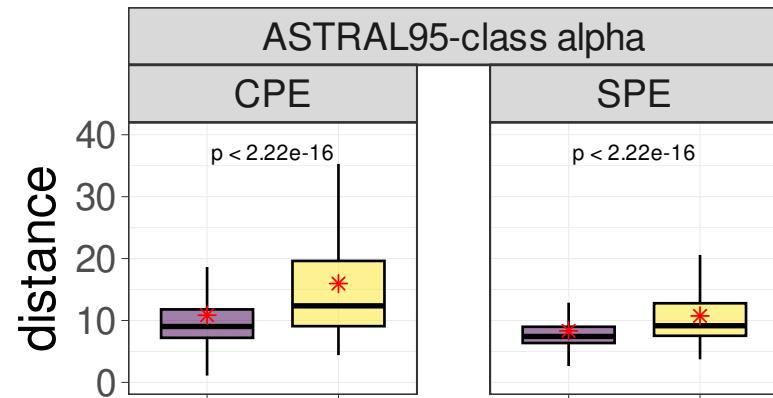
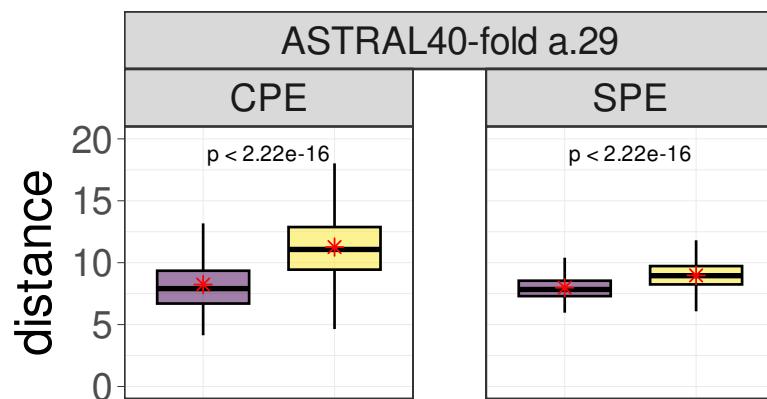
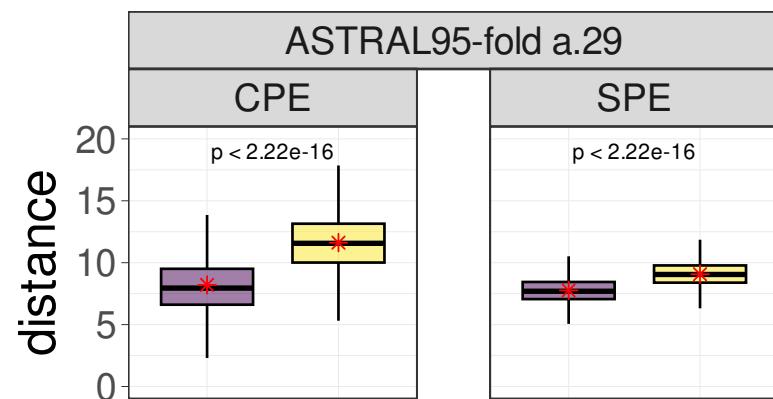
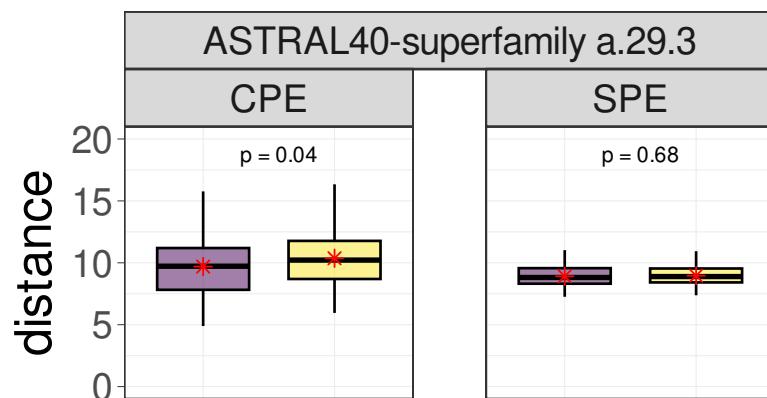
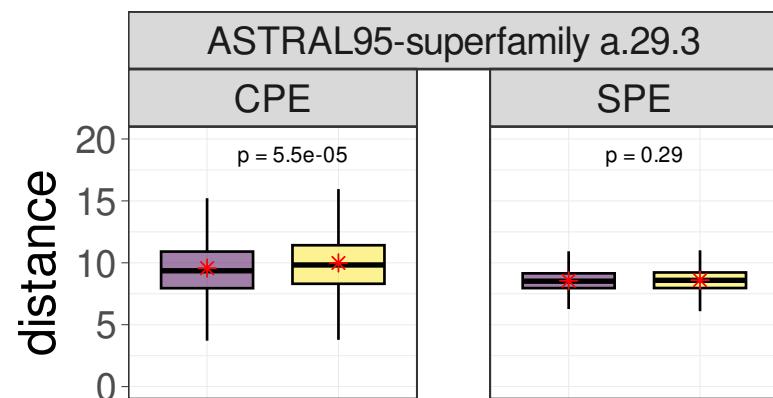


All_alpha All_beta alpha/beta alpha+beta

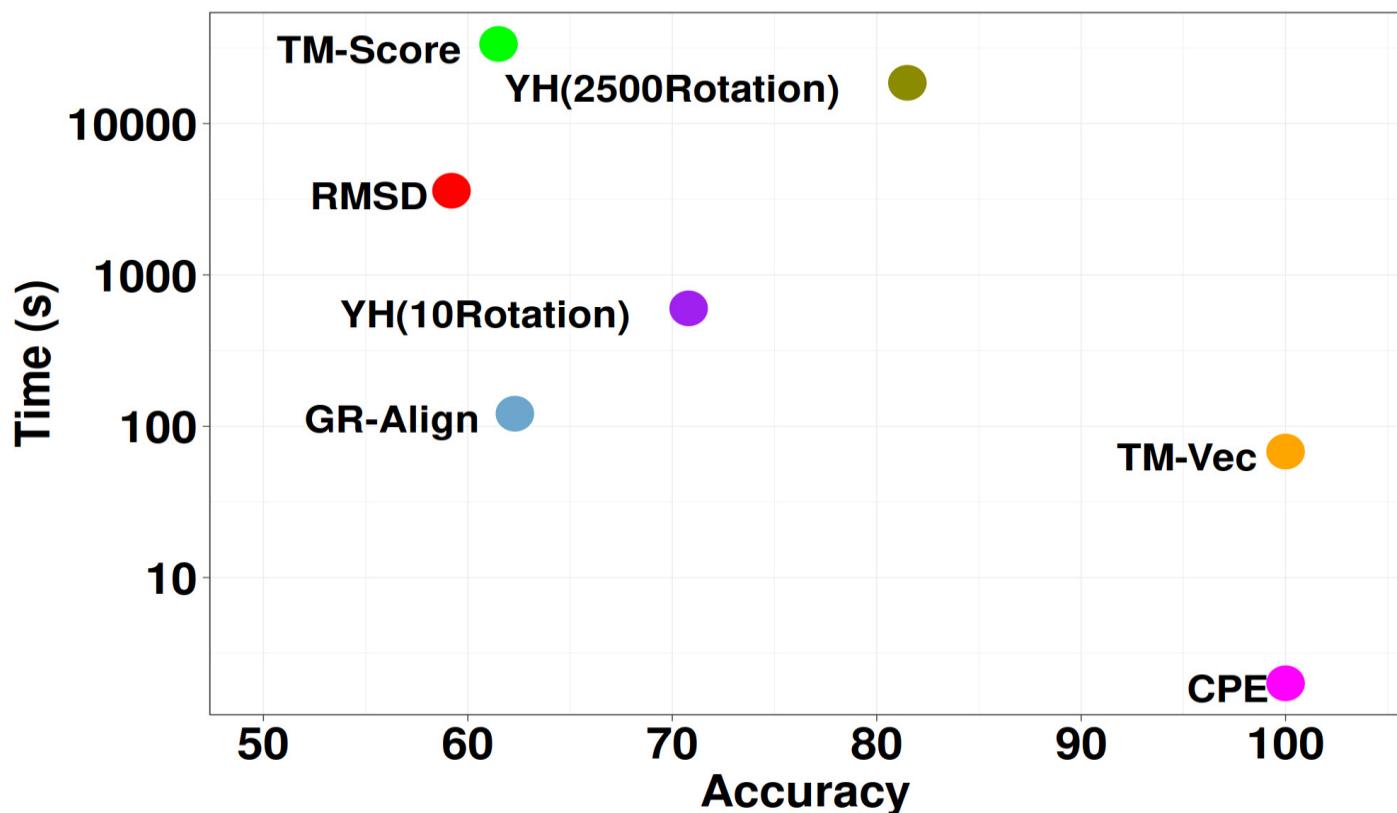
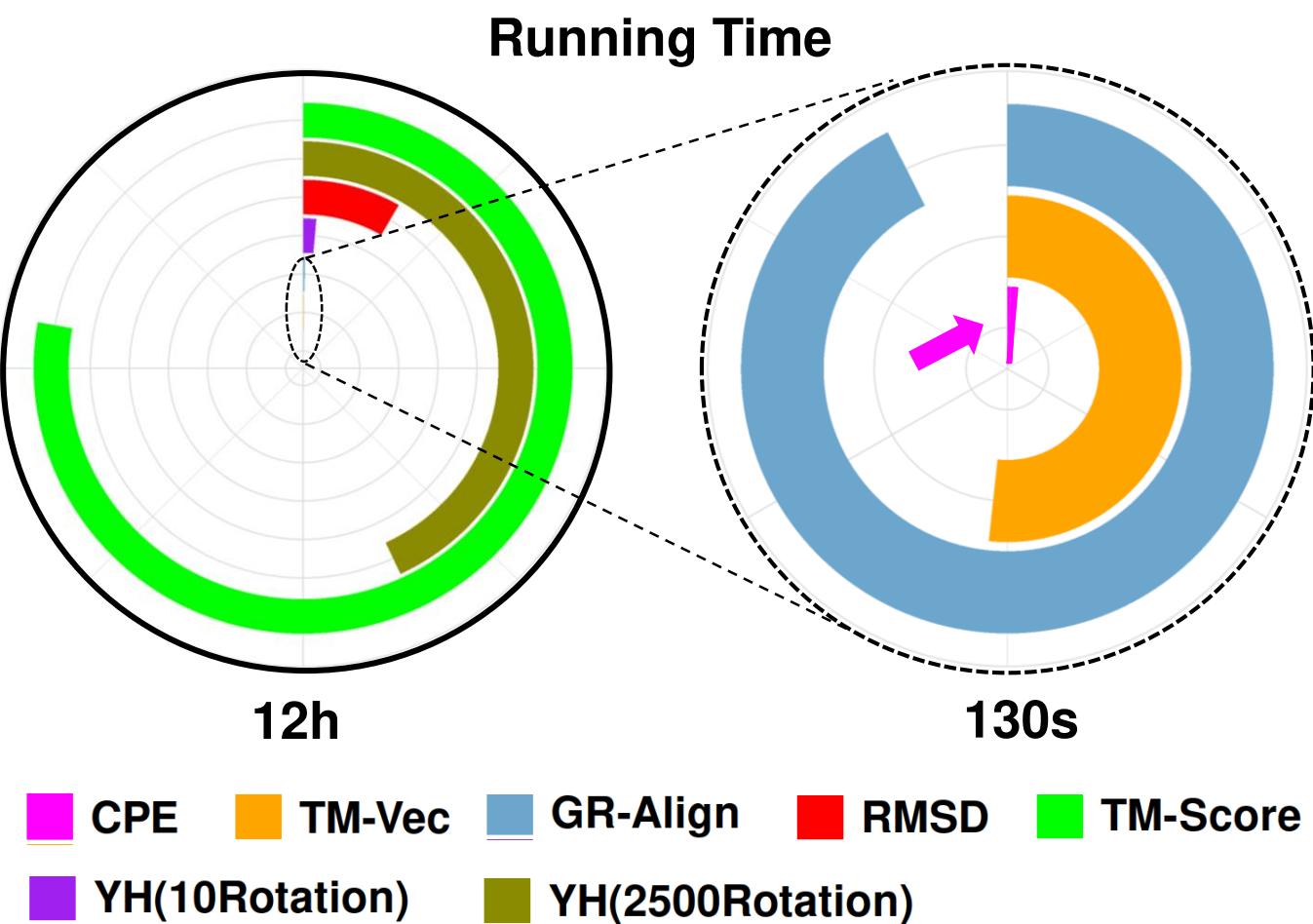


● All-alpha ● All-beta

A**B****C**

A**B****C****D****E****F**

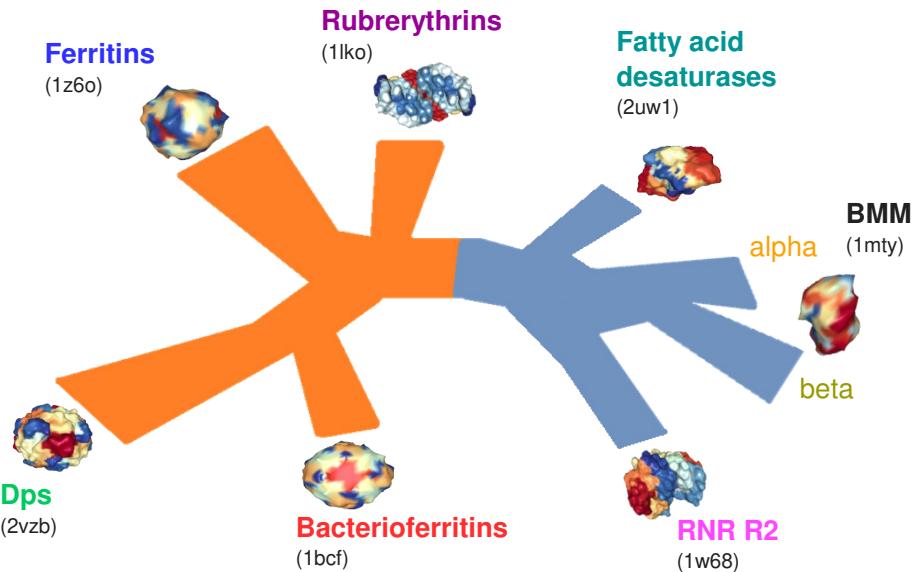
■ within ■ between

A**B**

A

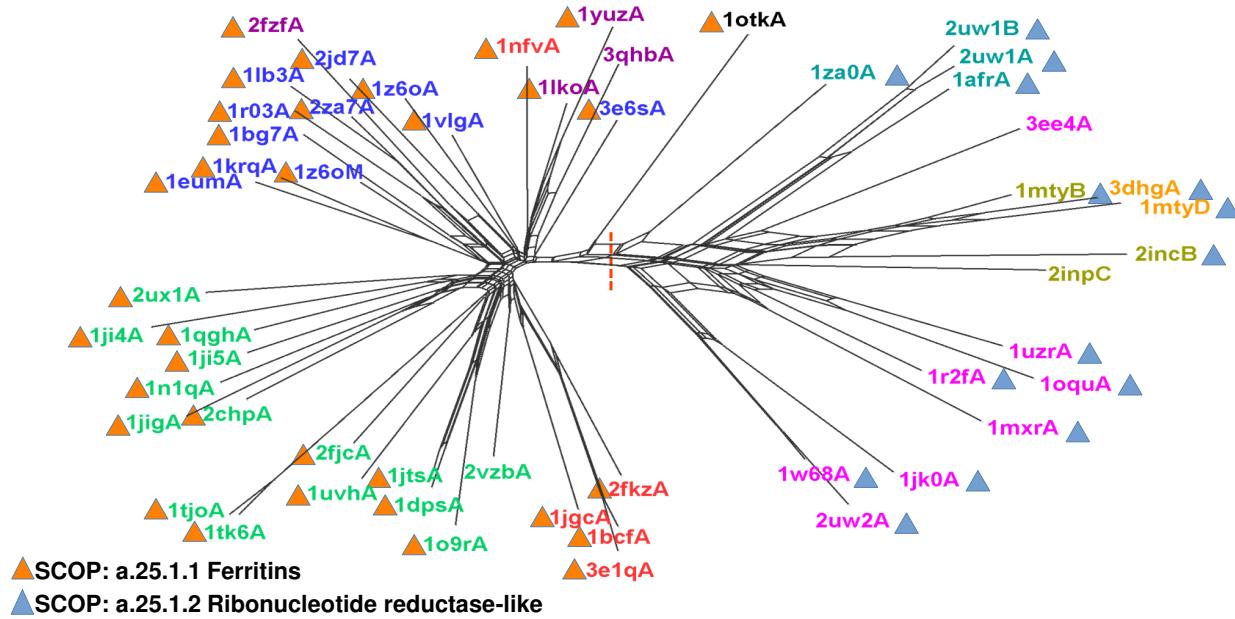
Manual Annotation

Subgroup	Pfam
Bacterioferirtin	
Dps	00210 Ferritin
Ferritin	
BMM alpha	02332
BMM beta	Phenol_hydrox
Fatty acid desaturase	03405 FA_desaturase
RNR R2	00268 Ribonuc_red_sm
Rubrerythrins	02915 Rubrerythrins
1tokA	05138 PaaA_PaaC



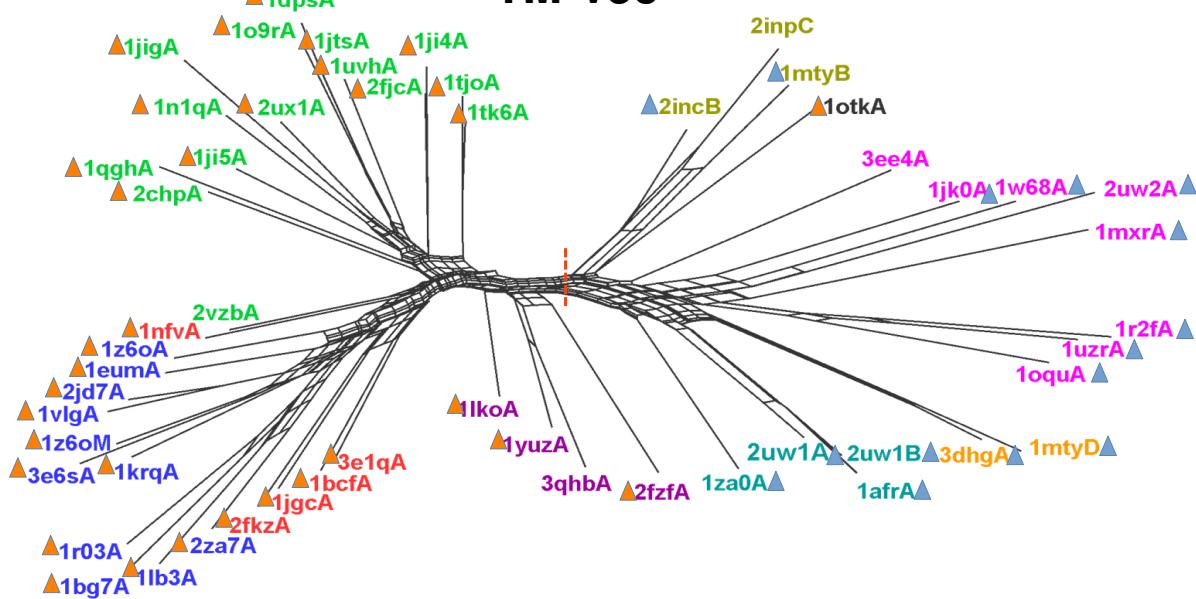
B

SPE



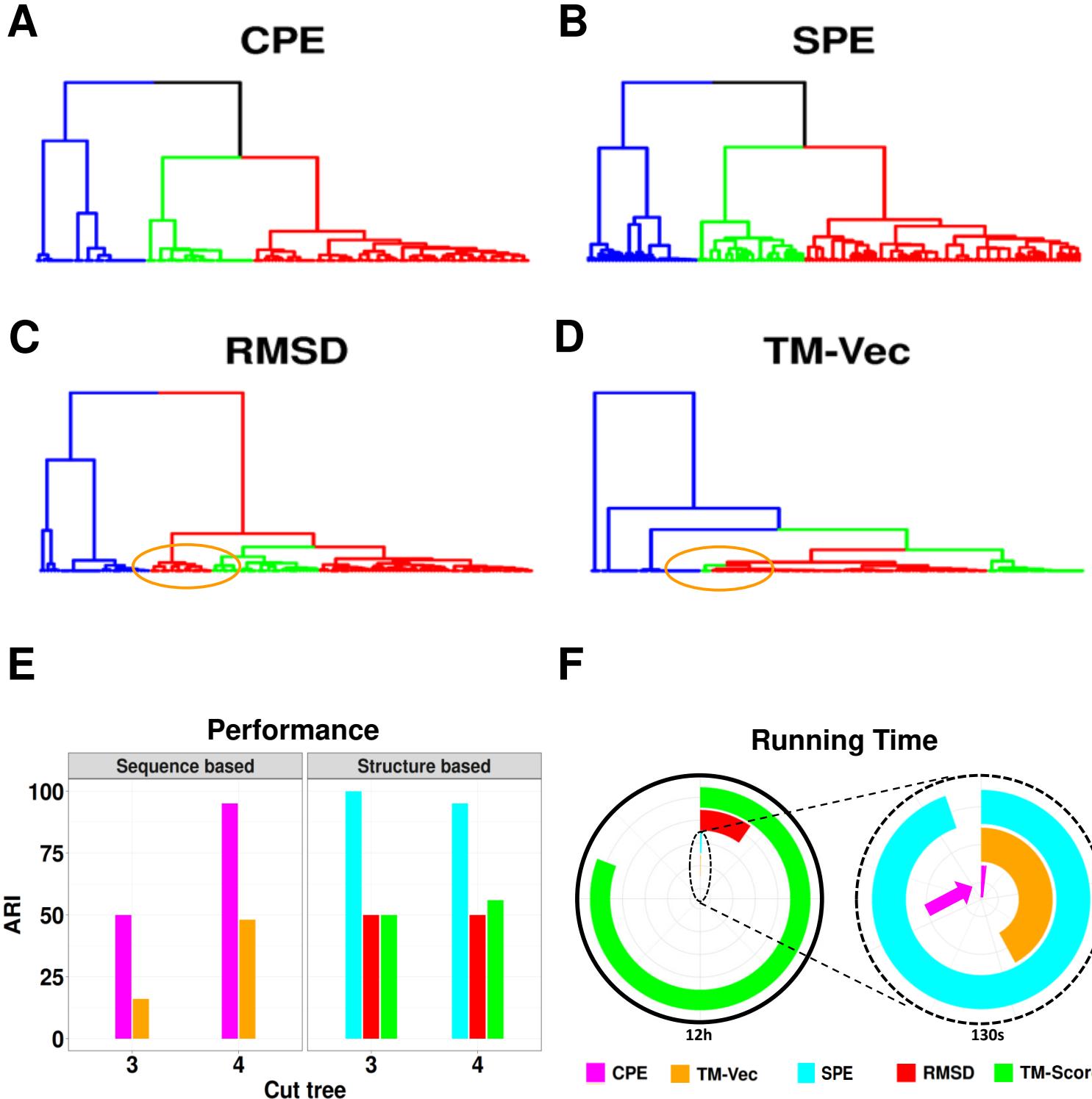
C

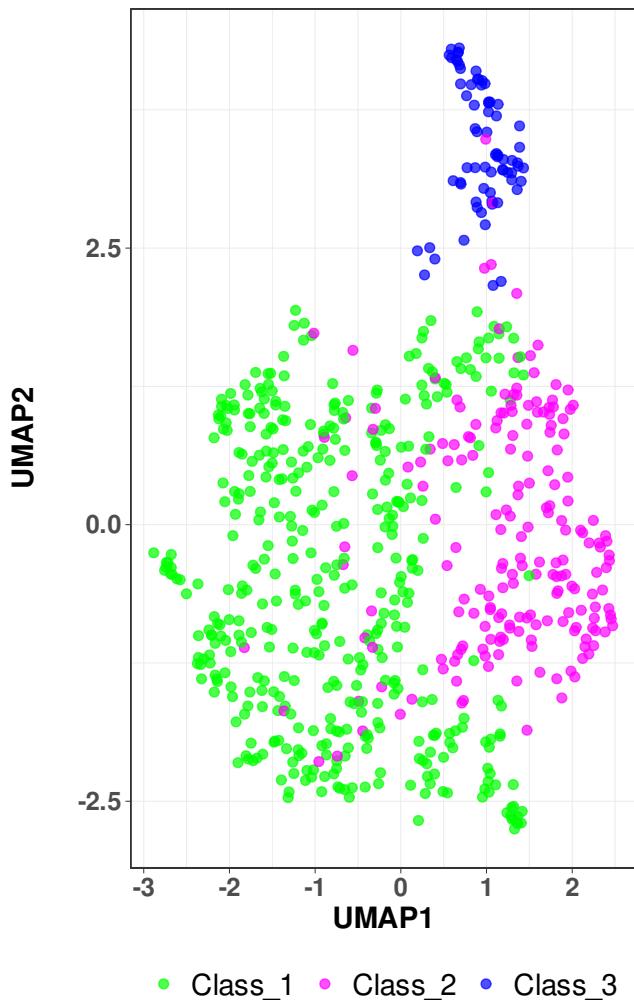
TM-Vec



Evolution of Spike protein

MERS_Cov → SARS_Cov → SARS_Cov2



A**B**