1    A general substitution matrix for structural phylogenetics.
2
3    Sriram G Garg[1] and Georg KA Hochberg[1,2,3]
4
5    1 Max Planck Institute for Terrestrial Microbiology, Karl-von-Frisch-Straße 10, 35043
6    Marburg, Germany.
7    2 Center for Synthetic Microbiology (SYNMIKRO), Philipps-University Marburg; Karl-von-
8    Frisch-Str. 14, 35043 Marburg, Germany
9    3 Department of Chemistry, Philipps-University Marburg; Hans-Meerwein-Str. 4, 35043
10   Marburg, Germany
11
12   Correspondence to: georg.hochberg@mpi-marburg.mpg.de
13

## 14   Abstract

15

16   Sequence-based maximum likelihood (ML) phylogenetics is a widely used method for
17   inferring evolutionary relationships, which has illuminated the evolutionary histories of
18   proteins and the organisms that harbour them. But modern implementations with
19   sophisticated models of sequence evolution struggle to resolve deep evolutionary
20   relationships, which can be obscured by excessive sequence divergence and substitution
21   saturation. Structural phylogenetics has emerged as a promising alternative, because
22   protein structure evolves much more slowly than protein sequences. Recent developments
23   protein structure prediction using AI have made it possible to predict protein structures for
24   entire protein families, and then to translate these structures into a sequence
25   representation - the 3Di structural alphabet - that can in theory be directly fed into existing
26   sequence based phylogenetic software. To unlock the full potential of this idea, however,
27   requires the inference of a general substitution matrix for structural phylogenetics, which
28   has so far been missing. Here we infer this matrix from large datasets of protein structures
29   and show that it results in a better fit to empirical datasets that previous approaches. We
30   then use this matrix to re-visit the question of the root of the tree of life. Using structural
31   phylogenies of universal paralogs, we provide the first unambiguous evidence for a root
32   between and archaea and bacteria. Finally, we discuss some practical and conceptual
33   limitations of structural phylogenetics. Our 3Di substitution matrix provides a starting point
34   for revisiting many deep phylogenetic problems that have so far been extremely difficult to
35   solve.
36

37   **Keywords:** Phylogenetics, Maximum likelihood, Structural phylogenetics, evolution,
38   substitution models
39

## 40   Introduction

41

42   The field of phylogenetics has evolved from relying on morphological comparisons to
43   sophisticated sequence-based analyses (Whelan et al., 2001). The advent of
44   computational methods marked a turning point, introducing a range of algorithms from
45   Neighbour-Joining (NJ) (Saitou and Nei, 1987) and Maximum Parsimony (MP) (Farris, 1970;
46   Fitch 1971) to Maximum Likelihood (ML) (Felsenstein, 1981) and Bayesian inferences
47   (Rannala and Yang, 1996; Mau and Newton, 1997) on nucleotide and amino acid

48  sequences. Each methodological leap has brought with it a deeper understanding of
49  evolutionary history through better trees. Among the various phylogenetic methods,
50  Maximum Likelihood (ML) approaches have emerged as particularly powerful tools for
51  modelling evolutionary processes (Posada and Crandall, 2021). The flexibility and
52  robustness of ML techniques have made them indispensable for contemporary
53  phylogenetic studies, especially those tackling large datasets or seeking to resolve deep
54  evolutionary relationships. But especially deep, sequenced-based phylogenetics remains
55  difficult. Substitution saturation is a particular challenge, in which each site in the
56  alignment has accumulated multiple substitutions over a branch of interest (Brown, 1982,
57  Phillippe and Forterre 1999). Depending on the accuracy of the substitution model of
58  sequence evolution, saturation can lead to spurious phylogenetic signals and artefacts in
59  phylogenetic trees (Felsenstein, 2003). The problem of saturation cannot always be solved
60  by adding more sequences (Philippe et al., 2011) or better models of sequence evolution.

62  Saturation is a relevant problem for the identification of the root of the tree of life. It is
63  traditionally placed on the branch between bacteria and archaea (Gouy et al., 2015), which
64  has important implications for the nature of the Last Universal Common Ancestor (LUCA).
65  This inference is based on paralog rooting with universally duplicated genes, where the
66  paralogs reciprocally root each other (Iwabe et al., 1989). Although this root is tacitly
67  accepted by the majority of biologists, the paralog trees it is based on are riddled with
68  potential problems. In all previous attempts, the branch between universal paralogs
69  remains so long as to be probably saturated (Brown and Doolittle, 1995; Philippe and
70  Forterre, 1999; Gouy et al., 2015; Mahendrarajah et al., 2023). This means that the root
71  position within each paralog might be mostly determined by the preferences of the
72  substitution model, rather than real phylogenetic signal, which has been erased almost
73  entirely. Some phylogenetics therefore still consider the root of the tree of life an unsolved
74  problem (Gouy et al., 2015).

76  Structural phylogenetics offers a potentially powerful alternative to traditional sequence-
77  based approaches. Structures evolve much more slowly than sequences, and if a model
78  for structural evolution could be inferred, this could help resolve phylogenies that are
79  beyond the reach of sequenced-based methods. Early attempts at this idea were limited by
80  the lack of high-quality protein structures or reliable methods of scoring multiple sequence
81  alignments of protein structures (Johnson, Šali, et al., 1990; Johnson et al., 1990; Balaji et
82  al., 2001; Balaji and Srinivasan, 2001). This changed with the advent of artificial intelligence
83  models than can predict protein structures with good accuracy (Jumper et al., 2021; Varadi
84  et al., 2023). The availability of a large database of structures has prompted researchers to
85  mould this novel source of information for identification of structural homologs in a process
86  similar to BLAST. Chief among these tools is FoldSeek which translates the 3D information
87  in predicted and experimentally determined structures into 20 unique characters the
88  authors call the 3Di alphabet (Kempen et al., 2023). The advantage of using an alphabet of
89  20 characters is that it enables the direct use of these 3Di characters in conventional
90  implementations of amino acid-based likelihood methods.

92  The conversion of a large dataset of 3D structures into the 3Di alphabet allows the
93  computation of a scoring matrix like the BLOSUM scoring matrix commonly employed by
94  Multiple Sequence Alignment programs (Kempen et al., 2023). This scoring matrix enables

95  the quick identification of structural homologs of proteins which has been very successful
96  in the identification of divergent orthologs. Such a scoring system also allows to compute a
97  similarity score (*fident* in case of FoldTree) which can then be used to compute Neighbour
98  Joining (NJ) trees as demonstrated by FoldTree (Moi et al., 2023). Furthermore, one could
99  also calculate a substitution matrix from this BLOSUM style scoring matrix which can be
100 directly implemented in ML approaches such as in the case of 3DiPhy (Puente-Lelievre et
101 al., 2024).  Neither of these approaches correspond to standard maximum likelihood
102 phylogenetics for amino acids: FoldTree's neighbour joining method is fast and simple but
103 inherits all limitations of classical neighbour joining in that it relies on the true distance
104 between sequences being close to their observed distance (an assumption that is often
105 violated in realistic datasets) and it does not account for among site rate variation
106 (Mihaescu et al., 2009). 3DiPhy does use a full likelihood model, which can account for
107 these phenomena however, its substitution matrix is derived from a BLOSUM-like
108 alignment scoring matrix. Such matrices are constructed by counting co-occurrences of
109 particular characters in sequence pairs, rather than inferring their contents using maximum
110 likelihood (Le and Gascuel, 2008). In standard sequence phylogenetics the BLOSUM matrix
111 has long been superseded by empirical models which are inferred in a full phylogenetic
112 likelihood framework, and generally result in a much better fit to empirical data (Le and
113 Gascuel, 2008).
114
115 These features of existing structural phylogenetics frameworks motivated us to infer a new
116 substitution model using a phylogenetic maximum likelihood framework. This substitution
117 model can in theory be directly inferred from each alignment in the form of a General Time
118 Reversible (GTR) model but inferring a substitution matrix for a 20-letter alphabet from a
119 single multiple sequence alignment is difficult and prone to overfitting. For conventional
120 protein models, this problem is solved by combining large numbers of protein alignments
121 and inferring from them one substitution model that best describes all the data. Once
122 computed, this general model, also denoted as $Q$, can then be used for individual protein
123 families, which avoids overfitting using GTR.  Here we make use of AlphaFold and a recently
124 developed protein large language model to infer a general substitution matrix for structural
125 phylogenetics. We show that this Q-matrix outperforms all previous methods to use 3Di
126 characters to infer ML phylogenies. Finally, we use our Q-matrix to re-infer the phylogenies
127 of universal paralogs and photosystems to settle long-standing questions in deep evolution
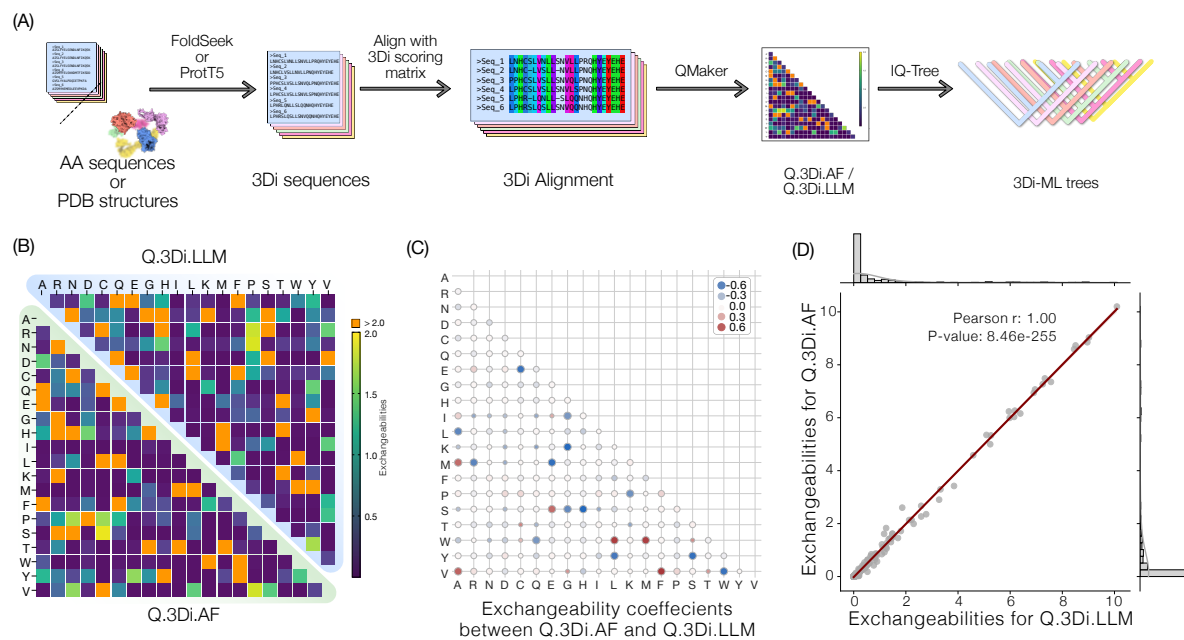128 that previously suffered from saturation.
129
130
131

*Figure 1: (A) Overview of the pipeline employed in the manuscript. Briefly, Amino acid (AA) or PDB structures were translated into 3Di characters using FoldSeek or the bilingual ProtT5 model. These 3Di characters are aligned with MAFFT using the 3Di scoring matrix before being used to estimate the general substitution Q-matrix using QMaker which were subsequently used to estimate 3Di-ML trees using IQ-tree. (B) Lower triangular portion is a representation of the Q-matrix estimated from 1660 AF clusters while the upper triangular section denotes the Q matrix estimated from 6653 PFAM clusters translated to 3Di alphabet using the ProtT5 bilingual language model. In both cases values higher than 2 are coloured orange. (C) Ratio of exchangeabilities between the Q.3Di.AF and the Q.3Di.LLM matrix. Each square represents the value $(m_1^{ij} - m_2^{ij})/(m_1^{ij} + m_2^{ij})$ where m1 and m2 represent Q.3Di.AF and Q.3Di.LLM respectively. (D) Pearsons correlation between the exchangeabilities of the two matrices indicating very little differences between the two matrices*

## Results and Discussion

*Estimation of the 3Di Q-matrix*

We set out to compute a general Q-matrix for structural phylogenetics. Given a large enough dataset, this is straightforward to achieve using the QMaker routine of IQ tree (Minh et al., 2021). We used two strategies to gather a large dataset of protein families and their predicted structures. Our first goal was to use the set of 6653 protein families that was used to infer a Q matrix in the initial study by Minh *et al*. To avoid having to predict AlphaFold models for every sequence in this large database, we opted to use a recently developed bilingual large language model, ProtT5. This model was trained to directly translate between an amino acid sequence and its corresponding 3Di sequence, without having to infer an AlphaFold model (Heinzinger et al., 2023). We used this method to translate all sequences in the PFAM dataset from AA-sequences to 3Di. The ProtT5 model is not perfect, as it introduces some randomness into the 3Di translation, meaning that translating to a 3Di sequence from the same input amino acid sequence results in a slightly different prediction (Supplementary Figure 1A-C). In addition, when comparing 3Di sequences extracted from AlphFold2 structures to the same 3Di sequence predicted with the LLM model, we found large numbers of sequences in which the AlphaFold and LLM predictions had low pairwise identities (Supplementary Figure 1D-F, Supplementary Figure 2). In order to safeguard against potential errors in estimating the substitution model using incorrectly translated 3Di

165 sequences we also estimated a separate Q matrix using 3Di sequences extracted from
166 AlphaFold predictions. We employed FoldSeek to cluster the SwissProt AlphaFold
167 Database. These 1660 AF clusters (hereafter AF-db) were used along with 3Di translation
168 of the 6653 protein families (hereafter Pfam-db), for the QMaker pipeline. Crucially, both
169 sets of 3Di sequences were then aligned using the alignment program *mafft* using the 3Di
170 scoring matrix from (Kempen et al., 2023) instead of the standard BLOSUM62 matrix used
171 for amino acid alignments.
172
173 We then estimated a tree for each of 3Di Multiple Sequence Alignments (MSAs) in our two
174 datasets using the GTR20 model despite the concern of model overfitting given the unique
175 nature of the 3Di alphabet and the lack of other models that could serve as the initial
176 starting model. These initial trees were then used to estimate a single Q-matrix that best
177 explains the respective sets of MSAs as described in the QMaker pipeline (Minh et al., 2021).
178 This resulted in two Q-matrixes hereafter denoted as Q.3Di.AF and Q.3Di.LLM. The two Q-
179 matrices estimated were very similar with minor differences in exchangeabilities (Figure 1C)
180 with a Pearsons correlation of 1 (Figure 1D). We then checked if these matrices are
181 preferred by IQ-Tree's *modelfinder* over the GTR20 or the previously published 3DiPhy
182 model using a test set of 6653 3Di MSAs from PFAM that were not used for estimating the
183 Q-matrix. Indeed, the 3DiPhy model is only preferred in 278 MSAs over 6267 MSAs that
184 prefer either the Q.3Di.AF or the Q.3Di.LLM model, which are practically the same (Figure
185 1B-D). This increased our confidence that we had successfully captured the mechanism of
186 change describing the mutability in the structural alphabet across a wide range of proteins.
187 In the analyses of specific protein families that follow, IQ-Tree's *modelfinder* predominantly
188 chose Q.3Di.AF over Q.3Di.LLM or GTR20 according to the Corrected Akaike Criterion
189 (AICc). Generally, we encourage future users of these matrices to always test if using
190 Q.3Di.AF changes any conclusions in cases where Q.3Di.LLM is the better fit model. This is
191 because the AF matrix is much less affected by the misprediction issues than the LLM
192 (which we discuss further below).
193

| Model | AICc | AIC | BIC |
|---|---|---|---|
| Q.3Di.AF | 2342 | 2065 | 2309 |
| Q.3Di.LLM | 3925 | 2958 | 3697 |
| 3DiPhy | 278 | 322 | 267 |
| GTR20 | 108 | 1308 | 380 |
| *Total* | *6653* | *6653* | *6653* |

194
195 Table 1: Number of trees that preferred each model/Q-matrix as identified using *modelfinder* from IQ-Tree
196 according to corrected Akaike Information Criterion (AICc), Akaike Information Criterion (AIC) and Bayesian
197 Information Criterion (BIC)
198
199 *Rooting the ToL using structural phylogenetics*
200
201 Rooting the tree of life is a particularly challenging problem owing to the lack of outgroups
202 that can reliably root phylogenetic trees. Paralog rooting is a powerful method which uses
203 phylogenetic trees with duplicated genes that reciprocally root each other. In most cases
204 the paralogs root each other along the same branch recovering an unambiguous root for
205 the species tree containing the paralogs. However, in cases of highly divergent paralogs,

the two paralogs sometimes do not agree on the same root (Figure 2A). We tested if our new matrix can help improve trees used to root the tree of life using two universal paralogs that have been previously used for this purpose: Elongation factors and catalytic and non-catalytic subunits of the rotary ATPase. We begin with the Elongation factor phylogeny. Elongation factor EF-Tu/EF-2 delivers aminoacyl-tRNAs to the A-site of the ribosome while the Elongation Factor EF-G/EF-1A catalyses the translocation of the peptidyl-tRNA (Miller, 1972). Both paralogs are conserved across the tree of life, making them an ideal candidate for paralog rooting (Baldauf et al., 1996; Philippe and Forterre, 1999; Gouy et al., 2015). In all previous attempts to root the tree of life using EF-G and EF-Tu, the branch separating the paralogs is extremely long and potentially completely saturated, which implies that the position of the root within each paralog might be determined entirely by the substitution model and not by any synapomorphies between the paralogs. In addition, the two paralogs do not root each other consistently increasing the uncertainty.

To test if our new matrix can help solve this problem, we first assembled a dataset of 1076 homologs of EF-Tu and EF-G. In an amino acid-based ML tree we also recover a very long branch (Branch length (BL) = 3.284) between the two paralogs albeit still separating the bacteria and archaea (Figure 2A). In line with previous phylogenies, this tree recovers different roots for the tree of life in the two paralogs: between bacteria and archaea plus eukaryotes, and between archaea and bacteria plus eukaryotes (Figure 2B). We then extracted 3Di sequences from 1076 AlphaFold predictions using FoldSeek (see methods) and utilized our new Q.3Di.AF Q-matrix as the substitution model, to estimate a new tree of the EF-G and EF-tu paralogs. This recovered a phylogenetic tree with the length of the branch separating the paralogs far below 1 (0.186). Crucially, the root position is now consistent in both the paralogs and indicates a root between archaea and bacteria for life (Figure 2C). The archaea in both paralogs remain paraphyletic, which is consistent with the two-domain tree of life.
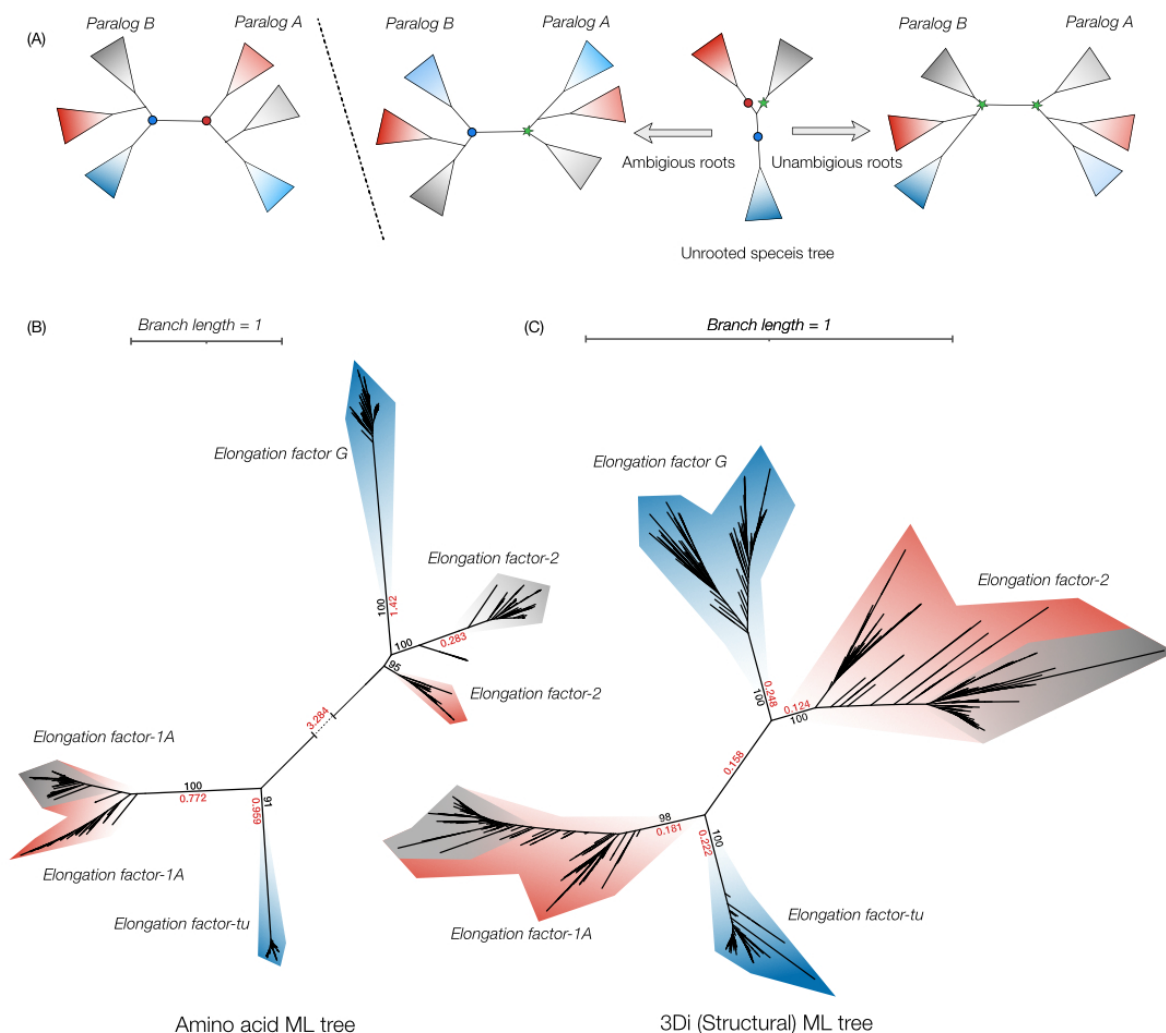
Figure 2: *(A) A schematic representation of paralog rooting. Three possible root positions are shown with the "true" root depicted with a green star and two other possible roots with circles. In the scenario where paralogous rooting is successful both paralog subtrees reciprocally root each other (right). Other possible scenarios are also shown where the paralog subtrees are ambiguously rooted (left).* (B) Amino acid ML tree containing 1076 EF-tu and EF-G homologs from eukaryotes, bacteria and archaea. The mitochondrial and plastid encoded copies are not included. Note that the branch separating the EF-tu and EF-G is broken for illustration. *(C)* 3Di structural ML tree estimated 3Di sequences and the Q.3Di.AF model from the predicted AlphaFold structures of 1069 EF-tu and EF-G homologs. In both cases blue, red and grey clades represent bacteria, archaea and eukaryotes respectively. Numbers in red, black indicate branch lengths and ultrafast bootstrap supports respectively.

Another universally conserved paralogous gene family used to root the tree of life are the catalytic and non-catalytic subunits of the rotary ATPase. The head group of the rotary ATPase is a hexamer consisting of two subunits, only one of which is catalytic (Figure 3A). The bacterial and mitochondrial ATPases are called the $F_0F_1$-ATPases, and their subunits are called F1-alpha and F1-beta for the non-catalytic and catalytic subunits respectively (Grüber et al., 2001). The archaeal ATPase is called the V-ATPase and shares a similar architecture with a non-catalytic and a catalytic subunit in its headgroup (Figure3A). Owing to the endosymbiotic event between archaea and bacteria at eukaryogenesis, the eukaryotes and archaea also share this ATPase which in eukaryotes is in the vacuole, where it functions to acidify lysosomes (Gogarten et al., 1989). The archaeal/eukaryotic subunits

255    are named V1-beta and V1-alpha for the non-catalytic and catalytic subunits respectively
256    (Grüber et al., 2001; Cross and Müller 2004). A recent analysis on rooting the ToL using the
257    ATPase subunits (Mahendrarajah et al., 2023) recovers a tree that separates the four major
258    subunits with extremely long basal branches (Figure 3B).  This tree is consistent with the
259    idea that the catalytic and non-catalytic subunits originated before the divergence of
260    archaea and bacteria, and roots the tree of life between these two domains. The same study
261    also identified an early transfer of the archaeal non-catalytic subunit into bacteria,
262    however, the catalytic counterpart to this transfer was not recovered in the catalytic sub-
263    tree suggesting multiple transfer events (Figure 3B).
264
265    As before, we predicted AlphaFold structures for all 1520 sequences and extracted the 3Di
266    sequences using FoldSeek and calculated a 3Di (structural) ML tree with the Q.3Di.AF as
267    the substitution model. While the tree in this case looks remarkably like the amino acid ML
268    tree, the 3Di structural ML tree has significantly shorter branches (Figure 3C). This new
269    topology also reconfirms the root of ToL as between the archaea and bacteria.
270    Furthermore, in the 3Di tree the early transfer of the archaeal ATPase subunits is recovered
271    basal in both catalytic and the non-catalytic subtrees suggesting a single early transfer from
272    archaea to bacteria. Together with the Elongation factors, our results bolster support for the
273    two-domain tree of life with the eukaryotes branching within archaea. In both these cases
274    it is evident that structural phylogenetics can resolve deep phylogenies and recover
275    consistent groupings within the paralogs despite large divergences in amino acid
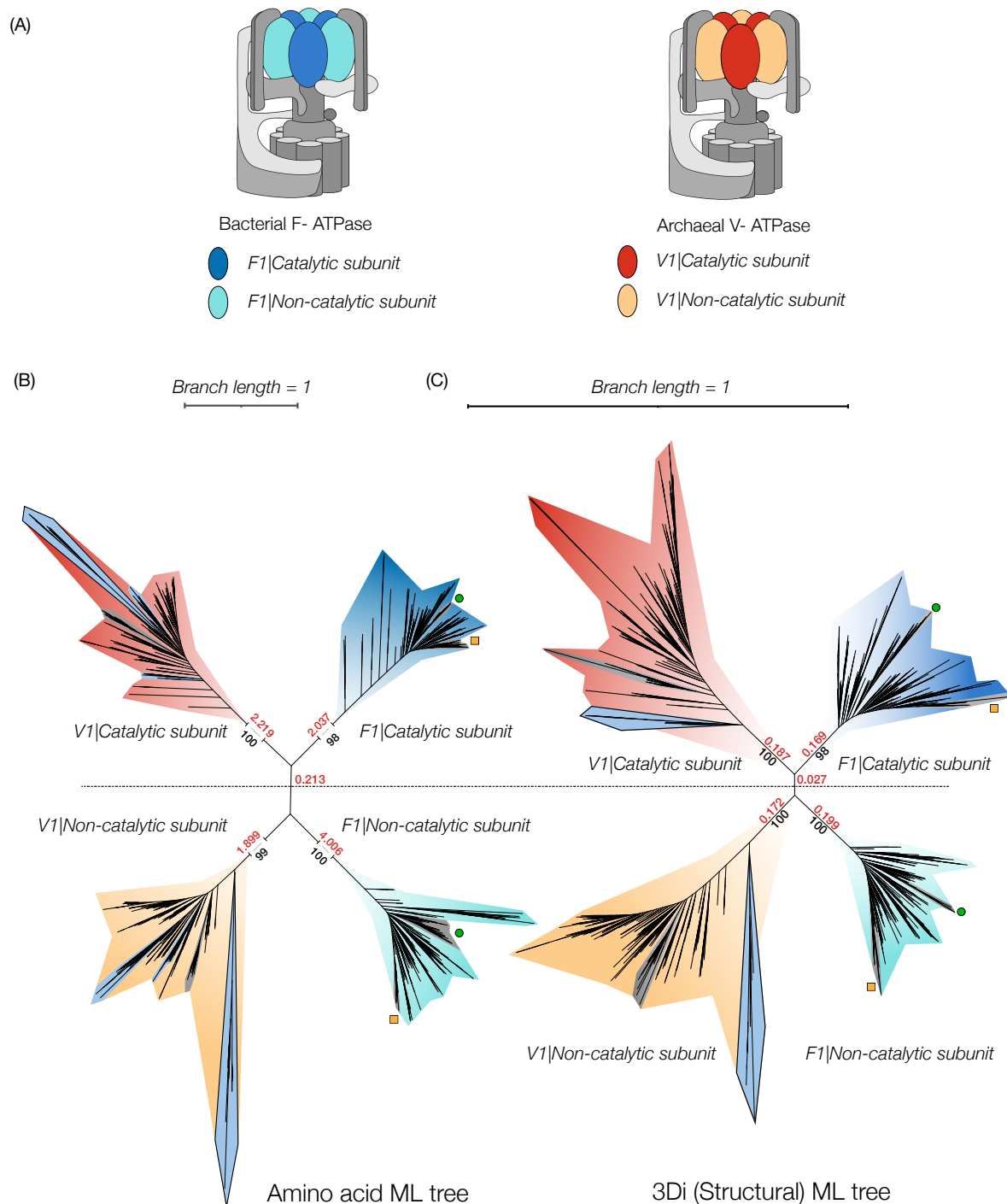276    sequences.

*Figure 3: (A) A schematic representation of the bacterial and archaeal ATPase highlighting the subunits under investigation. They are represented using the same colours in the phylogenetic trees. (B) Amino acid ML tree of 1520 sequences across the ToL reproduced from Mahendrarajah et al., 2023 of the catalytic and non-catalytic subunits of bacterial, archaeal, and eukaryotic rotary ATPase. The early branching transfer from bacteria and archaea in the non-catalytic V1 clade is highlighted in blue with a black outline. The corresponding clade in the V1 catalytic clade branches deep inside of the archaeal sequences and is highlighted similarly. (C) 3Di structural tree estimated using the Q.3Di.AF model. Sequences assigned to the early transfer from the archaeal clade to bacteria are highlight as in (B), but now this transfer is inferred for both the catalytic and non-catalytic subunits. Numbers in red, black indicate branch lengths and ultrafast bootstrap supports respectively. In both cases grey clades represent eukaryotes. The green circles and orange squares indicate cyanobacterial and proteobacterial contributions in eukaryotes representing the plastid and mitochondrial ATPases.*

291

292    *Evolution of photosystems RCI and RCII*

293

294    The issue of saturation is not exclusive to tree-of-life problems but to all evolutionarily
295    divergent proteins that share remote homology in sequence. The origin of oxygenic
296    photosynthesis is another event that impacted the overall geochemistry of the planet and
297    has been the subject of contentious debate. Photosynthesis can be classified into two
298    major types: anoxygenic photosynthesis, which uses either reaction centre II (RCII) or
299    reaction centre I (RCI), but never both together and oxygenic photosynthesis which uses
300    both reaction centres I and II (RCI and RCII) coupled to a water splitting reaction that leads
301    to the formation of oxygen (Hohmann-Marriott and Blankenship, 2011). One set of theories
302    suggests that anoxygenic photosynthesis evolved first and later developed into oxygenic
303    photosynthesis (Martin et al., 2018). An alternative view favours oxygenic photosynthesis to
304    have evolved first, with anoxygenic phototrophs having lost either RCI or RCII.  One piece of
305    evidence for the latter view is the lack of any bacterial group that harbours the anoxic
306    versions of both RCI and RCII, which is thought to be a necessary precursor to oxygenic
307    photosynthesis (Sánchez-Baracaldo and Cardona, 2020). Until recently, members of the
308    *chloroflexota* phylum have only been known to harbour anoxic RCII. This changed when a
309    *chloroflexota* group, *Ca. Chloroheliales,* was identified that contains RC1 (Tsuji et al.,
310    2024). This still falls short of proving that anoxic RCI and RCII have existed together in the
311    same genome however, one possible interpretation of these data is that an ancestral
312    *Chloroflexus* might have contained both, leading to differential losses in extant lineages of
313    *Chloroflexi*. This would support the idea that anoxic photosynthesis may have come first, if
314    these photosystems are close relatives of the photosystems that were eventually
315    transferred into cyanobacteria

316

317    The phylogenetic tree based on amino acids of RCI containing *Chloroflexi* does not place
318    their RCI sequences as close relatives to those of cyanobacteria (4A, re-inferred for this
319    study). But this tree suffers from extremely long branches, and we wondered whether this
320    placement is the result of long branch attraction. We therefore set out to re-infer this tree
321    using 3Di characters and our structural substitution matrix (Figure 4B). This shortened all
322    relevant internal branches to lengths well below one but yielded the same topology as the
323    amino acid tree. This confirms the authors' original inferences and leaves the evolution of
324    oxygenic photosynthesis an unsolved problem for now.
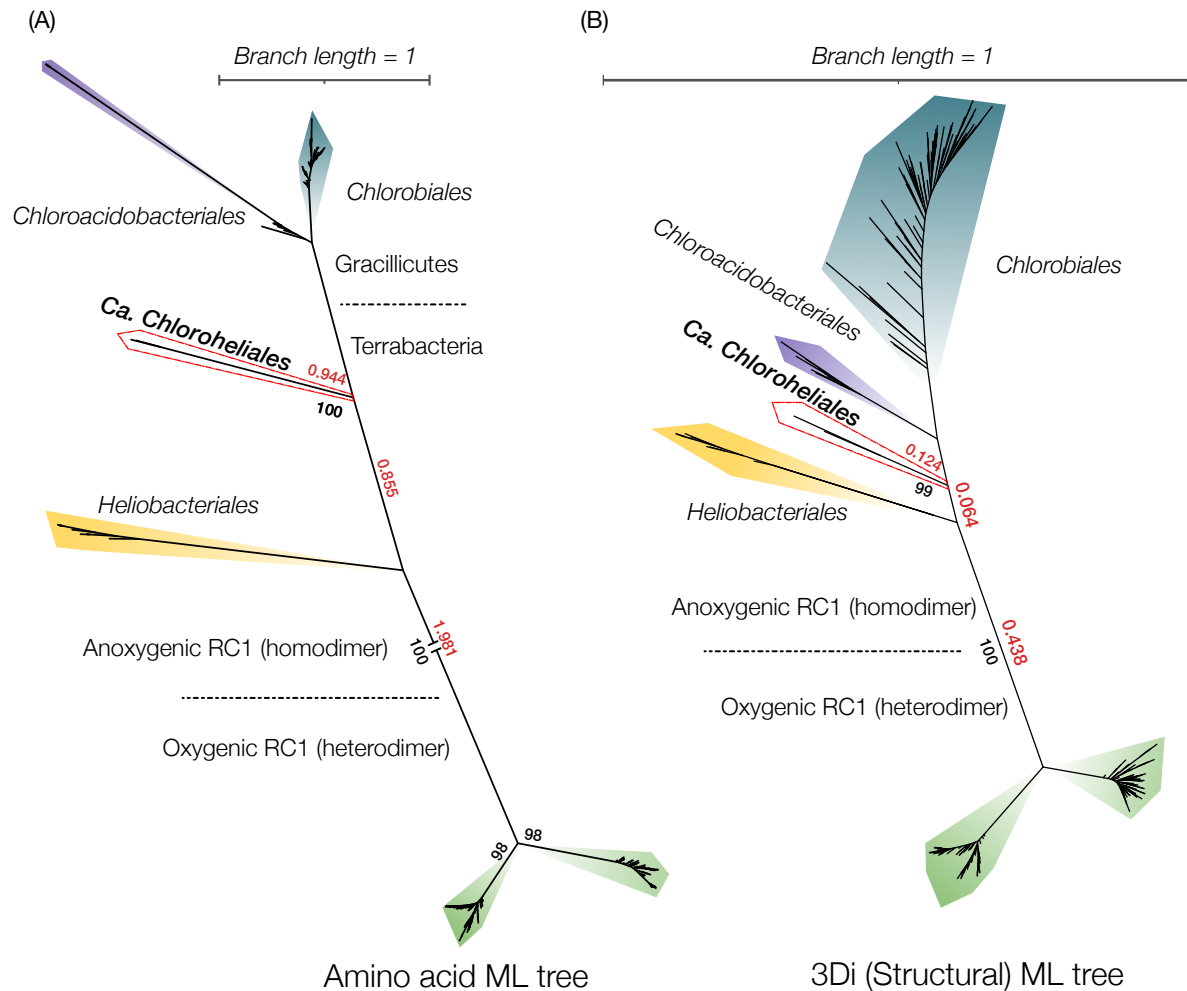
325

326

Figure 4: *(A)* Amino acid ML tree of 321 RC1 protein sequences. Note that long branches are broken as indicated for illustration. *(B)* 3Di structural ML tree of 297 3Di sequences from AlphaFold structures using the Q.3Di.AF model. Numbers in red, black indicate branch lengths and ultrafast bootstrap supports respectively.

Our work in this manuscript and that of others (Moi et al., 2023; Puente-Lelievre et al., 2024) clearly points to the utility of structural phylogenetics in cases where structures can be predicted reliably and with one possible structure per sequence. There are several practical and conceptual caveats that come with using this method, which we will briefly elaborate on. We present these caveats in the spirit of critical optimism about the utility and impact of this new method.

*Prediction accuracy of LLMs*

Structural phylogenies can only ever be as good as the predicted models that are used to derive 3Di sequences. Predicting large numbers of sequences with AlphaFold is computationally costly and potentially prohibitive for many interested users. Using bilingual Protein LLMs like ProtT5 may seem like an obvious solution, because it removes the computationally expensive requirement of predicting the AF structures of a large number of protein clusters not only in the Q-matrix estimation, but also for tree inference of single protein families with a lot of members. Encouragingly, the Q-matrix estimated from 3Di

349  sequences derived from AlphaFold structures (Q.3di.AF) is very similar to the one from
350  PFAM clusters translated using ProtT5 (Q.3Di.LLM) despite their low accuracy compared to
351  AlphaFold predictions (Figure 1C,1D, Supplementary Figure 1). This could be due to the fact
352  that the LLM has issues when dealing with long repeated stretches which in some cases
353  leads to possible register shifts of structural motifs (Supplementary Figure 2). These register
354  shifts can be dealt with during a 3Di sequence alignment using the 3Di scoring matrix, which
355  we did not perform for our pairwise identity calculations (as the two sequences are the
356  same length). It is also possible that prediction errors average out when inferring a Q-matrix
357  for thousands of protein families, even if there are substantial errors in the alignments of
358  any one family. It is clear, however, that ProtT5 translations are not reliable for inferring
359  individual trees. We tested this by using ProtT5 derived 3Di alignments for the three protein
360  families we investigated here. In two out of three cases we recovered phylogenies that
361  either were biologically improbable (Supplementary Figure 4) and/or erroneous with non-
362  sensical topologies (Supplementary Figure 5). Most of these issues stem from the faulty
363  prediction of 3Di sequences. While we did not observe this problem here when using
364  AlphaFold structures, we expect similar issues when using structures that are not
365  confidently predicted by AlphaFold. For now, reliable tree inference only seems possible
366  using AlphaFold generated structures and therefore comes with a significant
367  computational overhead. Better language and structure prediction models are certain to be
368  available in the future and they should make structural phylogenetics more widely
369  accessible.
370
371  *Fold-switching and conformational variability*
372
373  Many proteins undergo conformational changes and some even switch their folds entirely
374  as part of their functions (Chang et al., 2015). Previous analyses using AlphaFold suggests
375  that it can sometimes predict structures in different conformations despite having a strong
376  bias towards one dominant conformation (Chakravarty and Porter 2022; Sala et al., 2023;
377  Wayment-Steele et al., 2023). Since this can lead to different 3Di sequences for the same
378  protein, depending on which conformation it is predicted in, we wondered if this could lead
379  to spurious grouping according to conformation rather than genealogical relationships on
380  3Di phylogenies. We examined two proteins for this purpose. One is KaiB, which is known
381  to fold-switch as part of its catalytic cycle, involving a drastic change of a helix to a beta-
382  sheet (Chang et al., 2015; Zhang et al., 2023). The other is the RNA Polymerase III subunit
383  Rpc10, which undergoes a conformational change during its function in gene transcription
384  (Girbig et al., 2021).
385
386  To test how much this issue can affect 3Di trees, we constructed a worst-case scenario for
387  both proteins. In both cases, we modelled each sequence on the two distinct
388  conformations using homology modelling and inferred their 3Di sequences using FoldSeek.
389  For tree inference, we then randomly chose the 3Di sequence of one of the two possible
390  conformations for each protein, such that approximately half our sequences were
391  predicted in one conformation, and the other half in the other conformation. For both KaiB
392  and Rpc10 we found that the phylogenetic tree splits the two conformational states with
393  long branches (Figure 4B, 4D) as opposed to a 3Di structural tree which was inferred from
394  3Di sequences reflecting a single conformation (Figure 4A, 4C). This highlights a severe
395  limitation of structural phylogenetics where the presence of multiple predicted

396    conformations can generate spurious branches and relationships. Here we concocted an
397    extreme case by forcing sequences randomly into distinct conformations. However, in
398    cases where only a small minority of proteins within the analyses share a different
399    conformation, these artefacts can lead to false conclusions. It is therefore very important
400    to assess the conformational homogeneity of the predicted sequences before inferring a
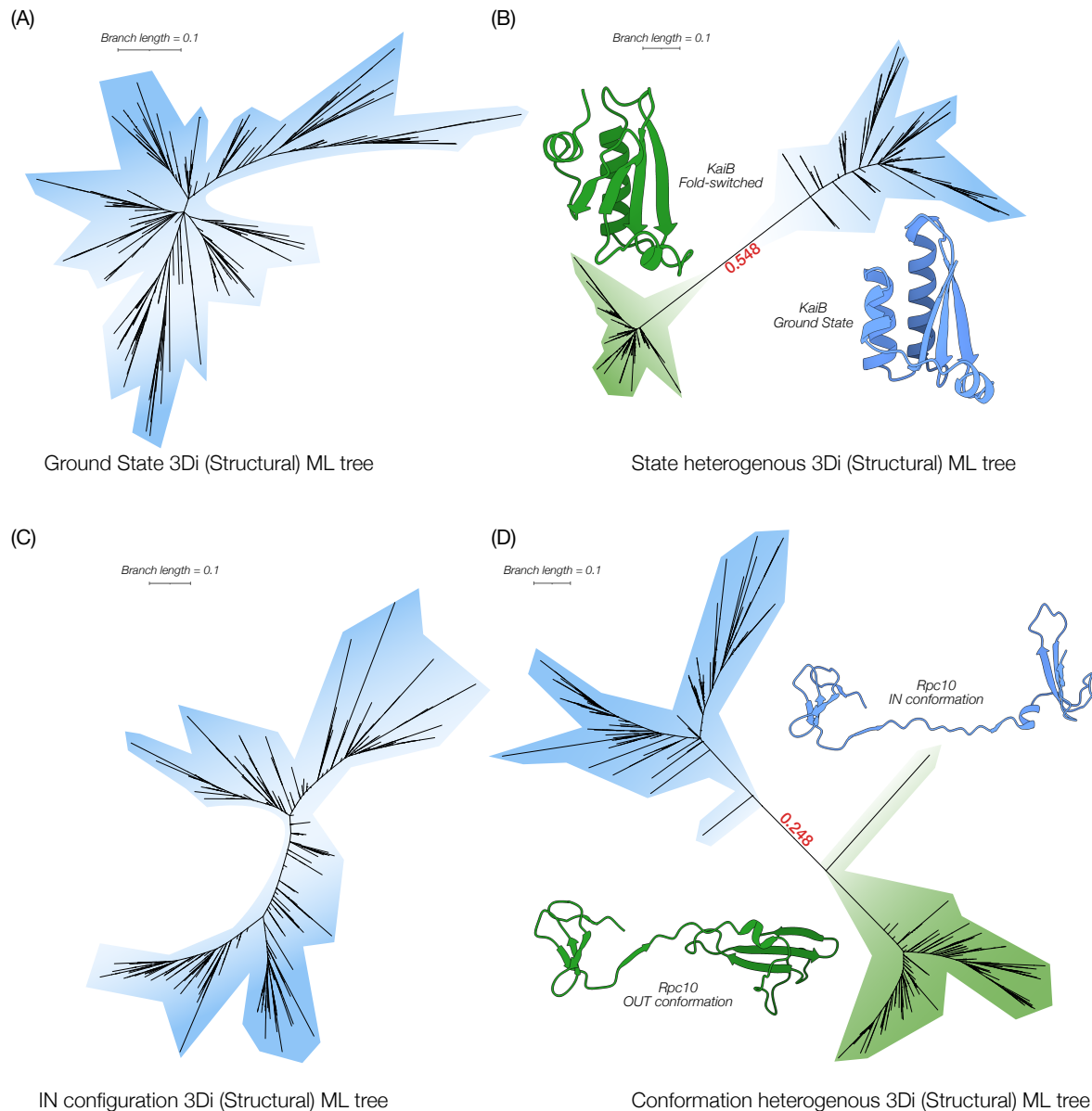401    3Di tree.



402
403    Figure 5: (A) 3Di structural ML tree constructed from KaiB proteins modelled in the ground state. (B) 3Di
404    structural ML tree constructed from approximately 50% of the KaiB proteins modelled in the ground state
405    (blue) and the other 50% modelled in the fold switched state (green). (C) 3Di structural ML tree constructed
406    from RPC10 proteins modelled in the IN conformation (D) 3Di structural ML tree constructed from
407    approximately 50% of the RPC10 proteins modelled in the OUT conformation (blue) and the other 50%
408    modelled in the IN conformation (green). In both cases the distinct conformations form monophyletic groups
409    in contrast to their placements in (A) and (C) respectively.
410
411    *Site-independence in structural alignments*
412

413   One of the main assumptions of maximum likelihood is site independence, which allows
414   the likelihood to be computed independently for all sites in the alignment (Liò and Goldman
415   1998). It has been obvious for a long time that this is not a realistic assumption. Epistasis
416   between amino acids is a well demonstrated phenomenon and quite extensive among
417   proteins (Starr and Thornton 2016). This violates the site-independence assumption of
418   maximum likelihood phylogenetics, even though it has been argued that increasing the
419   number of sites normally associated with a protein sequence or increasing the number of
420   proteins used for a concatenated alignment averages out the signal in most cases (Starr
421   and Thornton, 2016; Magee et al., 2021). In the case of the structural phylogenetics and 3Di
422   alphabet however, this assumption is explicitly violated since each letter corresponds to at
423   least 6 other amino acid positions in 3D space. It is for example not clear to us that it is even
424   possible for a single substitution to occur at the level of 3Di characters, because of the
425   structural dependence between sites. In a sense, structural phylogenetics makes the ugly
426   truth of model violation explicit in its alphabet. Whether or not this approach becomes
427   widely accepted in evolutionary biology will depend on investigating the consequences of
428   this violation, which is beyond the scope of this manuscript.
429
430   *Information loss*
431
432   The 3Di alphabet compresses information that would be present in amino acids. This is the
433   very reason for its utility in deep phylogenetics, because it overcomes the saturation
434   problem. But it also makes evolution on short time-scales is harder to resolve using these
435   models, and relationships at the very tips of 3Di trees probably much less reliable than in
436   an amino acid or DNA tree (Mutti et al.,2024). A potential solution is to use partitioned
437   models, in which a tree is inferred from both 3Di and amino acid alignments
438   simultaneously, using different substitution models for the partitions (Puente-Lelievre et
439   al., 2024). To make this approach work, however, one would have to allow the structural
440   partition to also have a different set of branch lengths than the amino acid partition (Lopez
441   et al., 2002), which the first use of this approach did not yet include (Puente-Lelievre et al.,
442   2024). Such a heterotachous model presents a difficult optimization problem, which in our
443   hands leads to impractically long run times on our datasets. Another question is the size of
444   the alphabet. 3Di uses 20 characters because this allows simple integration with existing
445   phylogenetic software. It is not yet clear that whether this is even close to the optimal
446   number of characters for structural phylogenetics. Larger alphabets could perhaps retain
447   more short-term information. They would, however, make the inference of substitution
448   matrices much harder.
449
450   *Structural phylogenetics and the future of deep history*
451
452   Our work complements and builds on other recent tools that utilise the 3Di alphabet for
453   structural phylogenetics (Moi et al., 2023; Puente-Lelievre et al., 2024).  Our structural Q-
454   matrices should make it much easier to infer structure-based trees for anyone familiar with
455   maximum likelihood phylogenetics. Newly developed online tools for the generation of 3Di
456   alignments should further lower the bar for adoption (Gilchrist et al., 2024). As with every
457   new method, it is difficult to know exactly what impact structural phylogenetics will have.
458   For now, we see its main utility in solving difficult rooting problems involving distant
459   outgroups that amino acid phylogenies cannot solve with any degree of confidence.  This

460 will help polarize the direction of evolutionary change in the emergence of many important
461 functions. Better resolved deep protein phylogenies will also improve our reconstructions
462 of the gene content of ancient organisms (The Last Universal Common Ancestor, the Last
463 Eukaryotic Common Ancestor, and the Last Archaeal Common Ancestor, for example).
464
465 For now, these methods will not be useful for ancestral sequence reconstruction, because
466 3Di sequences cannot be back translated into a unique amino acid sequence (Heinzinger
467 et al., 2023). Even though our matrix allows us to infer 3Di sequences at internal nodes of
468 structural phylogenies, it is at present not possible to then turn these reconstructed 3Di
469 sequences into resurrected proteins composed of amino acids. It may, however, be
470 possible to restrict a set of plausible amino acid reconstructions at one particular node on
471 an amino acid phylogeny to a subset that agrees with the reconstructed ancestral 3Di
472 sequence at the corresponding node on a structural phylogeny.
473
474 The true impact of viewing the past through the glacial change in the structure of proteins
475 will only emerge when this method is robustly tested and becomes widely adopted in
476 evolutionary biology. We hope the matrix inferred here will be a first step in this process.
477
478 ## Methods
479
480 *Datasets for QMaker*
481 The SwissProt AlphaFold database (https://alphafold.ebi.ac.uk/download) was
482 downloaded and then clustered with FoldSeek (https://github.com/steineggerlab/foldseek)
483 *easy-cluster* with default settings and a coverage of 80%. This yielded 1660 clusters which
484 contained at least 50 members and a maximum of 2500 members. Databases of PDB
485 structures were then created and 3Di sequences were subsequently extracted from these
486 1660 clusters using FoldSeek as previously described. The PFAM sets were taken from
487 Minh et al., 2021 which contained 6655 protein families used for training the Q-matrix and
488 a further 6653 families were used for testing. In the case of PFAM families the amino acid
489 FASTA files were directly translated to the 3Di alphabet using the scripts provided by
490 Heinzinger et al.,, 2023 (https://github.com/mheinzinger/ProstT5).
491
492 *Q-matrix estimation*
493 Both the AF-db and PFAM-db sets of 3Di sequnces were aligned using *ginsi* method within
494 Mafft (v7.515) and the 3Di scoring matrix from FoldSeek using the –-*aamatrix* flag
495 implemented within mafft. The 3Di MSAs thus generated were then used in the QMaker
496 routine as described in Minh et al.,, 2021 (http://www.iqtree.org/doc/Estimating-amino-
497 acid-substitution-models). Briefly, for each MSA the best fit substitution model was
498 initialised with GTR20 along with the best RHAS model to account for rate-heterogeneity
499 using ModelFinder (Kalyaanamoorthy et al., 2017). In the Next step we estimate a joint
500 reversible Q-matrix for all the 3Di MSAs as described.
501
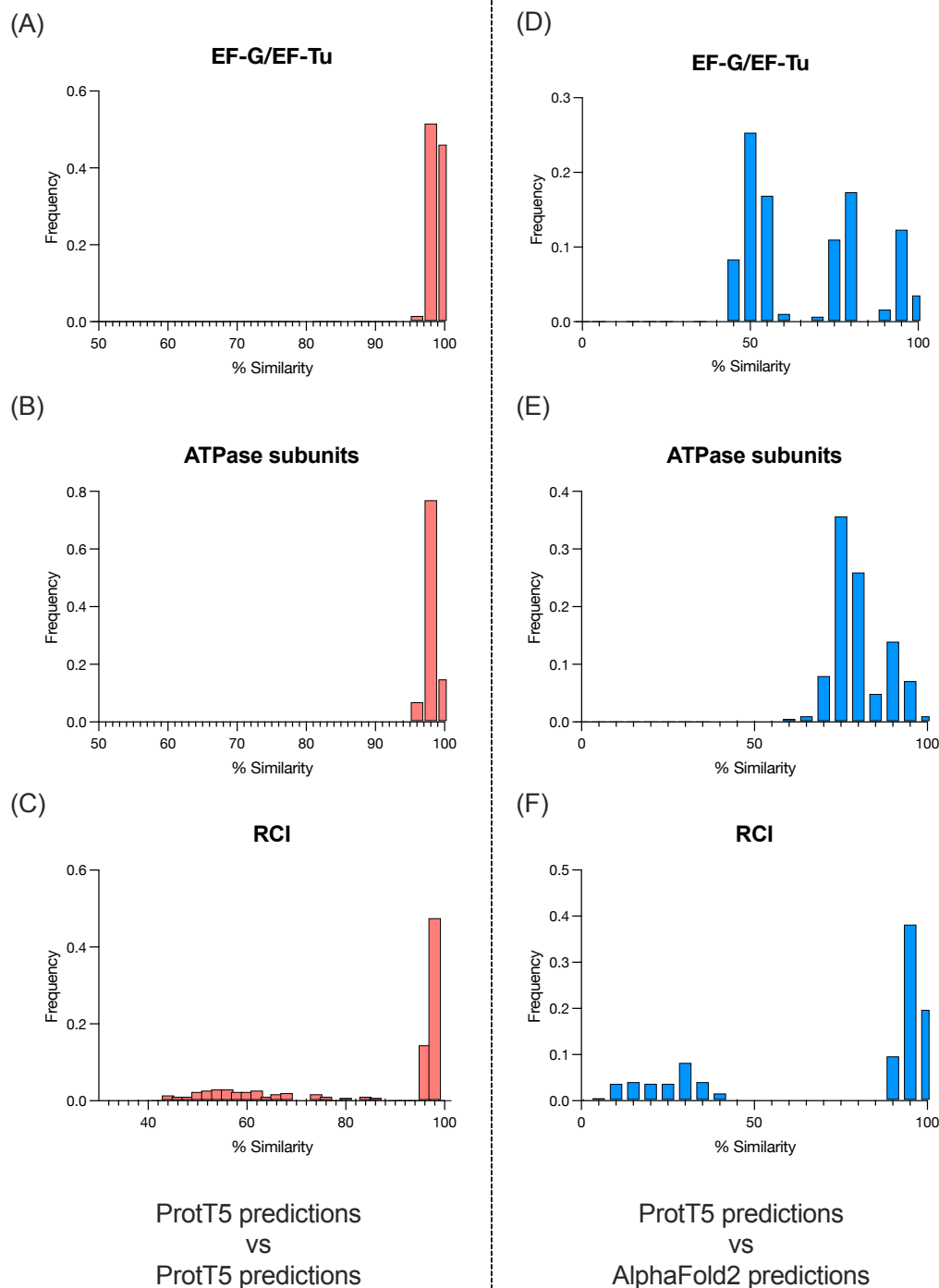502 *Individual Protein/3Di sequences and Phylogenetic tree reconstructions*
503 Elongation factors and Reaction Centre I homologs were identified using BLAST against the
504 NCBI non-redundant (*nr*) database, and then filtered with a minimum similarity threshold
505 of 50% and an e-value cutoff of 1E-5. For the ATPase phylogeny was exactly reproduced
506 from Mahendrarajah et al.,, 2023 and the same sequences used for the 3Di sequences. The

507 amino acid sequences were aligned using *linsi* and then subsequently trimmed using TrimAl
508 (v1.4) (Capella-Gutiérrez et al., 2009) with the *-automated1* setting. Trimmed amino acid
509 alignments were then used for maximum likelihood tree estimation using IQ-tree with the
510 best-fit model suggested by ModelFinder. 3Di sequences for individual proteins trees were
511 extracted from PDB files from individual AlphaFold (v2.2.0) predictions. The best ranked
512 AlphaFold models were used to create a database using FoldSeek which allowed us to
513 extract 3Di sequences from the PDB structures. For 3Di sequences translated from ProtT5,
514 the model was queried as described in Heinzinger et al., 2023 using amino acid sequences
515 as input. All 3Di sequences were aligned with Mafft (*ginsi*) using the --aamatrix option
516 specifying the 3Di scoring matrix provided by van Kempen et al., as part of FoldSeek. 3Di
517 MSAs were then used to estimate the structural ML tree as described above. For individual
518 3Di tree reconstructions IQ-tree (v2.3.0) was used to identify the best-fit model (Q.3Di.AF,
519 Q.3Di.LLM or GTR20) according to AICc, along with rate-heterogeneity using ModelFinder.
520 Both amino acid and 3Di trees were estimated with 10000 Ultrafast bootstraps (-bb) and
521 10000 iterations for SH-test (-alrt).
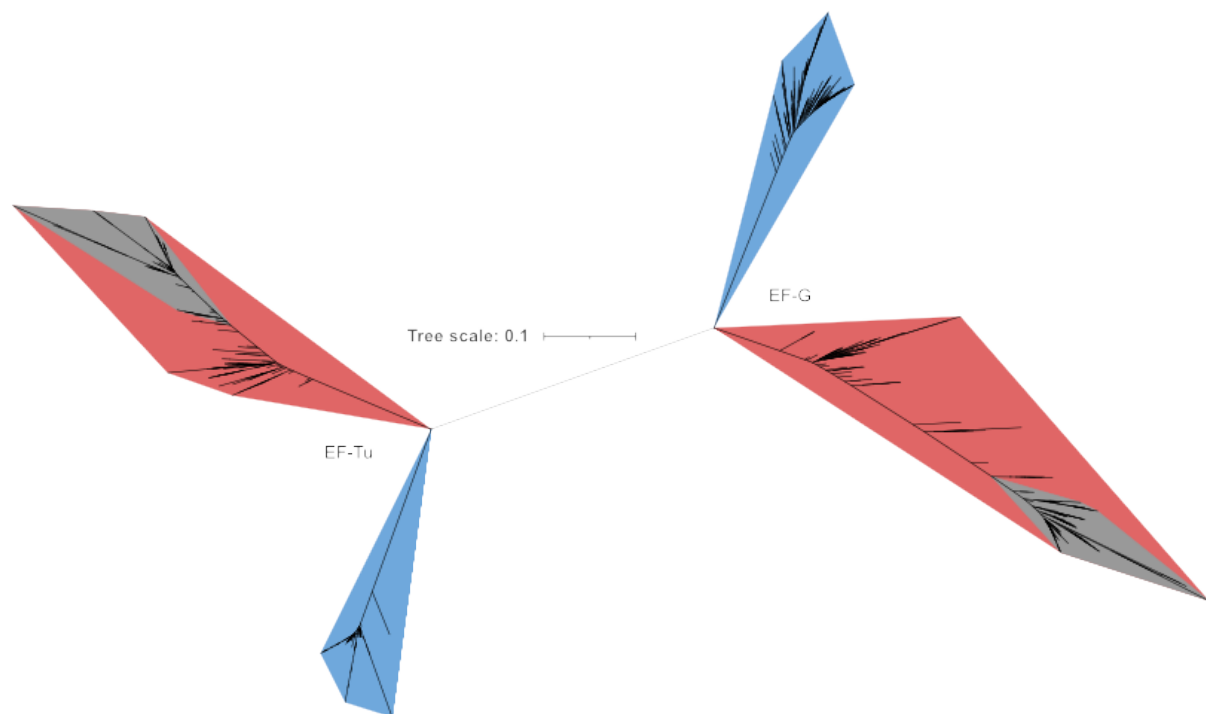522
523 *Homology Modelling*
524 For the KaiB and RPC10 proteins homologs were identified via BLAST as described above.
525 Then PDB structures or AlphaFold structures of the two conformations in question were
526 used as a template in SWISS-MODEL (Waterhouse et al., 2018). KaiB was modelled using
527 the PDB structure 2QKE in the ground state and 5JYT in the fold-switched state from
528 *Thermosynechococcus elongatus*. The RPC10 was homology modelled on the PdB
529 structure 7AE1 in the OUT conformation and 7AE3 in the IN conformation as described in
530 (Girbig et al., 2021). 3Di sequences were extracted from both sets of states/conformations
531 and then randomly sampled to generate a set composed approximately 50% of 3Di
532 sequences from PDB of KaiB and RPC10 in one of the two states/conformation. ML trees
533 were then estimated using these proteins sequences as described above.
534
535
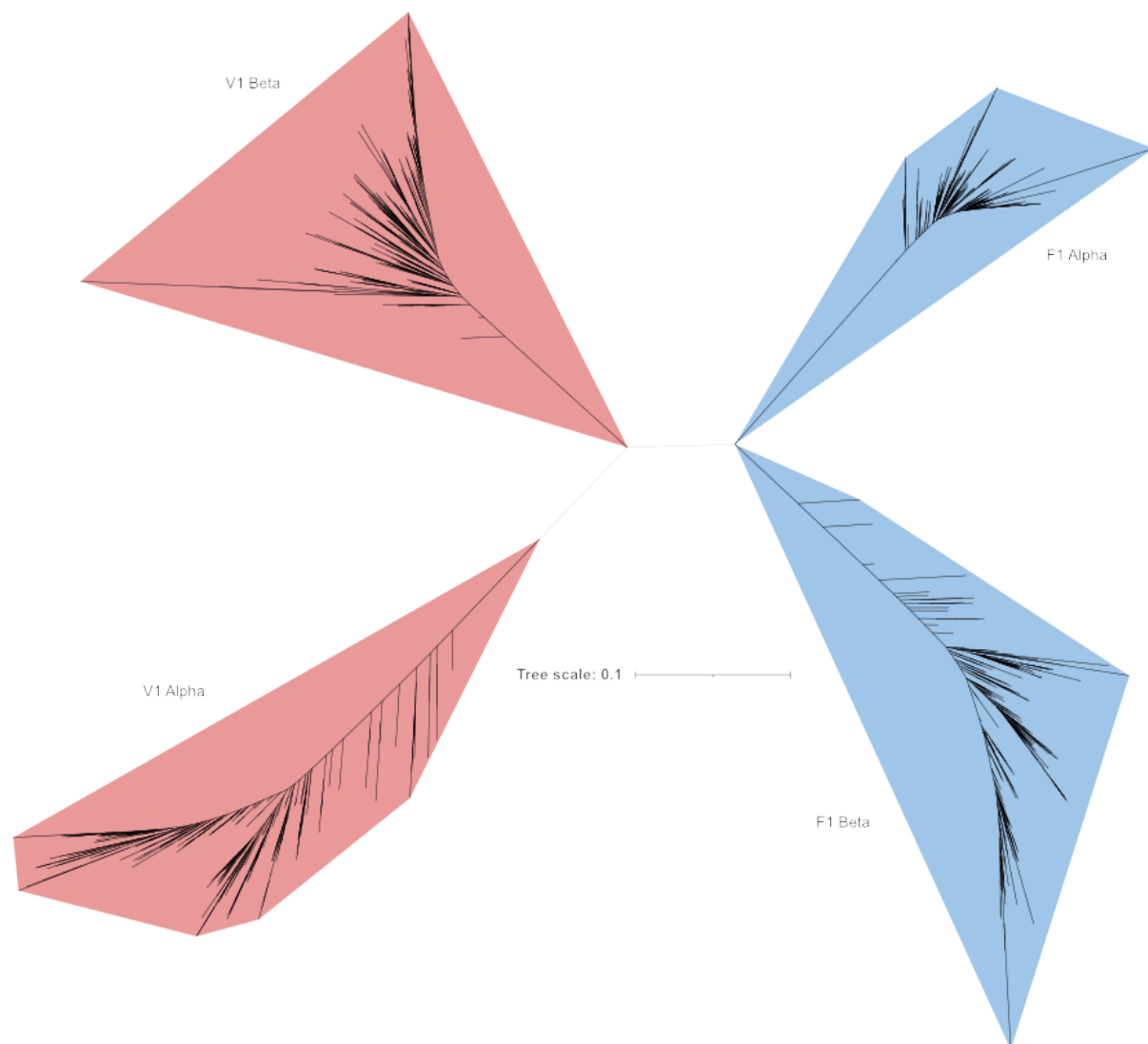536 **Supplementary Figure Legends**

**Supplementary Figure 1:** (A-C) Average Percentage similarity between 10 independent rounds of 3Di translations using the ProtT5 model for Elongation factors, ATPase subunits and the Reaction Centre I proteins respectively. (D-F) Percentage similarity between 3Di translation using the ProtT5 model and 3Di sequences extracted from AlphaFold predicted structures for Elongation factors, ATPase subunits and the Reaction Centre I proteins respectively. In all cases percentage similarities were calculated based on the BLOSUM style 3Di scoring matrix on unaligned sequences. Results show that the ProtT5 model is more precise than it is accurate when compared to AlphaFold predictions in all the three cases tested
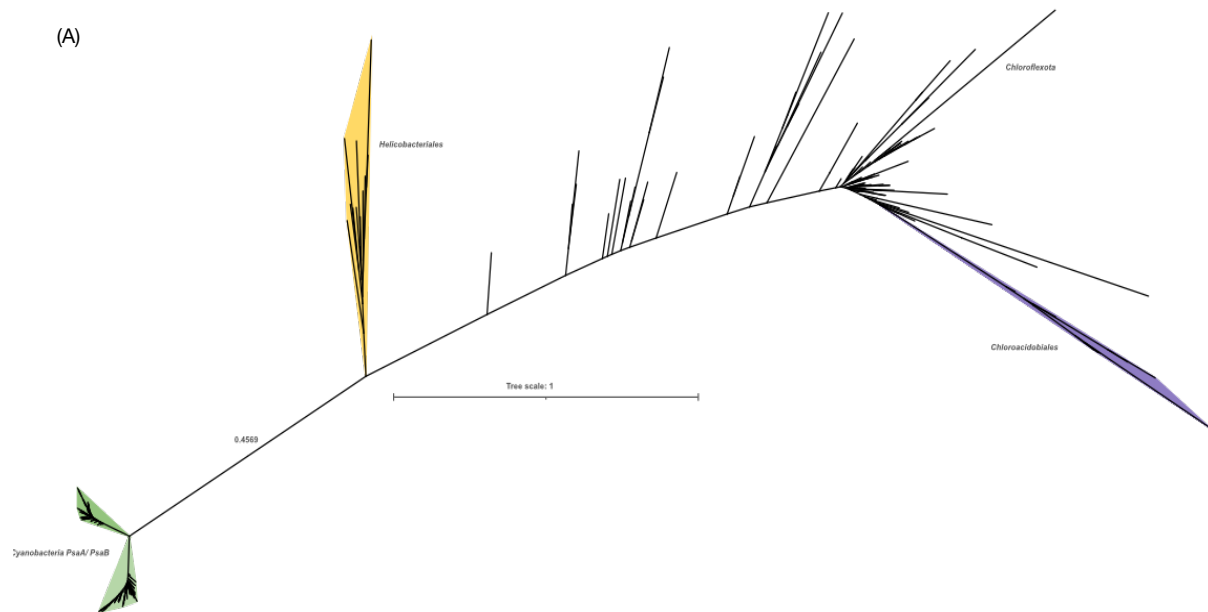
548
549 **Supplementary Figure 2:** (A) Sequences shown are 3Di translations of representative RC1
550 proteins prefixed with "ProtT5_" while the 3Di sequences extracted from AlphaFold
551 structures are suffixed with "_AF". The first observation is the obvious dissimilarities
552 between the ProtT5 (LLM) predictions and the AF extractions. The second observation is
553 that in some cases this erroneous insertions of stretches of 3Di characters are responsible
554 for mismatches observed downstream. This suggests that while the predictions of the
555 ProtT5 LLM is wrong, an alignment guided using the 3Di scoring matrix will be able to align
556 portions of the predictions. This explains the remarkable similarities between the Q-
557 matrices estimated from the AlphaFold predictions and the LLM.
558

559
560 **Supplementary Figure 3:** (A) 3Di (structural) ML tree on 3Di translations using ProtT5 of
561 Elongation factor proteins. Red, Blue, and Grey represent Archaeal, Bacterial, and
562 Eukaryotic groups respectively. This particular tree recovers the two-domain topology for
563 the tree of life albeit consistent with the 3Di (structural) ML tree estimated from 3Di
564 sequences extracted from AlphaFold structures.

565
566 **Supplementary Figure 4:** (A) 3Di (structural) ML tree on 3Di translations using ProtT5 of
567 ATPase subunits. Red Blue, and Grey represent Archaeal, Bacterial, and Eukaryotic groups
568 respectively. V1 Alpha and F1 Beta are the catalytic subunits while V1 Beta and F1 Alpha
569 are non-catalytic. This tree recovers a root for the tree of life between archaea and bacteria.
570 It does, however, groups the respective catalytic and non-catalytic subunits of bacteria
571 together, as well as the catalytic and non-catalytic subunits of archaea. This would require
572 an independent loss of catalytic activities in one of the subunits in both the groups. This is
573 inconsistent with currently established theories on the origin of the rotary ATPase. For
574 comparison, our structural tree derived from AlphaFold predictions (Figure 3C) groups
575 archaeal and bacterial catalytic subunits as one monophyletic group and the non-catalytic
576 subunits as another.

577
578 **Supplementary Figure 5:** (A) 3Di (structural) ML tree on 3Di translations using ProtT5 of
579 ATPase subunits. This particular tree is highly inconsistent and does not recover the split
580 between chlorobiales and chloroacidobiales. This is also evident from the particularly low
581 similarities between the 3Di translations using ProtT5 and the 3Di sequences extracted
582 form AlphaFold structures. Such cases highlight the importance of the quality of the
583 structure predictions.
584

585 **Author Contributions**
586 SGG and GKAH conceptualized, designed and wrote the manuscript. SGG performed the
587 computations
588

593

594 **Data availability**
595 All datasets, trees and alignments are available at
596 https://edmond.mpg.de/privateurl.xhtml?token=624d9e21-f33d-408b-8a81-
597 93d9ad020426 for review and will be made fully public on publication.
598

599 **References**
600

601 Balaji S, Srinivasan N. 2001. Use of a database of structural alignments and phylogenetic
602 trees in investigating the relationship between sequence and structural variability among
603 homologous proteins. *Protein Eng*. 14:219–226.

604 Balaji S, Sujatha S, Kumar SSC, Srinivasan N. 2001. PALI—a database of Phylogeny and
605     ALIgnment of homologous protein structures. *Nucleic Acids Res.* 29:61–65.

606  Baldauf SL, Palmer JD, Doolittle WF. 1996. The root of the universal tree and the origin
607      of eukaryotes based on elongation factor phylogeny. *Proc. Natl. Acad. Sci.* 93:7749–
608      7754.

609  Brown JR, Doolittle WF. 1995. Root of the universal tree of life based on ancient
610      aminoacyl-tRNA synthetase gene duplications. *Proc. Natl. Acad. Sci.* 92:2441–2445.

611  Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for automated
612      alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25:1972–
613      1973.

614  Chakravarty D, Porter LL. 2022. AlphaFold2 fails to predict protein fold switching.
615      *Protein Sci.* 31:e4353.

616  Chang Y-G, Cohen SE, Phong C, Myers WK, Kim Y-I, Tseng R, Lin J, Zhang L, Boyd JS, Lee
617      Y, et al. 2015. A protein fold switch joins the circadian oscillator to clock output in
618      cyanobacteria. *Science* 349:324–328.

619  Cross RL, Müller V. 2004. The evolution of A-, F-, and V-type ATP synthases and
620      ATPases: reversals in function and changes in the $H^+$/ATP coupling ratio. *FEBS Lett.*
621      576:1–4.

622  Farris JS. 1970. Methods for computing wagner trees. *Syst. Biol.* 19:83–92.

623  Felsenstein J. 1981. Evolutionary trees from DNA sequences: A maximum likelihood
624      approach. *J. Mol. Evol.* 17:368–376.

625  Felsenstein J. 2003 Inferring phylogenies. *Sinauer Associates, Sunderland,*
626      *Massachusetts*.

627  Fitch WM. 1971. Toward defining the course of evolution: minimum change for a
628      specific tree topology. *Syst. Zoöl.* 20:406.

629  Gilchrist CLM, Mirdita M, Steinegger M. 2024. Multiple protein structure alignment at
630      scale with FoldMason. *bioRxiv*:2024.08.01.606130.

631  Girbig M, Misiaszek AD, Vorländer MK, Lafita A, Grötsch H, Baudin F, Bateman A, Müller
632      CW. 2021. Cryo-EM structures of human RNA polymerase III in its unbound and
633      transcribing states. *Nat. Struct. Mol. Biol.* 28:210–219.

634  Gogarten JP, Kibak H, Dittrich P, Taiz L, Bowman EJ, Bowman BJ, Manolson MF, Poole
635      RJ, Date T, Oshima T, et al. 1989. Evolution of the vacuolar $H^+$-ATPase: implications
636      for the origin of eukaryotes. *Proc. Natl. Acad. Sci. USA.* 17:6661-5.

637  Gouy R, Baurain D, Philippe H. 2015. Rooting the tree of life: the phylogenetic jury is still
638      out. *Philos. Trans. R. Soc. B: Biol. Sci.* 370:20140329.

639 Grüber G, Wieczorek H, Harvey WR, Müller V. 2001. Structure–function relationships of
640     A-, F- and V-ATPases. *J. Exp. Biol.* 204:2597–2605.

641 Heinzinger M, Weissenow K, Sanchez JG, Henkel A, Steinegger M, Rost B. 2023. ProstT5:
642     bilingual language model for protein sequence and structure.
643     *bioRxiv*:2023.07.23.550085.

644 Hohmann-Marriott MF, Blankenship RE. 2011. Evolution of photosynthesis. *Plant Biol.*
645     62:515–548.

646 Iwabe N, Kuma K, Hasegawa M, Osawa S, Miyata T. 1989. Evolutionary relationship of
647     archaebacteria, eubacteria, and eukaryotes inferred from phylogenetic trees of
648     duplicated genes. *Proc. Natl. Acad. Sci.* 86:9355–9359.

649 Johnson MS, Šali A, Blundell TL. 1990. [42] Phylogenetic relationships from three-
650     dimensional protein structures. *Methods Enzym.* 183:670–690.

651 Johnson MS, Sutcliffe MJ, Blundell TL. 1990. Molecular anatomy: Phyletic relationships
652     derived from three-dimensional structures of proteins. *J. Mol. Evol.* 30:43–59.

653 Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, Tunyasuvunakool K,
654     Bates R, Žídek A, Potapenko A, et al. 2021. Highly accurate protein structure
655     prediction with AlphaFold. *Nature* 596:583–589.

656 Kalyaanamoorthy S, Minh BQ, Wong TKF, Haeseler A von, Jermiin LS. 2017.
657     ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat.*
658     *Methods* 14:587–589.

659 Kempen M van, Kim SS, Tumescheit C, Mirdita M, Lee J, Gilchrist CLM, Söding J,
660     Steinegger M. 2023. Fast and accurate protein structure search with Foldseek. *Nat.*
661     *Biotechnol.*:1–4.

662 Le SQ, Gascuel O. 2008. An improved general amino acid replacement matrix. *Mol.*
663     *Biol. Evol.* 25:1307–1320.

664 Liò P, Goldman N. 1998. Models of molecular evolution and phylogeny. *Genome Res.*
665     8:1233–1244.

666 Lopez P, Casane D, Philippe H. 2002. Heterotachy, an important process of protein
667     evolution. *Mol. Biol. Evol.* 19:1–7.

668 Magee AF, Hilton SK, DeWitt WS. 2021. Robustness of phylogenetic inference to model
669     misspecification caused by pairwise epistasis. *Mol. Biol. Evol.* 38:4603–4615.

670 Mahendrarajah TA, Moody ERR, Schrempf D, Szánthó LL, Dombrowski N, Davín AA,
671     Pisani D, Donoghue PCJ, Szöllősi GJ, Williams TA, et al. 2023. ATP synthase evolution
672     on a cross-braced dated tree of life. *Nat. Commun.* 14:7456.

673   Martin WF, Bryant DA, Beatty JT.  2018. A physiological perspective on the origin and
674        evolution of photosynthesis. *FEMS Microbiol. Rev.* 2:205-231

675   Mau B, Newton MA. 1997. Phylogenetic inference for binary data on dendograms using
676        markov chain monte carlo. *J. Comput. Graph. Stat.* 6:122–131.

677   Mihaescu R, Levy D, Pachter L. 2009. Why neighbor-joining works. *Algorithmica* 54:1–
678        24.

679   Miller DL. 1972. Elongation factors EF Tu and EF G interact at related sites on
680        ribosomes. *Proc. Natl. Acad. Sci. USA.* 3:752-5.

681   Minh BQ, Dang CC, Vinh LS, Lanfear R. 2021. QMaker: Fast and accurate method to
682        estimate empirical models of protein evolution. *Syst. Biol.* 70:syab010-.

683   Moi D, Bernard C, Steinegger M, Nevers Y, Langleib M, Dessimoz C. 2023. Structural
684        phylogenetics unravels the evolutionary diversification of communication systems in
685        gram-positive bacteria and their viruses. *bioRxiv*:2023.09.19.558401.

686   Mutti G, Ocaña-Pallarés E, Gabaldón T. 2024. Newly developed structure-based
687        methods do not outperform standard sequence-based methods for large-scale
688        phylogenomics. *bioRxiv*:2024.08.02.606352

689   Philippe H, Brinkmann H, Lavrov DV, Littlewood DTJ, Manuel M, Wörheide G, Baurain D.
690        2011. Resolving difficult phylogenetic questions: Why more sequences are not
691        enough. *PLoS Biol.* 9:e1000602.

692   Philippe H, Forterre P. 1999. The rooting of the universal tree of life is not reliable. *J. Mol.
693        Evol.* 49:509–523.

694   Posada D, Crandall KA. 2021. Felsenstein phylogenetic likelihood. *J. Mol. Evol.* 89:134–
695        145.

696   Puente-Lelievre C, Malik AJ, Douglas J, Ascher D, Baker M, Allison J, Poole A, Lundin D,
697        Fullmer M, Bouckert R, et al. 2024. Tertiary-interaction characters enable fast,
698        model-based structural phylogenetics beyond the twilight zone.
699        *bioRxiv*:2023.12.12.571181.

700   Rannala B, Yang Z. 1996. Probability distribution of molecular evolutionary trees: A new
701        method of phylogenetic inference. *J. Mol. Evol.* 43:304–311.

702   Saitou N, Nei M. 1987. The neighbor-joining method: a new method for reconstructing
703        phylogenetic trees. *Mol. Biol. Evol.* 4:406–425.

704   Sala D, Engelberger F, Mchaourab HS, Meiler J. 2023. Modeling conformational states
705        of proteins with AlphaFold. *Curr. Opin. Struct. Biol.* 81:102645.

706 Sánchez-Baracaldo, P. and Cardona, T. 2020. On the origin of oxygenic photosynthesis
707     and cyanobacteria. *New Phytol.* 225: 1440-1446.

708 Starr TN, Thornton JW. 2016. Epistasis in protein evolution. *Protein Sci.* 25:1204–1218.

709 Tsuji JM, Shaw NA, Nagashima S, Venkiteswaran JJ, Schiff SL, Watanabe T, Fukui M,
710     Hanada S, Tank M, Neufeld JD. 2024. Anoxygenic phototroph of the Chloroflexota
711     uses a type I reaction centre. *Nature* 627:915–922.

712 Varadi M, Bertoni D, Magana P, Paramval U, Pidruchna I, Radhakrishnan M, Tsenkov M,
713     Nair S, Mirdita M, Yeo J, et al. 2023. AlphaFold Protein Structure Database in 2024:
714     providing structure coverage for over 214 million protein sequences. *Nucleic Acids*
715     *Res.* 52:D368–D375.

716 Waterhouse A, Bertoni M, Bienert S, Studer G, Tauriello G, Gumienny R, Heer FT, Beer
717     TAP de, Rempfer C, Bordoli L, et al. 2018. SWISS-MODEL: homology modelling of
718     protein structures and complexes. *Nucleic Acids Res.* 46:W296–W303.

719 Wayment-Steele HK, Ojoawo A, Otten R, Apitz JM, Pitsawong W, Hömberger M,
720     Ovchinnikov S, Colwell L, Kern D. 2023. Predicting multiple conformations via
721     sequence clustering and AlphaFold2. *Nature*:1–3.

722 Whelan S, Liò P, Goldman N. 2001. Molecular phylogenetics: state-of-the-art methods
723     for looking into the past. *Trends Genet.* 17:262–272.

724 Zhang N, Guan W, Cui S, Ai N. 2023. Crowded environments tune the fold-switching in
725     metamorphic proteins. *Commun. Chem.* 6:117.

726
727
728
729
730
731
732