

# 1 What is an archaeon and are the 2 Archaea really unique?

3 **Ajith Harish**<sup>1\*</sup>

\*For correspondence:  
ajith.harish@gmail.com

4 <sup>1</sup>Department of Cell and Molecular Biology, Structural and Molecular Biology Program,  
5 Uppsala University, Uppsala, Sweden

---

7 **Abstract** The recognition of the group Archaea 40 years ago stimulated research in microbial  
8 evolution and molecular systematics that prompted a new classificatory scheme to organize  
9 biodiversity. Advances in DNA sequencing techniques have since significantly improved the  
10 genomic representation of the archaeal biodiversity. In addition, advances in phylogenetic  
11 modeling that facilitate large-scale phylogenomics have resolved many recalcitrant branches of the  
12 Tree of Life. Despite the technical advances and an expanded taxonomic representation, two  
13 important aspects of the origins and evolution of the Archaea remain controversial, even as we  
14 celebrate the 40th anniversary of the monumental discovery. The issues concern (i) the uniqueness  
15 (monophyly) of the Archaea, and (ii) the evolutionary relationships of the Archaea to the Bacteria  
16 and the Eukarya; both of these are relevant to the deep structure of the Tree of Life. The  
17 uncertainty is primarily due to a scarcity of information in standard datasets—the core-genes  
18 datasets—to reliably resolve the conflicts. These conflicts can be resolved efficiently by employing  
19 complex genomic features and genome-scale evolution models—a distinct class of phylogenomic  
20 characters and evolution models—that can be employed routinely to maximize the use of genome  
21 sequences as well as to minimize uncertainties in tests of evolutionary hypotheses.

---

## 23 Introduction

24 The recognition of the Archaea as the so-called “third form of life” was made possible in part by a  
25 new technology for sequence analysis, oligonucleotide cataloging, developed by Fredrik Sanger and  
26 colleagues in the 1960s (1, 2). Carl Woese’s insight of using this method, and the choice of the small  
27 subunit ribosomal RNA (16S/SSU rRNA) as a phylogenetic marker, not only put microorganisms  
28 on a phylogenetic map (or tree), but also revolutionized the field of molecular systematics that  
29 Zuckerkandl and Pauling has previously alluded to (3). Comparative analysis of organism-specific  
30 (oligonucleotide) sequence-signatures in SSU rRNA led to the recognition of a distinct group of  
31 microorganisms (2, 4). Initially referred to as Archaeobacteria, these unusual organisms had  
32 ‘oligonucleotide signatures’ distinct from other bacteria (Eubacteria), and they were later found to  
33 be different from those of Eukarya (eukaryotes) as well. Many other features, including molecular,  
34 biochemical as well as ecological, corroborated the uniqueness of the Archaea. Thus the archaeal  
35 concept was established (2).

36 The study of microbial diversity and evolution has come a long way since then: sequencing  
37 microbial genomes, and directly from the environment without the need for culturing is now  
38 routine (5, 6). This wealth of sequence information is exciting not only for cataloging and organizing  
39 biodiversity, but also to understand the ecology and evolution of microorganisms – archaea and  
40 bacteria as well as eukaryotes – that make up a vast majority of the planetary biodiversity. Since  
41 large-scale exploration by the means of environmental genome sequencing became possible almost  
42 a decade ago, there has also been a palpable excitement and anticipation of the discovery of a

43 fourth form of life or a “fourth domain” of life (7). The reference here is to a fourth form of cellular  
44 life, but not to viruses, which some have already proposed to be the fourth domain of the Tree of  
45 Life (ToL) (7, 8). If a fourth form of life were to be found, what would the distinguishing features be,  
46 and how could it be measured, defined and classified?

47 Rather than the discovery of a fourth domain, and contrary to the expectations, however, current  
48 discussion is centered around the return to a dichotomous classification of life (9-11), despite the  
49 rapid expansion of sequenced biodiversity – hundreds of novel phyla descriptions (12, 13). The  
50 proposed dichotomous classifications schemes, unfortunately, are in sharp contrast to each other,  
51 depending on: (i) whether the Archaea constitute a monophyletic group—a unique line of descent  
52 that is distinct from those of the Bacteria as well as the Eukarya; and (ii) whether the Archaea form  
53 a sister clade to the Eukarya or to the Bacteria. Both the issues stem from difficulties involved in  
54 resolving the deep branches of the ToL (10, 11, 14).

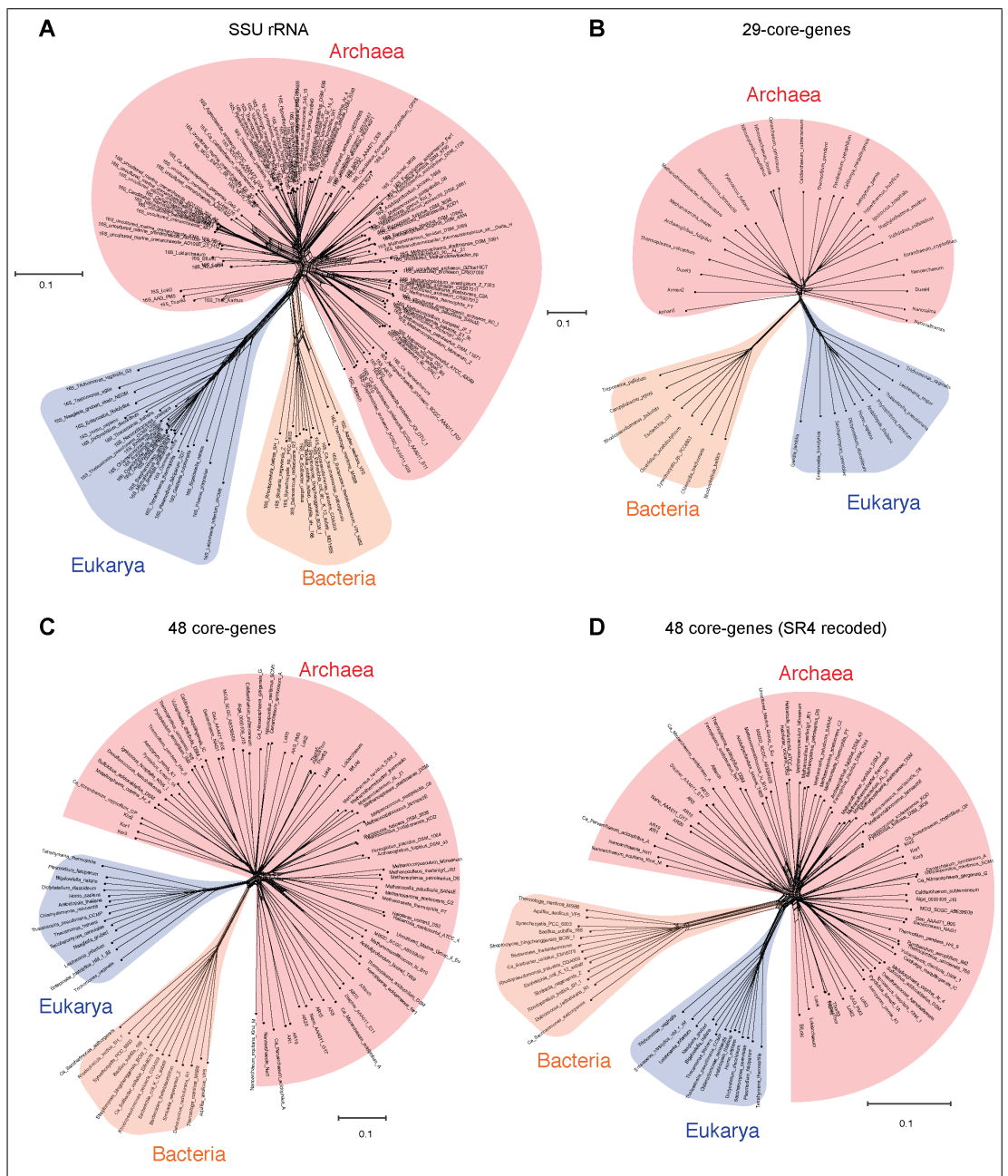
55 The twin issues, first recognized in the 80s based on single-gene (SSU rRNA) analyses, continue  
56 to be the subjects of a long-standing debate, which remains unresolved despite large-scale analyses  
57 of multi-gene datasets (5, 15-19). In addition to the choice of genes to be analyzed, the choice of  
58 the underlying character evolution model is at the core of contradictory results that either supports  
59 the Three-domains tree (5, 19) or the Eocyte tree (17, 20). In many cases, adding more data, either  
60 as enhanced taxon (species) sampling or enhanced character (gene) sampling, or both, can resolve  
61 ambiguities (21, 22). However, as the taxonomic diversity and evolutionary distance increases  
62 among the taxa studied, the number of conserved marker-genes that can be used for phylogenomic  
63 analyses decreases. Accordingly, resolving the phylogenetic relationships of the Archaea, Bacteria  
64 and Eukarya is restricted to a small set of genes—50 at most—in spite of the large increase in the  
65 numbers of genomes sequenced and the associated development of sophisticated phylogenomic  
66 methods.

67 Based on a closer scrutiny of the recent phylogenomic datasets employed in the ongoing  
68 debate, I will show here that one of the reasons for this persistent ambiguity is that the ‘information’  
69 necessary to resolve these conflicts is practically nonexistent in the standard marker-genes (i.e. core-  
70 genes) datasets employed routinely for phylogenomics. Further, I discuss analytical approaches  
71 that maximize the use of the information that is in genome sequence data and simultaneously  
72 minimize phylogenetic uncertainties. In addition, I discuss simple but important, yet undervalued,  
73 aspects of phylogenetic hypothesis testing, which together with the new approaches hold promise  
74 to resolve these long-standing issues effectively.

## 75 Results

### 76 Information in core genes is inadequate to resolve the archaeal radiation

77 Data-display networks (DDNs) are useful to examine and visualize character conflicts in phylogenetic  
78 datasets, especially in the absence of prior knowledge about the source of such conflicts, ideally  
79 before downstream processing of the data for phylogenetic inference (23, 24). While congruent  
80 data will be displayed as a tree in a DDN, incongruences are displayed as reticulations in the tree.  
81 Fig. 1A shows a neighbor-net analysis of the SSU rRNA alignment used to resolve the phylogenetic  
82 position of the recently discovered Asgard archaea (20). The DDN is based on character distances  
83 calculated as the observed genetic distance (p-distance) of 1,462 characters, and shows the total  
84 amount of conflict in the dataset that is incongruent with character bipartitions (splits). The edge  
85 (branch) lengths in the DDN correspond to the support for the respective splits. Accordingly, two  
86 well-supported sets of splits for the Bacteria and the Eukarya are observed. The Archaea, however,  
87 does not form a distinct, well-resolved/well-supported group, and is unlikely to correspond to a  
88 monophyletic group in a phylogenetic tree.



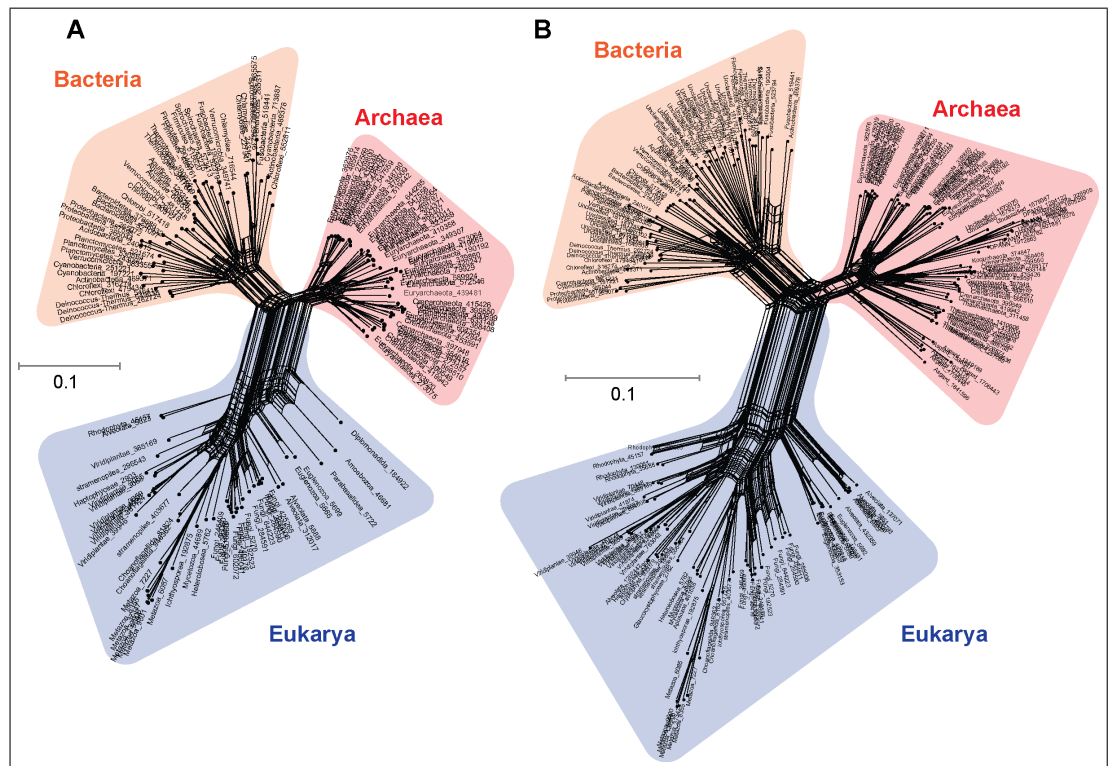
**Figure 1.** Data-display networks depicting the character conflicts in different datasets that employ different character types. (A) SSU rRNA alignment of 1,462 characters. Concatenated protein sequence alignment of (B) 29 core-genes, 8,563 characters; (C) 48 core-genes, 9,868 characters and (D) also 48 core-genes, 9,868 SR4 recoded characters (data simplified from 20 to 4 character-states). Each network is constructed from a neighbor-net analysis based on the observed genetic distance (p-distance) and displayed as an equal angle split network. Edge (branch) lengths correspond to the support for character bipartitions (splits), and reticulations in the tree correspond to character conflicts. Datasets in (A), (C) and (D) are from Ref. 20, and in (B) is from Ref. 17.

89 Likewise, the concatenated protein sequence alignment of the so-called 'genealogy defining core  
 90 of genes' (25) – a set of conserved single-copy genes – also does not support a unique archael lineage.  
 91 Fig. 1B is a DDN derived from a neighbor-net analysis of 8,563 characters in 29 concatenated core-  
 92 genes (17), while Fig. 1C,D is based on 9,868 characters in 44 concatenated core-genes (also from  
 93 (20)). Even taken together, none of the standard marker gene datasets are likely to support the  
 94 monophyly of the Archaea — a key assertion of the three-domains hypothesis (26). Simply put,  
 95 there is not enough information in the core-gene datasets to resolve the archael radiation, or to

96 determine whether the Archaea are really unique compared to the Bacteria and Eukarya. However,  
97 other complex features — including molecular, biochemical and phenotypic characters, as well  
98 as ecological adaptations — support the uniqueness of the Archaea. These idiosyncratic archaeal  
99 characters include the subunit composition of supramolecular complexes like the ribosome, DNA-  
100 and RNA-polymerases, biochemical composition of cell membranes, cell walls, and physiological  
101 adaptations to energy-starved environments, among other things (27, 28).

102 **Complex phylogenomic characters minimize uncertainties regarding the unique-**  
103 **ness of the Archaea**

104 A nucleotide is the smallest possible locus, and an amino acid is a proxy for a locus of a nucleotide  
105 triplet. Unlike the elementary amino acid- or nucleotide-characters in the core-genes dataset (Fig.1),  
106 the DDN in Fig. 2 is based on complex molecular characters – genomic loci that correspond to  
107 protein domains, typically ~200 amino acids (600 nucleotides) long. Neighbor-net analysis of protein-  
108 domain data coded as binary characters (presence/absence) is based on the Hamming distance  
109 (identical to the p-distance used in Fig.1). Here the Archaea also form a distinct well-supported  
110 cluster, as do the Bacteria and the Eukarya.



**Figure 2.** Data-display networks (DDN) depicting character conflicts among complex phylogenomic characters – genomic loci corresponding to protein-domains in this case. (A) Neighbor-net analysis based on Hamming distance (identical to the p-distance used in Fig.1) of 1,732 characters sampled from 141 species. (B) DDN based on an enriched taxon sampling of 81 additional species totaling 222 species and a modest increase to 1,738 characters. The dataset in (A) is from Ref. 10, which was updated with novel species to represent the recently described archaeal and bacterial species (5, 12, 20).

111 Fig 2A is a DDN based on the dataset that includes protein-domain cohorts of 141 species, used  
112 in a phylogenomic analysis to resolve the uncertainties at the root of the ToL (29). Compared to  
113 the data in Fig. 1, the taxonomic diversity sampled for the Bacteria and Eukarya is more extensive,  
114 but less extensive for the Archaea; it is composed of the traditional groups Euryarchaeota and  
115 Crenarchaeota. Fig. 2B is a DDN of an enriched sampling of 81 additional species, which includes  
116 representatives of the newly described archaeal groups: TACK (30), DPANN (5), and Asgard group

117 including the Lokiarchaeota (20). In addition, species sampling was enhanced with representatives  
118 from the candidate phyla described for Bacteria, and with unicellular species of Eukarya. The  
119 complete list species analyzed is in SI Table 1.

120 Notably, the extension of the protein-domain cohort was insignificant, from 1,732 to 1,738  
121 distinct domains (characters). Based on the well-supported splits in the DDN that form a distinct  
122 archaeal cluster, the Archaea are likely to be a monophyletic group (clade) in phylogenies inferred  
123 from these datasets.

#### 124 **Data quality affects model complexity required to explain phylogenetic datasets**

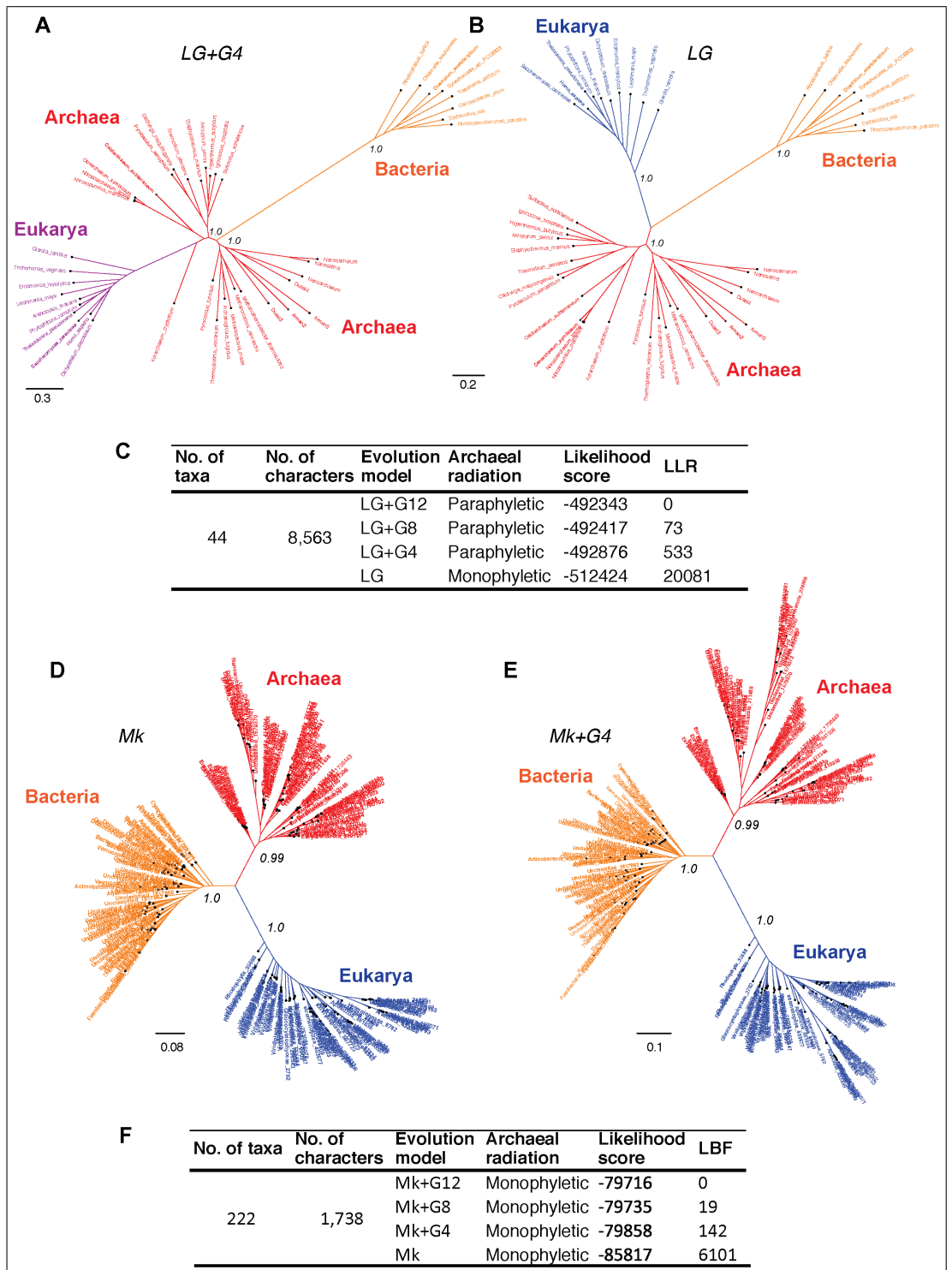
125 Resolving the paraphyly or monophyly of the Archaea is relevant to determining whether the Eocyte  
126 tree (Fig. 3A) or the Three-domains tree (Fig. 3B), respectively, is a better-supported hypothesis.  
127 Recovering the Eocyte tree typically requires implementing complex models of sequence evolution  
128 rather than their relatively simpler versions (11). In general, complex models tend to fit the data  
129 better. For instance, according to a model selection test for the 29 core-genes dataset, the LG  
130 model (31) of protein sequence evolution is a better-fitting model than other standard models,  
131 such as the WAG or JTT substitution model (SI-Table 2), as reported previously (17). Further, a  
132 relatively more complex version of the LG model, with multiple rate-categories was found to be a  
133 better-fitting model than the simpler single-rate-category model (Fig. 3C; SI-Table 2). The fit of the  
134 data is estimated as the likelihood of the best tree given the model.

135 A complex, multiple rate-categories model accounts for site-specific substitution rate variation.  
136 Substitution-rate heterogeneity across different sites in the multiple-sequence alignment (MSA)  
137 was approximated using a discrete Gamma model with 4, 8 or 12 rate categories (LG+G4, LG+G8 or  
138 LG+G12, respectively). The Archaea is consistent with a paraphyletic group in trees derived from the  
139 rate-heterogeneous versions of the LG model (Fig. 3A). Furthermore, the fit of the data improves  
140 with the increase in complexity of the substitution model (Fig. 3C). Model complexity increases  
141 with any increase in the number of rate categories and/or the associated numbers of parameters  
142 that need to be estimated. However, with a relatively simpler version – a rate-homogeneous LG  
143 model, in which the substitution-rates are approximated to a single rate-category, the Archaea are  
144 consistent with a monophyletic group (Fig. 3B).

145 In contrast, trees inferred from the protein-domain datasets are consistent with monophyly  
146 of the Archaea irrespective of the complexity of the underlying model (Fig. 3D-F). The Mk model  
147 (Markov k model) is the best-known probabilistic model of discrete character evolution, particularly  
148 of complex characters coded as binary-state characters (32, 33). Since the Mk model assumes a  
149 stochastic process of evolution, it is able to estimate multiple state changes along the same branch.  
150 Implementing a simpler rate-homogeneous version of the Mk model (Fig. 3D), as well as more  
151 complex rate-heterogeneous versions with 4, 8 or 12 rate categories (Mk+G4, Mk+G8 or Mk+G12,  
152 respectively), also recovered trees that are consistent with the monophyly of the Archaea (Fig. 3E)  
153 The tree derived from the Mk+G4 model is shown in Fig. 3E. While the tree derived from Mk+G8  
154 model is identical (SI-Fig. 1) to the Mk+G4 tree, the Mk+G12 tree is almost identical with minor  
155 differences in the bacterial sub-groups (SI-Fig. 2)

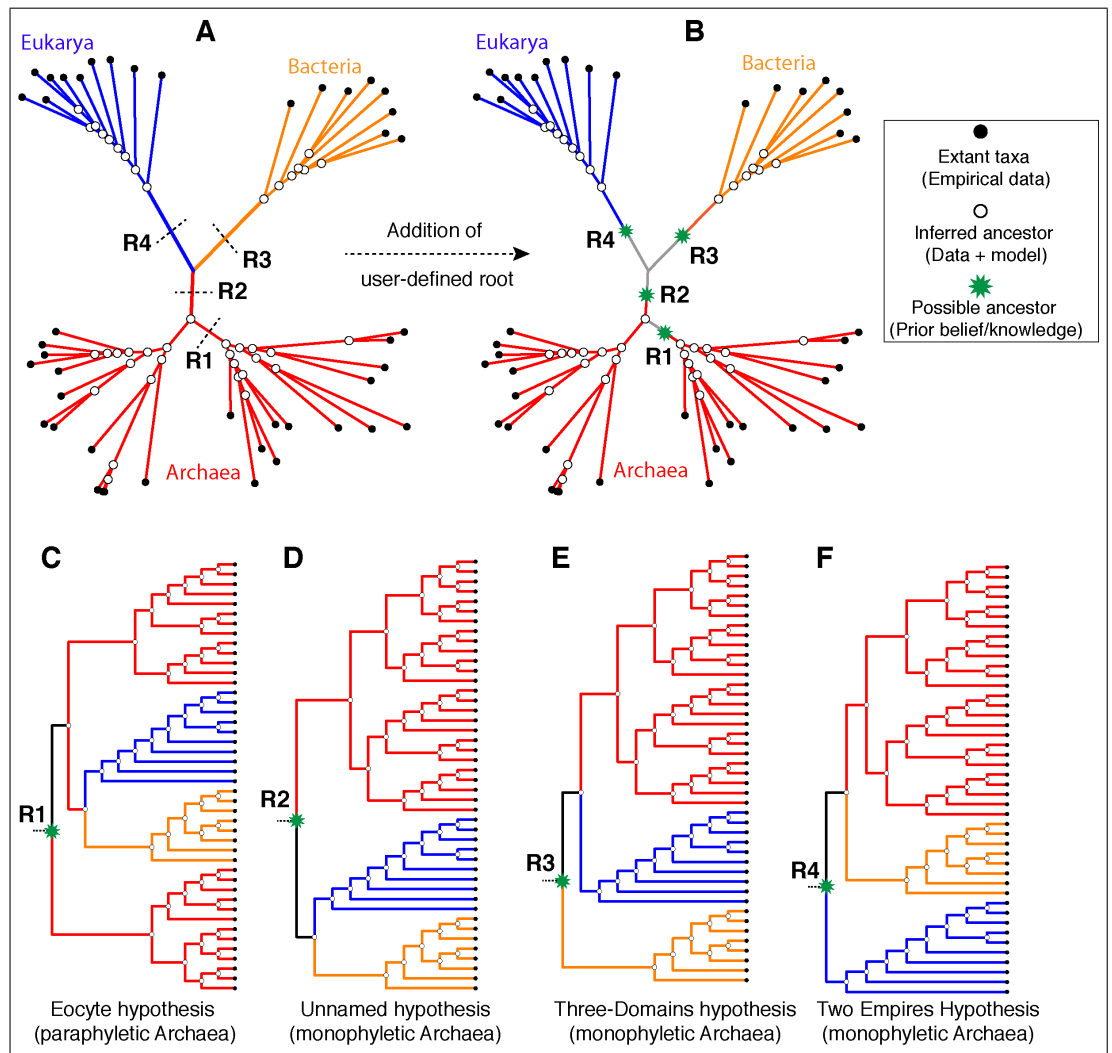
156 In all cases, bipartitions for Archaea show strong support with posterior probability (PP) of 0.99  
157 while that of Bacteria and Eukarya is supported with a PP of 1.0; in spite of substantially different  
158 fits of the data. The uniqueness of the Archaea is almost unambiguous in this case (but see next  
159 section).





**Figure 3.** Comparison of concatenated-gene trees derived from amino acid characters and genome trees derived from protein-domain characters. Branch support is shown only for the major branches. Scale bars represent the expected number of changes per character. (A), (B) Core-genes-tree derived from a better-fitting model (LG+G4) and a worse fitting mode (LG), respectively, of amino acid substitutions. (C) Model fit to data is ranked according the log likelihood ratio (LLR) scores. LLR scores are computed as the difference from the best-fitting model (LG+G12) of the likelihood scores estimated in PhyML. Thus, larger LLR values indicate less support for that model/tree relative to the most-likely model/tree. Substitution rate heterogeneity is approximated with 4, 8 or 12 rate categories in the complex models, but with a single rate category in the simpler model. (D), (E) are genome-trees derived from a better-fitting model (Mk+G4) and a worse fitting model (Mk), respectively. (F) Model fit to data is ranked according log Bayes factor (LBF) scores, which like LLR scores are the log odds of the hypotheses. LBF scores are computed as the difference in likelihood scores estimated in MrBayes.

160 **Siblings and cousins are indistinguishable when reversible models are employed**  
 161 Although a DDN is useful to identify and diagnose character conflicts in phylogenetic datasets and  
 162 to postulate evolutionary hypotheses, a DDN by itself cannot be interpreted as an evolutionary  
 163 network, because the edges do not necessarily represent evolutionary phenomena and the nodes  
 164 do not represent ancestors (23, 24). Therefore, evolutionary relationships cannot be inferred from  
 165 a DDN. Likewise, evolutionary relationships cannot be inferred from unrooted trees, even though  
 166 nodes in an unrooted tree do represent ancestors and an evolution model defines the branches  
 167 (see Fig. 4A).



**Figure 4.** Effect of alternative *ad hoc* rootings on the phylogenetic classification of archaeal biodiversity. (A) An unrooted tree is not fully resolved into bipartitions at the root of the tree (i.e. a polytomous rather than a dichotomous root branching) and thus precludes identification of sister group relationships. It is common practice to add a user-specified root *a posteriori* based on prior knowledge (or belief) of the investigator. Four possible (of many) rootings R1-R4 are shown. (B) Operationally, adding a root (rooting) *a posteriori* amounts to adding new information – a new bipartition and an ancestor as well as an evolutionary polarity – that is independent of the source data. (C-F) The different possible evolutionary relationships of the Archaea to other taxa, depending on the position of the root, are shown. Rooting is necessary to determine the recency of common ancestry as well the temporal order of key evolutionary transitions that define phylogenetic relationships.

168 An unrooted tree, unlike a rooted tree, is not an evolutionary (phylogenetic) tree *per se*, since it  
 169 is a minimally defined hypothesis of evolution or of relationships; it is, nevertheless, useful to rule

170 out many possible bipartitions and groups (34, 35). Given that a primary objective of phylogenetic  
171 analyses is to identify clades and the relationships between these clades, it is not possible to  
172 interpret an unrooted tree meaningfully without rooting the tree (see Fig. 4A). Identifying the root is  
173 essential to: (i) distinguish between ancestral and derived states of characters, (ii) determine the  
174 ancestor-descendant polarity of taxa, and (iii) diagnose clades and sister-group relationships (Fig.  
175 4). Yet, most phylogenetic software construct only unrooted trees, which are then consistent with  
176 several rooted trees (Fig. 4 C-F). However, an unrooted tree cannot be fully resolved into bipartitions,  
177 because an unresolved polytomy (a trifurcation in this case) exists near the root of the tree (Fig. 4A),  
178 which otherwise corresponds to the deepest split (root) in a rooted tree (Fig. 4, C-F).

179 Resolving the polytomy requires identifying the root of the tree. The identity of the root  
180 corresponds, in principle, to any one of the possible ancestors as follows:

- 181 i. Any one of the inferred-ancestors at the resolved bipartitions (open circles in Fig. 4A), or
- 182 ii. Any one of the yet-to-be-inferred-ancestors that lies along the stem-branches of the unre-  
183 solved polytomy (dashed lines in Fig. 4A) or along the internal-branches.

184 In the latter case, rooting the tree *a posteriori* on any of the branches amounts to inserting an ad-  
185 ditional bipartition and an ancestor that is neither inferred from the source data nor deduced from  
186 the underlying character evolution model. Since standard evolution models employed routinely  
187 cannot resolve the polytomy, rooting, and hence interpreting the Tree of Life depends on:

- 188 i. Prior knowledge — eg., fossils or a known sister-group (outgroup), or
- 189 ii. Prior beliefs/expectations of the investigators — eg., simple is primitive (36, 37), bacteria are  
190 primitive (38, 39), archaea are primitive (1), etc.

191 Both of these options are independent of the data used to infer the unrooted ToL. Some possible  
192 rootings and the resulting rooted-tree topologies are shown as cladograms in Fig. 4, C-F. If the root  
193 lies on any of the internal branches (e.g. R1 in Fig. 4A-C), or corresponds to one of the internal  
194 nodes, within the archael radiation, the Archaea would not constitute a unique clade (Fig. 4C).  
195 However, if the root lies on one of the stem-branches (R2/R3/R4 in Fig. 4 A, B), monophyly of  
196 the Archaea would be unambiguous (Fig. 4 D-F). Determining the evolutionary relationship of the  
197 Archaea to other taxa, though, requires identifying the root.

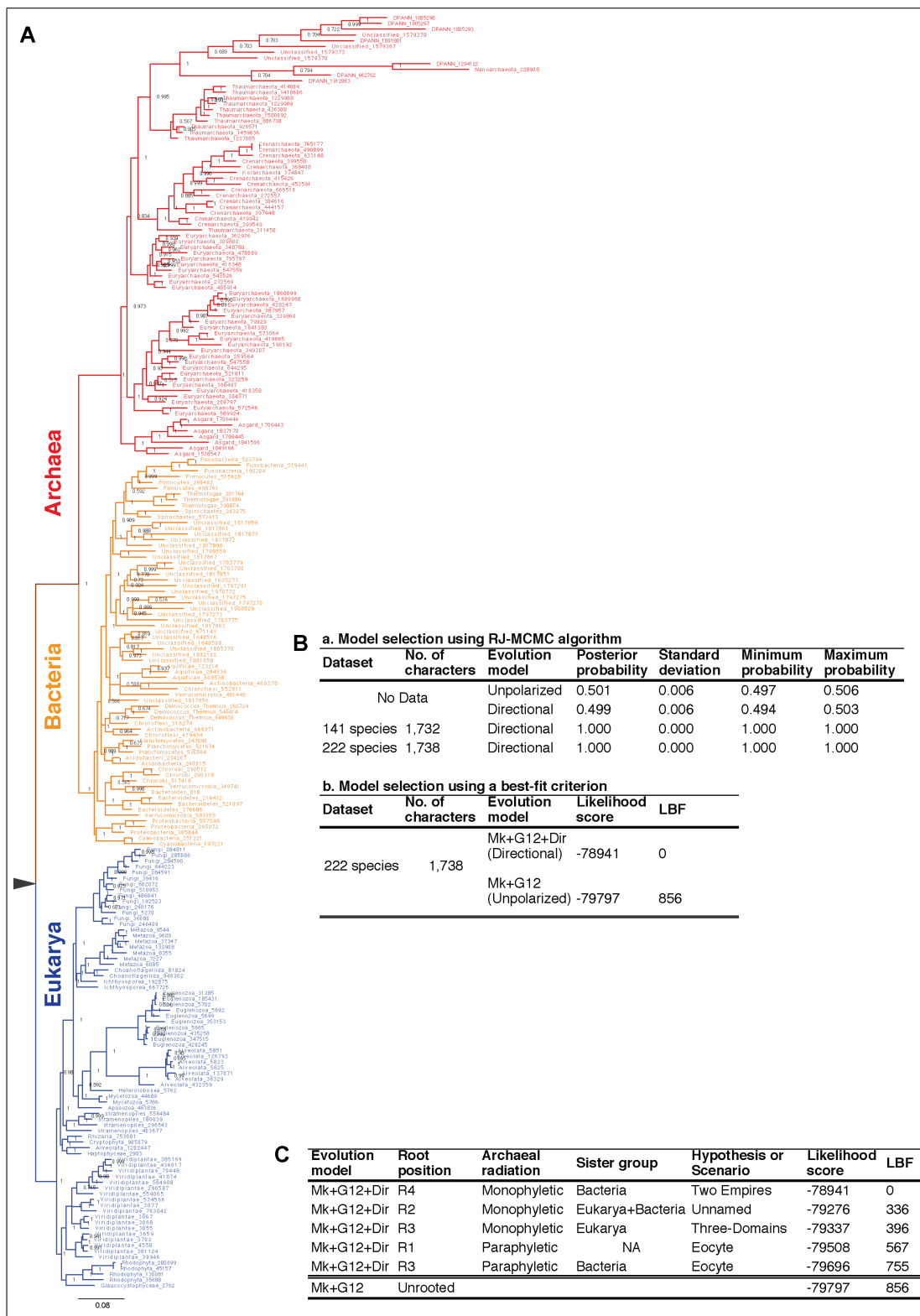
198 Directional evolution models, unlike reversible models, are able to identify the polarity of state  
199 transitions, and thus the root of a tree (40-42). Therefore, the uncertainty due to a polytomous root  
200 branching is not an issue (Fig 5A). Moreover, directional evolution models are useful to evaluate the  
201 empirical support for prior beliefs about the universal common ancestor (UCA) at the root of the  
202 ToL (29). A Bayesian model selection test implemented to detect directional trends (42) chooses the  
203 directional model, overwhelmingly (Fig. 5B), over the unpolarized model for the protein-domain  
204 dataset in Fig. 2B, as reported previously for the dataset in Fig. 2A (29). Further, the best-supported  
205 rooting corresponds to root R4 (Fig. 4F and Fig. 5A) — monophyly of the Archaea is maximally  
206 supported (PP of 1.0). Furthermore, the sister-group relationship of the Archaea to the Bacteria  
207 is maximally supported (PP 1.0). Accordingly, a higher order taxon, Akaryotes, proposed earlier  
208 (Forterre 1992) forms a well-supported clade. Thus Akaryotes (or Akarya) and Eukarya are sister  
209 clades that diverge from the UCA at the root of the ToL, also as reported previously (29).

210 Alternative rootings are much less likely, and are not supported (Fig. 5C). Accordingly, indepen-  
211 dent origin of the eukaryotes as well akaryotes is the best-supported scenario. The Three-domains  
212 tree (root R3, Fig. 4E) is  $10^{171}$  times less likely, and the scenario proposed by the Eocyte hypothesis  
213 (root R1, Fig. 5A) is highly unlikely. The common belief that simple is primitive, as well as beliefs that  
214 archaea are primitive or that archaea and bacteria evolved before eukaryotes, are not supported  
215 either.

## 216 **Employing complex molecular characters maximizes representation of orthologous, 217 non-recombining genomic loci, and thus phylogenetic signal**

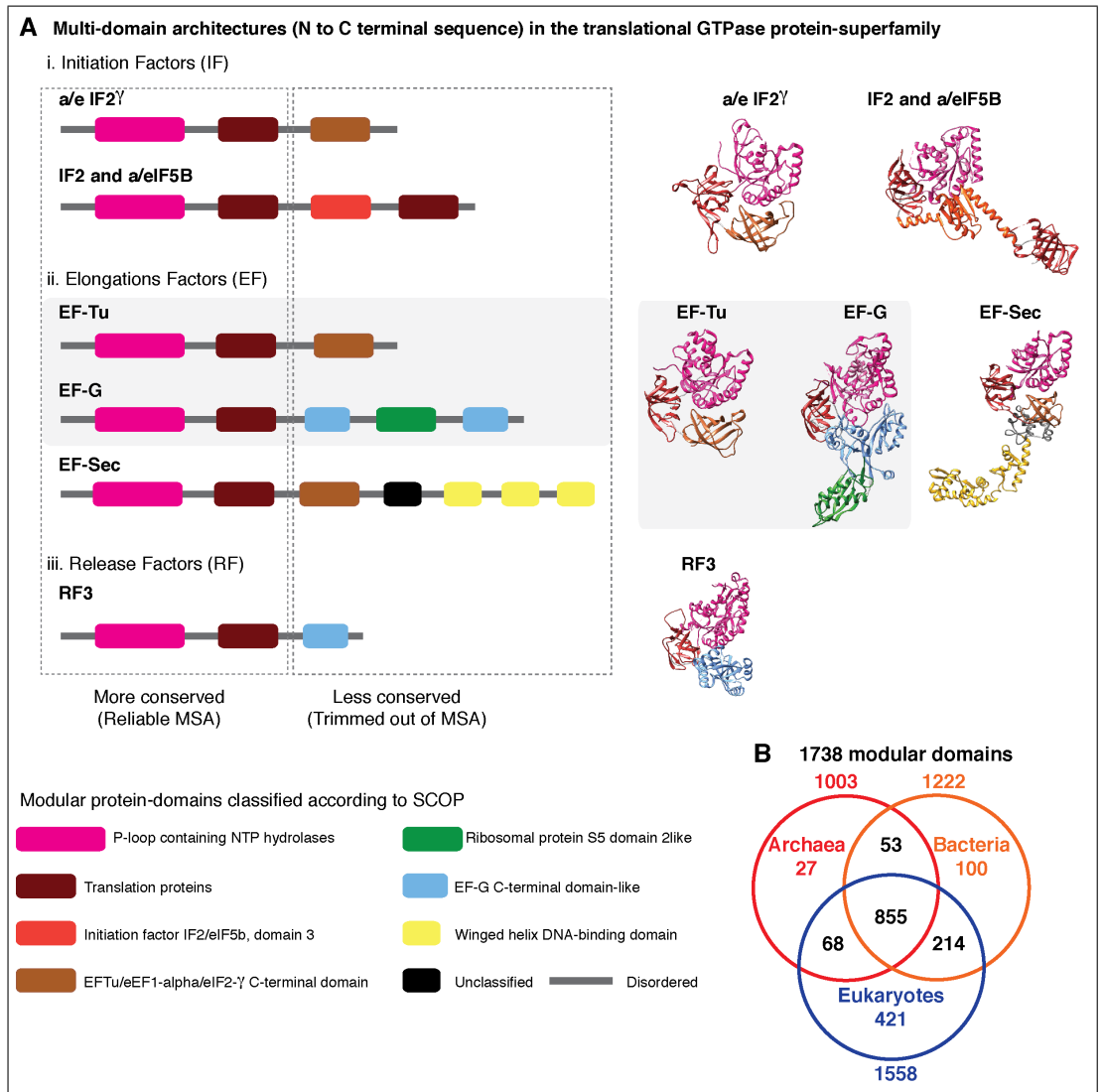
218 Genomic loci that can be aligned with high confidence using MSA algorithms are typically more  
219 conserved than those loci for which alignment uncertainty is high. Such ambiguously aligned





**Figure 5.** (A) Rooted tree of life inferred from patterns of inheritance of unique genomic-signatures. A dichotomous classification of the diversity of life such that Archaea is a sister group to Bacteria, which together constitute a clade of akaryotes (Akarya). Eukarya and Akarya are sister-clades that diverge from the root of the tree of life. Each clade is supported by the highest posterior probability of 1.0. The phylogeny supports a scenario of independent origins and descent of eukaryotes and akaryotes. (B) Model selection tests identify, overwhelmingly, directional evolution models to be better-fitting models. (C) Alternative rootings, and accordingly alternative classifications or scenarios for the origins of the major clades of life, are much less probable and not supported.

220 regions of sequences are routinely trimmed off before phylogenetic analyses (43). Typically,  
 221 the conserved well-aligned regions correspond to protein domains with highly ordered three-  
 222 dimensional (3D) structures with specific 3D folds (Fig. 6A). Regions of sequences that are trimmed  
 223 usually show higher variability in length, are less ordered and are known to accumulate insertion  
 224 and deletion (indel) mutations at a higher frequency than in the regions that correspond to folded  
 225 domains (44). These variable, structurally disordered regions, which flank the structurally ordered  
 226 domains, link different domains in multi-domain proteins (Fig. 6A). Multi-domain architecture (MDA),  
 227 the N-to-C terminal sequence of domain arrangement, is distinct for a protein family, and differs in  
 228 closely related protein families with similar functions (Fig. 6A). The variation in MDA also relates to  
 229 alignment uncertainties.



**Figure 6.** Alignment uncertainty in closely related proteins due to domain recombination. (A) Multi-domain architecture (MDA) of the translational GTPase superfamily based on recombination of 8 modular domains. 57 distinct families with varying MDAs are known, of which 6 canonical families are shown as a schematic on the left and the corresponding 3D folds on the right. Amino acid sequences of only 2 of the 8 conserved domains can be aligned with confidence for use in phylogenetic analysis. The length of the alignment varies from 200-300 amino acids depending on the sequence diversity sampled (14,76). The EF-Tu—EF-G paralogous pair employed as pseudo-outgroups for the classical rooting of the rRNA tree is highlighted. (B) Phyletic distribution of 1,738 out of the 2,000 distinct SCOP-domains sampled from 222 species used for phylogenetic analyses in the present study. About 70 percent of the domains are widely distributed across the sampled taxonomic diversity.

230 A closer look at the 29 core-genes dataset shows that the concatenated-MSA corresponds to a  
 231 total of 27 distinct protein domains or genomic loci (Table 1). The number of loci sampled from  
 232 different species varies between 20 and 27, since not all loci are found in all species. While some  
 233 loci are absent in some species, some loci are redundant. For instance, the P-loop NTP hydrolase  
 234 domain, one of the most prevalent protein domains, is represented up to 9 times in many species  
 235 (Table 1). Many central cellular functions are driven by the conformational changes in proteins  
 236 induced by the hydrolysis of nucleoside triphosphate (NTP) catalyzed by the P-loop domain. Out of  
 237 a total of 27 distinct domains, 7 are redundant, with two or more copies represented per species.  
 238 Similarly, 9 of the 50 domains have a redundant representation in the 44 core-gene dataset (Table 1).  
 239 The observed redundancy of the genomic loci in the core-genes alignments is inconsistent with the  
 240 common (and typically untested) assumption of using single-copy genes as a proxy for orthologous  
 241 loci sampled for phylogenetic analysis.

Dataset	No. of taxa	No. of unique genes	No. of unique domains	Redundant domains		No. of times redundant in each taxon	No. of taxa in which redundant
				SCOP Unique ID	Description		
29 core-genes dataset	44	29	27	52540	P-loop containing NTP hydrolases	9	29
						8	6
						7	2
				50447	Translation proteins	3	10
						2	13
				54211	Ribosomal protein S5 domain 2-like	3	33
						2	4
50249	Nucleic acid-binding proteins	2	17				
		2	34				
64484	beta and beta-prime subunits of DNA dependent RNA-polymerase	2	37				
48 core-genes dataset	96	48	50			5	3
						4	81
						3	11
						2	1
				50249	Nucleic acid-binding proteins	3	78
						2	18
				50104	Translation proteins SH3-like domain	3	15
						2	71
				50447	Translation proteins	3	88
						2	5
64484	beta and beta-prime subunits of DNA dependent RNA-polymerase	2	83				
52540	P-loop containing NTP hydrolases	2	40				
53067	Actin-like ATPase domain	2	90				
53137	Translational machinery components	2	93				
54211	Ribosomal protein S5 domain 2-like	2	88				
56053	Ribosomal protein L6	2					

**Figure 7.** Redundant representation of protein-domains in concatenated core-genes datasets. The P-loop NTP hydrolase domain is one of the most prevalent domain. Genomic loci corresponding to P-loop hydrolase domain are represented 8-9 times in each species in the single-copy genes employed from core-genes multiple sequence alignments. Redundant loci in the core-genes datasets vary depending on the genes and species sampled for phylogenomic analyses.

242 In contrast, the protein-domain datasets are composed of unique loci (Fig. 6B). Despite the  
 243 superficial similarity of the DDNs in Fig.1 and Fig.2, they are both qualitatively and quantitatively  
 244 different codings of genome sequences. As opposed to tracing the history of 30-50 loci in the  
 245 standard core-genes datasets (Fig. 1), up to 60 fold (1738 loci) more information can be represented  
 246 when genome sequences are coded as protein-domain characters (Fig. 2). Currently 2,000 unique  
 247 domains are described by SCOP (Structural Classification of Proteins) (45). The phyletic distribution  
 248 of 1,738 domains identified in the 222 representative species sampled here is shown in a Venn  
 249 diagram (Fig. 5B).

## 250 Discussion

### 251 Improving data quality can be more effective for resolving recalcitrant branches 252 than increasing model complexity

253 In the phylogenetic literature, the concept of data quality refers to the quality or the strength of  
 254 the phylogenetic signal that can be extracted from the data. The strength of the phylogenetic  
 255 signal is proportional to the confidence with which unique state-transitions can be determined for  
 256 a given set of characters on a given tree. Ideally, historically unique character transitions that entail

257 rare evolutionary innovations are desirable, to identify patterns of uniquely shared innovations  
258 (synapomorphies) among lineages. Synapomorphies are the diagnostic features used for assessing  
259 lineage-specific inheritance of evolutionary innovations. Therefore identifying character transitions  
260 that are likely to be low probability events is a basic requirement for the accuracy of phylogenetic  
261 analysis.

262 In their pioneering studies, Woese and colleagues identified unique features of the SSU rRNA  
263 – [oligonucleotide] “signatures” – that were six nucleotides or longer, to determine evolutionary  
264 relationships (2). An underlying assumption was that the probability of occurrence of the same set  
265 of oligomer signatures by chance, in non-homologous sequences, is low in a large molecule like  
266 SSU rRNA (1500-2000 nucleotides). Oligomers shorter than six nucleotides were statistically less  
267 likely to be efficient markers of homology (46). Thus SSU rRNA was an information-rich molecule to  
268 identify homologous signatures (characters) useful for phylogenetic analysis.

269 However, as sequencing of full-length rRNAs and statistical models of nucleotide substitution  
270 became common, complex oligomer-characters were replaced by elementary nucleotide-characters;  
271 and more recently by amino acid characters. Identifying rare or historically unique substitutions in  
272 empirical datasets has proven to be difficult (47, 48), consequently the uncertainty of resolving the  
273 deeper branches of the Tree of Life using marker-gene sequences remains high. A primary reason  
274 is the prevalence of phylogenetic noise (homoplasy) in primary sequence datasets (Figs 1), due to  
275 the characteristic redundancy of nucleotide and amino acid substitutions and the resulting difficulty  
276 in distinguishing phylogenetic noise from signal (homology) (49, 50). Better-fitting (or best-fitting)  
277 models are expected to extract phylogenetic signal more efficiently and thus explain the data better,  
278 but tend to be more complex than worse-fitting models (Fig. 3 C, F). Increasingly sophisticated  
279 statistical models that have been developed over the years have only marginally improved the  
280 situation (51, 52). Although increasing model complexity can correct errors of estimation and  
281 improve the fit of the data to the tree, it is not a solution to improve phylogenetic signal, especially  
282 when not present in the source data.

283 Character recoding is found to be effective in reducing the noise/redundancy in the data, and  
284 thus uncertainties in phylogenetic reconstructions. This is a form of data simplification wherein  
285 the number of amino acid alphabets is reduced to a smaller set of alphabets that are frequently  
286 substituted for each other, usually reduced from 20 to 6. Character recoding into reduced alphabets  
287 is useful in cases where compositional heterogeneity or substitution saturation is high. However,  
288 datasets in which phylogenetic noise is inherently limited are more desirable, to minimize ambi-  
289 guities. Like amino acids, protein domains are also modular alphabets, albeit higher order and  
290 more complex alphabets of proteins. Moreover, unlike the 20 standard amino acids, there are  
291 approximately 2,000 unique protein domains identified at present according to SCOP (45). The  
292 number is expected to increase; the theoretical estimates range between 4,000 and 10,000 distinct  
293 domain modules, depending on the classification scheme (53). Coding features as binary characters  
294 is the simplest possible representation of data for describing historically unique events.

295 The idea of ‘oligonucleotide-signatures’ used for estimating a gene phylogeny has been extended,  
296 naturally, to infer a genome phylogeny (54). The signatures were defined in terms of protein-coding  
297 genes that were shared among the Archaea. However, as proteins are mosaics of domains, domains  
298 are unique genomic signatures (Fig. 6). Protein domains defined by SCOP correspond to complex  
299 ‘multi-dimensional signatures’ defined by: (i) a unique 3D fold, (ii) a distinct sequence profile, and  
300 (iii) a characteristic function. Though domain recombination is frequent, substitution of one protein  
301 domain for another has not been observed in homologous proteins (Fig. 6). For phylogenomic  
302 applications protein domains are ‘sequence signatures’ that essentially correspond to single-copy  
303 orthologous loci when coded as binary-state characters (presence/absence). These sequence  
304 signatures are consistent with unique, non-recombining genomic loci, and are identified using  
305 sophisticated statistical models — profile hidden Markov models (pHMMs) (55, 56) — that can be  
306 used routinely to annotate and curate genome sequences in automated pipelines (57, 58).

307 For these reasons, protein domains are ideal molecular phylogenetic markers for which character-

308 homology can be validated through more than one property, statistically significant (i) sequence  
309 similarity, (ii) 3D structure similarity; and (iii) function similarity. In addition, employing genomic loci  
310 for protein domains maximizes the genomic information that can be employed for phylogenetic  
311 analysis. Even though many other genomic features are known to be useful markers (59), protein  
312 domains are the most conserved as well as most widely applicable genomic characters (Fig. 6B).

### 313 **Sorting vertical evolution (signal) and horizontal evolution (noise)**

314 Single-copy genes are employed as phylogenetic markers to minimize phylogenetic noise caused  
315 by reticulate evolution, including hybridization, introgression, recombination, horizontal transfer  
316 (HT), duplication-loss (DL), or incomplete lineage sorting (ILS) of genomic loci. However, the noise  
317 observed in the DDNs based on MSA of core-genes (Fig. 1) cannot be directly related to any of the  
318 above genome-scale reticulations, since the characters are individual nucleotides or amino acids.  
319 Apart from stochastic character conflicts, the observed conflicts are better explained by convergent  
320 substitutions, given the redundancy of substitutions. Convergent substitutions caused either due  
321 to stringent selection or by chance are a well-recognized form of homoplasy in gene-sequence data  
322 (47, 50, 60), and based on recent genome-scale analyses it is now known to be rampant (61, 62).

323 The observed noise in the DDNs based on protein-domain characters (Fig. 2), however, can be  
324 related directly to genome-scale reticulation processes and homoplasies. In general, homoplasy  
325 implies evolutionary convergence, parallelism or character reversals caused by multiple processes.  
326 In contrast, homology implies only one process: inheritance of traits that evolved in the common  
327 ancestor and were passed to its descendants. Operationally, tree-based assessment of homol-  
328 ogy requires tracing the phylogenetic continuity of characters (and states), whereas homoplasy  
329 manifests as discontinuities along the tree. Since clades are diagnosed on the basis of shared  
330 innovations (synapomorphies) and defined by ancestry (63, 64), accuracy of a phylogeny depends  
331 on an accurate assessment of homology — unambiguous identification of relative synapomorphies  
332 on a best fitting tree.

333 Identifying homoplasies caused by character reversals, i.e. reversal to ancestral states requires  
334 identification of the ancestral state of the characters under study. However, implementing reversible  
335 models precludes the estimation of ancestral states, in the absence of sister groups (outgroups)  
336 or other external references. Thus, the critical distinction between shared ancestral homology  
337 (symplesiomorphy) and shared derived homology (synapomorphy) is not possible with unrooted  
338 trees derived from standard reversible models. Hence, unrooted trees (Fig. 3) are not evolutionary  
339 (phylogenetic) trees *per se*, as they are uninformative about the evolutionary polarity (34, 35, 65).  
340 Thus, identifying the root (or root-state) is crucial to (i) determine the polarity of state transitions, (ii)  
341 identify synapomorphies, and (iii) diagnose clades.

342 Moreover, because clades are associated with the emergence and inheritance of evolutionary  
343 novelties, the discovery of clades is fundamental for describing and diagnosing sister group dif-  
344 ferences, which is a primary objective of modern systematics (66). A well-recognized deficiency of  
345 phylogenetic inference based on primary sequences is the abstraction of evolutionary ‘information’  
346 (54), often into less tangible quantitative measures. For instance, ‘information’ relevant to diag-  
347 nosing clades and support for clades is abstracted to branch lengths. Branch-length estimation is,  
348 ideally, a function of the source data and the underlying model. However, in the core-genes dataset  
349 the estimated branch lengths and the resulting tree is an expression of the model rather than of the  
350 data (Fig. 3 A, B). Some pertinent questions then are: should diagnosis of clades and the features  
351 by which clades are identified be delegated to, and restricted to, substitution mutations in a small  
352 set of loci and substitution models? Are substitution mutations in 40-50 loci more informative, or  
353 the birth and death of unique genomic loci more informative?

354 Proponents of the total evidence approach recommend that all relevant information — molecu-  
355 lar, biochemical, anatomical, morphological, fossils — should be used to reconstruct evolutionary  
356 history, yet genome sequences are the most widely applicable data at present (59, 67). Accordingly,  
357 phylogenetic classification is, in practice, a classification of genomes. There is no *a priori* theoretical

358 reason that phylogenetic inference should be restricted to a small set of genomic loci corresponding  
359 to the core genes, nor is there a reason for limiting phylogenetic models to interpreting patterns  
360 of substitution mutations alone. The ease of sequencing and the practical convenience of assembling  
361 large character matrices, by themselves, are no longer compelling reasons to adhere to the  
362 traditional marker gene analysis.

363 Annotations for reference genomes of homologous protein domains identified by SCOP and  
364 other protein-classification schemes, as well as tools for identifying corresponding sequence  
365 signatures, are readily available in public databases. An added advantage is that the biochemical  
366 function and molecular phenotype of the domains are readily accessible as well, through additional  
367 resources including protein data bank (PDB) and InterPro. For complex characters such as protein  
368 domains, character homology can be determined with high confidence using sophisticated statistical  
369 models (HMMs). Homology of a protein domain implies that the *de novo* evolution of a genomic  
370 locus corresponding to that protein domain is a unique historical event. Therefore, homoplasy  
371 due to convergences and parallelisms is highly improbable (68, 69). Although a handful of cases of  
372 convergent evolution of 3D structures is known, these instances relate to relatively simple 3D folds  
373 coded for by relatively simple sequence repeats (70).

374 However, the vast majority of domains identified by SCOP correspond to polypeptides that are  
375 on average 200 residues long with unique sequence profiles (57, 68). Thus, identifying homoplasy  
376 in the protein-domain datasets depends largely on estimating reversals, which in this case will  
377 be cases of secondary gains/losses; for instance gain-loss-regain events caused by DL-HT or HT.  
378 Such secondary gains are more likely to correspond to HT events than to convergent evolution, for  
379 reasons specified above. Instances of reversals are minimal, as seen from the strong directional  
380 trends detected in the data (Fig. 5B and Fig. 6B).

### 381 **Vertical and horizontal classification**

382 For decades, biologists have been faced with a choice between so-called horizontal (Linnean) and  
383 vertical (Darwinian) classification of biodiversity (71). The similarity of both schools of systematics  
384 concerns the identification of “signatures” or sets of characteristic features that codify evolutionary  
385 relationships (54, 63, 71). But the former emphasizes the unity of contemporary groups, i.e. those  
386 at a similar evolutionary state, and therefore separates ancestors from descendants, while the latter  
387 emphasizes the unity of the ancestors and separates descendants that diverge from a common  
388 ancestry (71). Vertical classification is more consistent with the concept of lineal descent, and  
389 is the predominant paradigm for which the operational methodology and the algorithmic logic  
390 were laid out as the principles of phylogenetic systematics (63, 72). Accordingly, determining the  
391 ancestor-descendant polarity, starting from the universal common ancestor (UCA) at the root of the  
392 Tree of Life, is crucial to accurately reconstructing the path of evolutionary descent.

393 The classical rooting of the (rRNA) ToL based on the EF-Tu—EF-G paralogous pair (73, 74) is known  
394 to be error-prone and highly ambiguous, due to LBA artifacts (14, 75). Remarkably, sequences  
395 corresponding to only one of the two conserved domains common to EF-Tu and EF-G ( 200 residues  
396 in the P-loop-containing NTP hydrolase domain (Fig. 5A)) can be aligned with confidence (14).  
397 Implementing better-fitting substitution models results in two alternative rootings (R1 and R4 in  
398 Fig. 5), which relate to distinct, irreconcilable scenarios (14) similar to scenarios in Fig 4C and 4F.  
399 Moreover, the EF-Tu—EF-G paralogous pair is only 2 of 57 known paralogs of the translational  
400 GTPase protein superfamily (76). Thus the assumption that EF-Tu—EF-G duplication is a unique  
401 event, which is essential for the paralogous outgroup-rooting method, is untenable.

402 In the absence of prior knowledge of outgroups or of fossils, rooting the Tree of Life is arguably  
403 one of the most difficult phylogenetic problems. Incorrect rooting may lead to profoundly misleading  
404 conclusions about evolutionary scenarios and taxonomic affinities, and it appears to be common  
405 in phylogenetic studies (77). Perhaps worse yet seems to be the preponderance of subjective *a*  
406 *posteriori* rooting based on untested preconceptions (e.g. (78, 79)) and scenario-driven erection  
407 of taxonomic ranks (e.g. (1, 30)) (80). The conventional practice of *a posteriori* rooting, wherein an



408 unrooted tree is converted into a rooted tree by adding an *ad hoc* root, encourages a subjective  
409 interpretation of the ToL. For example, the so-called bacterial rooting of the ToL (root R3; Fig. 4) is  
410 the preferred rooting hypothesis to interpret the ToL even though that rooting is not well supported  
411 (14).

### 412 **Untangling data bias, model bias and investigator bias (prior beliefs)**

413 Phylogenies, and hence the taxonomies and evolutionary scenarios they support, are falsifiable  
414 hypotheses. Statistical hypothesis testing is now an integral part of phylogenetic inference, to  
415 quantify the empirical evidence in support of the various plausible evolutionary scenarios. However,  
416 common statistical models implemented for phylogenomic analyses are limited to modeling varia-  
417 tion in patterns of point mutations, particularly substitution mutations. These statistical models are  
418 intimately linked to basic concepts of molecular evolution, such as the universal molecular clock  
419 (3), the universal chronometer (78), paralogous outgroup rooting (81), etc., which are gene-centric  
420 concepts that were developed to study the gene, during the age of the gene. Moreover, these  
421 idealized notions originated from the analyses of relatively small single-gene datasets.

422 Conventional phylogenomics of multi-locus datasets is a direct extension of the concepts and  
423 methods developed for single-locus datasets, which rely exclusively on substitution mutations  
424 (50). In contrast, the fundamental concepts of phylogenetic theory: homology, synapomorphy,  
425 homoplasy, character polarity, etc., even if idealized, are more generally applicable. And, apparently  
426 they are better suited for unique and complex genomic characters rather than for redundant,  
427 elementary sequence characters, with regards to determining both qualitative as well as statistical  
428 consistency of the data and the underlying assumptions.

429 Phylogenetic theory that was developed to trace the evolutionary history of organismal species,  
430 as well as related methods of discrete character analysis for classifying organismal families (63, 82),  
431 was adopted, although not entirely, to determine the evolution and classification of gene families (1,  
432 3). The discovery and initial description of the Archaea was based on the comparative analysis of a  
433 single-gene (rRNA) family. However, in spite of the large number of characters that can be analyzed,  
434 neither the rRNA genes nor multi-gene concatenations of core-genes have proved to be efficient  
435 phylogenetic markers to reliably resolve the evolutionary history and phylogenetic affinities of the  
436 Archaea (83, 84).

437 Uncertainties and errors in phylogenetic inference are primarily errors in adequately distinguish-  
438 ing homologous similarities from homoplastic similarities (34, 50, 85). Homologies, synapomorphies  
439 and homoplasies are qualitative inferences, yet are inherently statistical (probabilistic). The prob-  
440 abilistic framework (maximum likelihood and Bayesian methods) has proven to be powerful for  
441 quantifying uncertainties and testing alternative hypotheses. Log odds ratios, such as LLR and  
442 LBF, are measures of how one changes belief in a hypothesis in light of new evidence (86). Accord-  
443 ingly, directional evolution models are more optimal explanations of the observed distribution of  
444 genomic-characters, and such directional trends overwhelmingly support the monophyly of the  
445 Archaea, as well as the sisterhood of the Archaea and the Bacteria, i.e. monophyly of Akarya (Fig 6).

446 Data quality is at least as important as the evolution models that are posited to explain the  
447 data. Although sophisticated statistical tests for evaluating tree robustness, and for selecting  
448 character-evolution models, are becoming a standard feature of phylogenetic software (e.g. IQ-tree,  
449 MrBayes, Phylobayes), tests for character evaluation are not common. Routines for collecting and  
450 curating data upstream of phylogenetic analyses are rather eclectic. Besides, it is an open question  
451 as to whether qualitatively different datasets (as in Fig.1 and Fig.2) can be compared effectively.  
452 Nevertheless, employing DDNs and other tools of exploratory data analysis could be useful to  
453 identify conflicts that arise due to data collection and/or curation errors (23, 24).

### 454 **Conclusions**

455 The Tree of Life is primarily a phylogenetic classification that is invaluable to organize and to  
456 describe the evolution of biodiversity, explicated through evolutionary scenarios. Phylogenies are

457 hypotheses that mostly relate to extinct ancestors, while taxonomies are hypotheses that largely  
458 relate to extant species. Extant species contain distinct combinatorial mosaics of ancestral features  
459 (plesiomorphies) and evolutionary novelties (apomorphies). It is remarkable that the uniqueness of  
460 the Archaea was identified by the comparative analyses of oligonucleotide signatures in a single  
461 gene dataset (1). However the same is not true of the phylogenetic classification of the Archaea,  
462 based on marker-genes and reversible evolution models that rely exclusively on point mutations,  
463 specifically substitution mutations, which may not be ideal phylogenetic markers (59).

464 The Three-domains of Life hypothesis (26), which was initially based on the interpretation of an  
465 unrooted rRNA tree (of life) (1), was put forward largely to emphasize the uniqueness of the Archaea,  
466 ascribed to an exclusive lineal descent. Although many lines of evidence, molecular or otherwise,  
467 support the uniqueness of the Archaea, phylogenetic analysis of genomic signatures does not  
468 support the presumed primitive state of Archaea or Bacteria, and the common belief that Archaea  
469 and Bacteria are ancestors of Eukarya (1, 11, 39, 87). Models of evolution of genomic features  
470 support a Two-domains (or rather two empires) of Life hypothesis (9), as well as the independent  
471 origins and parallel descent of eukaryote and akaryote species (10, 14, 88, 89).

## 472 **Data and methods**

### 473 **Data collection and curation**

#### 474 **Marker domains datasets**

475 Character matrices of homologous protein-domains, coded as binary-state characters were assem-  
476 bled from genome annotations of SCOP-domains available through the SUPERFAMILY HMM library  
477 and genome assignments server; v. 1.75 (<http://supfam.org/SUPERFAMILY/>) (57, 90).

478 (i) 141-species dataset was obtained from a previous study (29)

479 (ii) The 141-species dataset was updated with representatives of novel species described recently,  
480 largely with archaeal species from TACK group (30), DPANN group (5) and Asgard group including  
481 the Lokiarchaeota (20). In addition, species sampling was enhanced with representatives from  
482 the candidate phyla (unclassified) described for bacterial species and with unicellular species of  
483 eukaryotes, to a total of 222 species. The complete list of the species with their respective Taxonomy  
484 IDs is available in SI Table 1.

485 When genome annotations were unavailable from SUPERFAMILY database, curated reference  
486 proteomes were obtained from the universal protein resource (<http://www.uniprot.org/proteomes/>).  
487 SCOP-domains were annotated using the HMM library and genome annotation tools and routines  
488 recommended by the SUPERFAMILY resource.

#### 489 **Marker genes datasets**

490 Marker gene datasets from previous studies were obtained as follows, (i) 29 core-genes align-  
491 ment(17) and (ii) SSU rRNA alignment and 48 core-genes alignments (20).

### 492 **Exploratory data analysis**

493 DDNs were constructed with SplitsTree v. 4.14. Split networks were computed using the Neigh-  
494 borNet method from the observed P-distances of the taxa for both nucleotide- and amino acid-  
495 characters. Split networks of the protein-domain characters were computed from Hamming  
496 distance, which is identical to the P-distance. The networks were drawn with the equal angle  
497 algorithm.

### 498 **Phylogenetic analyses**

499 Concatenated gene tree inference: Extensive analyses of the concatenated core-genes datasets  
500 are reported in the original studies (17, 20). Analysis here was restricted to the 29 core-genes  
501 dataset due its relatively small taxon sampling (44 species) compared to the 48 core-genes dataset  
502 (96 species) since there is little difference in data quality, but the computational time/resources

503 required is significantly lesser. Moreover, the general conclusions based on these datasets are  
504 consistent despite a smaller taxon sampling, particularly of archaeal species (26 as opposed to 64  
505 in the larger sampling).

506 Best-fitting amino acid substitution models were chosen using Smart Model Selection (SMS)  
507 (91) compatible with PhyML tree inference methods (92). Trees were estimated with a rate-  
508 homogeneous LG model as well as rate-heterogeneous versions of the LG model. Site-specific rate  
509 variation was approximated using the gamma distribution with 4, 8 and 12 rate categories, LG+G4,  
510 LG+G8 and LG+G12, respectively. More complex models (SI Table 2) that account for invariable sites  
511 (LG+GX+I) and/or models that compute alignment-specific state frequencies (LG+GX+F) were also  
512 used, but the trees inferred were identical to trees estimated from LG+GX models, and therefore not  
513 reported here. Log likelihoods ratio (LLR) was calculated as the difference in the raw log likelihoods  
514 for each model.

515 Genome tree inference: The Mk model (32) is the most widely implemented model for phyloge-  
516 netic inference in the probabilistic framework (maximum likelihood (ML) and Bayesian methods)  
517 applicable to complex features coded as binary characters. However, only the reversible model is  
518 implemented in ML methods at present. Both reversible and directional evolution models as well as  
519 model selection routines implemented in MrBayes 3.2 (42, 93) were used. The Metropolis-coupled  
520 MCMC algorithm was used with two chains, sampling every 500th generation. The first half of  
521 the generations was discarded as burn-in. MCMC sampling was run until convergence, unless  
522 mentioned otherwise. Convergence was assessed through the average standard deviation of  
523 split frequencies (ASDSF, less than 0.01) for tree topology and the potential scale reduction factor  
524 (PSRF, equal to 1.00) for scalar parameters, unless mentioned otherwise. Bayes factors for model  
525 comparison were calculated using the harmonic mean estimator in MrBayes. The log Bayes factor  
526 (LBF) was calculated as the difference in the log likelihoods for each model.

527 Convergence between independent runs was generally slower for directional models compared  
528 to the reversible models. When convergence was extremely slow (requiring more than 100 million  
529 generations) topology constraints corresponding to the clusters derived in the unrooted trees (Fig.  
530 3E) were applied to improve convergence rates. In general these clusters/constraints corresponded  
531 to named taxonomic groups e.g. Fungi, Metazoa, Crenarchaeota, etc. Convergence assessment  
532 between independent runs was relaxed for three specific cases that did not converge at the time of  
533 submission: the unrooted tree with Mk-uniform-rates model (ASDSF 0.05; PSRF 1.03), rooted trees  
534 corresponding to root-R2 (ASDSF 0.5; PSRF 1.04) and root-R3 (ASDSF 0.029; PSRF 1.03). In the three  
535 cases specified, the difference in bipartitions is in the shallow parts (minor branches) of the tree.  
536 For assessing well supported major branches of the tree, ASDSF values between 0.01 and 0.05 may  
537 be adequate, as recommended by the authors (94).

## 538 Funding

539 This research received no specific grant from any funding agency in the public, commercial, or  
540 not-for-profit sectors. Work by this author was partially supported by The Swedish Research Council  
541 (to Måns Ehrenberg) and the Knut and Alice Wallenberg Foundation, RiboCORE (to Måns Ehrenberg  
542 and Dan Andersson).

## 543 Acknowledgements

544 I am grateful to Charles (Chuck) Kurland and Måns Ehrenberg for support and encouragement.  
545 I thank Chuck Kurland and Siv Andersson for the discussions in general; Chuck for the many  
546 stimulating debates and Siv for inspiring the article title, in part; Seraina Klopstein for providing  
547 the algorithms for implementing the directional model in MrBayes and for helpful suggestions and  
548 Erling Wikman for help with computing equipment.

## References

- 549 1. Woese CR, Fox GE. Phylogenetic structure of the prokaryotic domain: the primary kingdoms.  
550 Proceedings of the National Academy of Sciences. 1977;74(11):5088-90.
- 551 2. Woese CR. The Archaeal Concept and the World it Lives in: A Retrospective. Photosynthesis  
552 Research. 2004;80(1):361-72.
- 553 3. Zuckerkandl E, Pauling L. Molecules as documents of evolutionary history. Journal of theoretical  
554 biology. 1965;8(2):357-66.
- 555 4. Ragan MA, Bernard G, Chan CX. Molecular phylogenetics before sequences. RNA Biology.  
556 2014;11(3):176-85.
- 557 5. Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, et al. Insights into the  
558 phylogeny and coding potential of microbial dark matter. Nature. 2013;499(7459):431-7.
- 559 6. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, et al. Structure and  
560 function of the global ocean microbiome. Science. 2015;348(6237):1261359.
- 561 7. Wu D, Wu M, Halpern A, Rusch DB, Yooseph S, Frazier M, et al. Stalking the Fourth Domain in  
562 Metagenomic Data: Searching for, Discovering, and Interpreting Novel, Deep Branches in Marker  
563 Gene Phylogenetic Trees. PLOS ONE. 2011;6(3):e18011.
- 564 8. Boyer M, Madoui M-A, Gimenez G, La Scola B, Raoult D. Phylogenetic and Phyletic Studies  
565 of Informational Genes in Genomes Highlight Existence of a 4th Domain of Life Including Giant  
566 Viruses. PLOS ONE. 2010;5(12):e15530.
- 567 9. Mayr E. Two empires or three? Proceedings of the National Academy of Sciences of the United  
568 States of America. 1998;95(17):9720-3.
- 569 10. Harish A, Tunlid A, Kurland CG. Rooted phylogeny of the three superkingdoms. Biochimie.  
570 2013;95(8):1593-604.
- 571 11. Williams TA, Foster PG, Cox CJ, Embley TM. An archaeal origin of eukaryotes supports only  
572 two primary domains of life. Nature. 2013;504(7479):231-6.
- 573 12. Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, et al. A new view of the  
574 tree of life. Nature Microbiology. 2016;1:16048.
- 575 13. Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, et al. Recovery  
576 of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. Nature  
577 Microbiology. 2017.
- 578 14. Gouy R, Baurain D, Philippe H. Rooting the tree of life: the phylogenetic jury is still out. Phil  
579 Trans R Soc B. 2015;370(1678):20140329.
- 580 15. Lake JA. An alternative to archaeobacterial dogma. Nature. 1986;319(6055):626-.
- 581 16. Tourasse NJ, Gouy M. Accounting for evolutionary rate variation among sequence sites con-  
582 sistent changes universal phylogenies deduced from rRNA and protein-coding genes. Molecular  
583 phylogenetics and evolution. 1999;13(1):159-68.
- 584 17. Williams TA, Embley TM. Archaeal "dark matter" and the origin of eukaryotes. Genome  
585 Biology and Evolution. 2014;6(3):474-81.
- 586 18. Spang A, Saw JH, Jørgensen SL, Zaremba-Niedzwiedzka K, Martijn J, Lind AE, et al. Complex  
587 archaea that bridge the gap between prokaryotes and eukaryotes. Nature. 2015;521(7551):173-9.
- 588 19. Da Cunha V, Gaia M, Gadelle D, Nasir A, Forterre P. Lokiarchaea are close relatives of Eur-  
589 yarchaeota, not bridging the gap between prokaryotes and eukaryotes. PLOS Genetics. 2017;13(6):e1006810.
- 590 20. Zaremba-Niedzwiedzka K, Caceres EF, Saw JH, Bäckström D, Juzokaite L, Vancaester E, et al.  
591 Asgard archaea illuminate the origin of eukaryotic cellular complexity. Nature. 2017;541(7637):353-  
592 8.
- 593 21. Zwickl DJ, Hillis DM. Increased taxon sampling greatly reduces phylogenetic error. Systematic  
594 Biology. 2002;51(4):588-98.
- 595 22. Salichos L, Rokas A. Inferring ancient divergences requires genes with strong phylogenetic  
596 signals. Nature. 2013;497(7449):327-31.
- 597

- 598 23. Morrison DA. Using data-display networks for exploratory data analysis in phylogenetic  
599 studies. *Molecular Biology and Evolution*. 2009;27(5):1044-57.
- 600 24. Huson DH, Bryant D. Application of phylogenetic networks in evolutionary studies. *Mol Biol*  
601 *Evol*. 2006;23.
- 602 25. Woese CR. On the evolution of cells. *Proceedings of the National Academy of Sciences of the*  
603 *United States of America*. 2002;99(13):8742-7.
- 604 26. Woese CR, Kandler O, Wheelis ML. Towards a natural system of organisms: Proposal for the  
605 domains Archaea, Bacteria, and Eucarya. *Proceedings of the National Academy of Sciences of the*  
606 *United States of America*. 1990;87(12):4576-9.
- 607 27. Garrett RA. Molecular evolution: The uniqueness of Archaeobacteria. *Nature*. 1985;318:233-5.
- 608 28. Valentine DL. Adaptations to energy stress dictate the ecology and evolution of the Archaea.  
609 *Nature Reviews Microbiology*. 2007;5(4):316-23.
- 610 29. Harish A, Kurland CG. Akaryotes and Eukaryotes are independent descendants of a universal  
611 common ancestor. *Biochimie*. 2017;138:168-83.
- 612 30. Guy L, Ettema TJG. The archaeal TACK superphylum and the origin of eukaryotes. *Trends in*  
613 *microbiology*. 2011;19(12):580-7.
- 614 31. Le SQ, Gascuel O. An Improved General Amino Acid Replacement Matrix. *Molecular Biology*  
615 *and Evolution*. 2008;25(7):1307-20.
- 616 32. Lewis PO. A Likelihood Approach to Estimating Phylogeny from Discrete Morphological  
617 Character Data. *Systematic Biology*. 2001;50(6):913-25.
- 618 33. Wright AM, Hillis DM. Bayesian Analysis Using a Simple Likelihood Model Outperforms Parsi-  
619 mony for Estimation of Phylogeny from Discrete Morphological Data. *PLOS ONE*. 2014;9(10):e109210.
- 620 34. Morrison DA. Phylogenetic Analyses of Parasites in the New Millennium. *Advances in*  
621 *Parasitology*2006. p. 1-124.
- 622 35. Wiley EO, Lieberman BS. *Phylogenetics: theory and practice of phylogenetic systematics*:  
623 John Wiley & Sons; 2011.
- 624 36. Whittaker RH. New concepts of kingdoms of organisms. *Science*. 1969;163(3863):150-60.
- 625 37. Nasir A, Kim K, Caetano-Anolles G. Giant viruses coexisted with the cellular ancestors and  
626 represent a distinct supergroup along with superkingdoms Archaea, Bacteria and Eukarya. *BMC*  
627 *Evolutionary Biology*. 2012;12(1):156.
- 628 38. Stanier RY, Niel Cv. The concept of a bacterium. *Archives of Microbiology*. 1962;42(1):17-35.
- 629 39. Sagan L. On the origin of mitosing cells. *Journal of theoretical biology*. 1967;14(3):225-75.
- 630 40. Yang Z, Roberts D. On the use of nucleic acid sequences to infer early branchings in the tree  
631 of life. *Molecular Biology and Evolution*. 1995;12(3):451-8.
- 632 41. Huelsenbeck JP, Bollback JP, Levine AM. Inferring the root of a phylogenetic tree. *Systematic*  
633 *biology*. 2002;51(1):32-43.
- 634 42. Klopfstein S, Vilhelmsen L, Ronquist F. A Nonstationary Markov Model Detects Directional  
635 Evolution in Hymenopteran Morphology. *Systematic Biology*. 2015;64(6):1089-103.
- 636 43. Criscuolo A, Gribaldo S. BMGE (Block Mapping and Gathering with Entropy): a new soft-  
637 ware for selection of phylogenetic informative regions from multiple sequence alignments. *BMC*  
638 *Evolutionary Biology*. 2010;10(1):210.
- 639 44. Light S, Sagit R, Sachenkova O, Ekman D, Elofsson A. Protein expansion is primarily due to  
640 indels in intrinsically disordered regions. *Molecular Biology and Evolution*. 2013;30(12):2645-53.
- 641 45. Murzin AG, Brenner SE, Hubbard T, Chothia C. SCOP: A structural classification of pro-  
642 teins database for the investigation of sequences and structures. *Journal of Molecular Biology*.  
643 1995;247(4):536-40.
- 644 46. Woese CR, Fox GE, Zablen L, Uchida T, Bonen L, Pechman K, et al. Conservation of primary  
645 structure in 16S ribosomal RNA. *Nature*. 1975;254(5495):83-6.
- 646 47. Rokas A, Carroll SB. Frequent and widespread parallel evolution of protein sequences.  
647 *Molecular Biology and Evolution*. 2008;25(9):1943-53.

- 648 48. Parker J, Tsagkogeorga G, Cotton JA, Liu Y, Provero P, Stupka E, et al. Genome-wide signatures  
649 of convergent evolution in echolocating mammals. *Nature*. 2013;502(7470):228-31.
- 650 49. Rokas A, Carroll SB. Bushes in the tree of life. *PLoS Biology*. 2006;4(11):1899-904.
- 651 50. Philippe H, Roure B. Difficult phylogenetic questions: more data, maybe; better methods,  
652 certainly. *BMC Biology*. 2011;9(1):1-4.
- 653 51. Shen X-X, Hittinger CT, Rokas A. Contentious relationships in phylogenomic studies can be  
654 driven by a handful of genes. *Nature ecology & evolution*. 2017;1(5):0126.
- 655 52. Springer MS, Gatesy J. On the importance of homology in the age of phylogenomics.  
656 *Systematics and Biodiversity*. 2017:1-19.
- 657 53. Govindarajan S, Recabarren R, Goldstein RA. Estimating the total number of protein folds.  
658 *Proteins: Structure, Function and Genetics*. 1999;35(4):408-14.
- 659 54. Graham DE, Overbeek R, Olsen GJ, Woese CR. An archaeal genomic signature. *Proceedings*  
660 *of the National Academy of Sciences*. 2000;97(7):3304-8.
- 661 55. Park J, Karplus K, Barrett C, Hughey R, Haussler D, Hubbard T, et al. Sequence comparisons  
662 using multiple sequences detect three times as many remote homologues as pairwise methods.  
663 *Journal of Molecular Biology*. 1998;284(4):1201-10.
- 664 56. Eddy SR. Accelerated profile HMM searches. *PLoS Computational Biology*. 2011;7(10).
- 665 57. Gough J, Karplus K, Hughey R, Chothia C. Assignment of homology to genome sequences  
666 using a library of hidden Markov models that represent all proteins of known structure. *Journal of*  
667 *Molecular Biology*. 2001;313(4):903-19.
- 668 58. Fang H, Oates ME, Pethica RB, Greenwood JM, Sardar AJ, Rackham OJL, et al. A daily-updated  
669 tree of (sequenced) life as a reference for genome research. *Scientific Reports*. 2013;3.
- 670 59. Rokas A, Holland PWH. Rare genomic changes as a tool for phylogenetics. *Trends in Ecology*  
671 *and Evolution*. 2000;15(11):454-9.
- 672 60. Castoe TA, de Koning AJ, Pollock DD. Adaptive molecular convergence: Molecular evolution  
673 versus molecular phylogenetics. *Communicative and Integrative Biology*. 2010;3(1):12-7.
- 674 61. Liu Y, Cotton JA, Shen B, Han X, Rossiter SJ, Zhang S. Convergent sequence evolution between  
675 echolocating bats and dolphins. *Current Biology*. 2010;20(2):R53-R4.
- 676 62. Foote AD, Liu Y, Thomas GWC, Vinar T, Alfoldi J, Deng J, et al. Convergent evolution of the  
677 genomes of marine mammals. *Nat Genet*. 2015;advance online publication.
- 678 63. Hennig W. Phylogenetic systematics. *Annual review of entomology*. 1965;10(1):97-116.
- 679 64. Padian K, Lindberg DR, Polly PD. Cladistics and the fossil record: the uses of history. *Annual*  
680 *Review of Earth and Planetary Sciences*. 1994;22:63-91.
- 681 65. Lienau EK, DeSalle R. Is the microbial tree of life verificationist? *Cladistics*. 2010;26(2):195-201.
- 682 66. Sanderson MJ. Where have all the clades gone? A systematist's take in Inferring Phylogenies.  
683 *Evolution*. 2005;59(9):2056-8.
- 684 67. Wheeler Q, Assis L, Rieppel O. Phylogenetics: Heed the father of cladistics. *Nature*.  
685 2013;496(7445):295-6.
- 686 68. Pethica RB, Levitt M, Gough J. Evolutionarily consistent families in SCOP: Sequence, structure  
687 and function. *BMC Structural Biology*. 2012;12.
- 688 69. Mackin KA, Roy RA, Theobald DL. An empirical test of convergent evolution in rhodopsins.  
689 *Molecular Biology and Evolution*. 2014;31(1):85-95.
- 690 70. Mistry J, Finn RD, Eddy SR, Bateman A, Punta M. Challenges in homology search: HMMER3  
691 and convergent evolution of coiled-coil regions. *Nucleic acids research*. 2013;41(12):e121-e.
- 692 71. Simpson GG. *The Principles of Classification and a Classification of Mammals*. *Bull Amer*  
693 *Museum Nat History*. 1945;85:xvi+350.
- 694 72. Felsenstein J. *Inferring phylogenies*. Sunderland, MA: Sinauer Associates; 2004.
- 695 73. Iwabe N, Kuma K, Hasegawa M, Osawa S, Miyata T. Evolutionary relationship of archaebacte-  
696 ria, eubacteria, and eukaryotes inferred from phylogenetic trees of duplicated genes. *Proceedings*  
697 *of the National Academy of Sciences*. 1989;86(23):9355-9.



- 698 74. Baldauf SL, Palmer JD, Doolittle WF. The root of the universal tree and the origin of eukaryotes  
699 based on elongation factor phylogeny. *Proceedings of the National Academy of Sciences of the*  
700 *United States of America.* 1996;93(15):7749-54.
- 701 75. Forterre P, Philippe H. Where is the root of the universal tree of life? *BioEssays.* 1999;21(10):871-  
702 9.
- 703 76. Atkinson GC. The evolutionary and functional diversity of classical and lesser-known cyto-  
704 plasmic and organellar translational GTPases across the tree of life. *BMC Genomics.* 2015;16(1):78.
- 705 77. Graham SW, Olmstead RG, Barrett SCH. Rooting Phylogenetic Trees with Distant Outgroups:  
706 A Case Study from the Commelinoid Monocots. *Molecular Biology and Evolution.* 2002;19(10):1769-  
707 81.
- 708 78. Woese CR. Bacterial evolution. *Microbiological reviews.* 1987;51(2):221.
- 709 79. Nasir A, Caetano-Anollés G. A phylogenomic data-driven exploration of viral origins and  
710 evolution. *Science Advances.* 2015;1(8).
- 711 80. Gribaldo S, Brochier-Armanet C. Time for order in microbial systematics. *Trends in microbi-*  
712 *ology.* 2012;20(5):209-10.
- 713 81. Schwartz R, Dayhoff M. Origins of prokaryotes, eukaryotes, mitochondria, and chloroplasts.  
714 *Science.* 1978;199(4327):395-403.
- 715 82. Darwin C. *On the Origin of Species by Means of Natural Selection, or the Preservation of*  
716 *Favoured Races in the Struggle for Life.* London: John Murray; 1859.
- 717 83. Gribaldo S, Poole AM, Daubin V, Forterre P, Brochier-Armanet C. The origin of eukaryotes  
718 and their relationship with the Archaea: are we at a phylogenomic impasse? *Nat Rev Micro.*  
719 2010;8(10):743-52.
- 720 84. Gupta RS. Impact of genomics on the understanding of microbial evolution and classification:  
721 the importance of Darwin's views on classification. *FEMS microbiology reviews.* 2016;40(4):520-53.
- 722 85. Avise JC, Robinson TJ. Hemipecty: a new term in the lexicon of phylogenetics. *Systematic*  
723 *Biology.* 2008;57(3):503-7.
- 724 86. Huelsenbeck JP, Larget B, Alfaro ME. Bayesian Phylogenetic Model Selection Using Reversible  
725 Jump Markov Chain Monte Carlo. *Molecular Biology and Evolution.* 2004;21(6):1123-33.
- 726 87. Woese CR. Interpreting the universal phylogenetic tree. *Proceedings of the National Academy*  
727 *of Sciences.* 2000;97(15):8392-6.
- 728 88. Brinkmann H, Philippe H. Archaea sister group of Bacteria? Indications from tree reconstruc-  
729 tion artifacts in ancient phylogenies. *Molecular biology and evolution.* 1999;16(6):817-25.
- 730 89. Harish A, Kurland CG. Mitochondria are not captive bacteria. *Journal of Theoretical Biology.*  
731 2017;434:88-98.
- 732 90. Oates ME, Stahlhacker J, Vavoulis DV, Smithers B, Rackham OJL, Sardar AJ, et al. The SUPER-  
733 FAMILY 1.75 database in 2014: A doubling of data. *Nucleic Acids Research.* 2015;43(D1):D227-D33.
- 734 91. Lefort V, Longueville J-E, Gascuel O. SMS: Smart Model Selection in PhyML. *Molecular Biology*  
735 *and Evolution.* 2017:msx149.
- 736 92. Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms and  
737 methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0.  
738 *Systematic biology.* 2010;59(3):307-21.
- 739 93. Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, Höhna S, et al. MrBayes 3.2:  
740 Efficient Bayesian Phylogenetic Inference and Model Choice Across a Large Model Space. *Systematic*  
741 *Biology.* 2012;61(3):539-42.
- 742 94. Ronquist F, Huelsenbeck J, Teslenko M. MrBayes version 3.2 manual: tutorials and model  
743 summaries. Available with the software distribution at [mrbayessourceforgenet.com/mr32\\_manual.pdf](http://mrbayessourceforgenet.com/mr32_manual.pdf).  
744 2011.

**Supplementary Information for**

**What is an archaeon and are the Archaea really unique?**

Ajith Harish

Department of Cell and Molecular Biology, Section of Structural and Molecular Biology,  
Uppsala University, Uppsala, Sweden

Table of contents

SI Table 1. List of species analyzed in Fig. 2B, Fig. 3D,E and Fig. 5A	pages 1-3
SI Table 2. List of amino acid substitution models compared	page 4
SI Figure 1. Unrooted genome trees derived from complex models	page 5

SI Table 1. List of organisms analyzed in Fig. 2B, Fig. 3D,E as well as Fig. 5A

Tree label	Taxonomy ID	Scientific name	Taxonomic group	Superkingdom
Asgard_1538547	1538547	<i>Lokiarchaeum</i> sp. GC14_75	Asgard	Archaea
Asgard_1706443	1706443	<i>Candidatus Thorarchaeota archaeon SMTZ-45</i>	Asgard	Archaea
Asgard_1706444	1706444	<i>Candidatus Thorarchaeota archaeon SMTZ1-45</i>	Asgard	Archaea
Asgard_1706445	1706445	<i>Candidatus Thorarchaeota archaeon SMTZ1-83</i>	Asgard	Archaea
Asgard_1837170	1837170	<i>Candidatus Thorarchaeota archaeon AB_25</i>	Asgard	Archaea
Asgard_1841596	1841596	<i>Candidatus Heimdallarchaeota archaeon AB_125</i>	Asgard	Archaea
Asgard_1849166	1849166	<i>Candidatus Lokiarchaeota archaeon CR_4</i>	Asgard	Archaea
Crenarchaeota_272557	272557	<i>Aeropyrum pernix</i>	Crenarchaeota	Archaea
Crenarchaeota_368408	368408	<i>Thermofilum pendens</i>	Crenarchaeota	Archaea
Crenarchaeota_384616	384616	<i>Pyrobaculum islandicum</i>	Crenarchaeota	Archaea
Crenarchaeota_397948	397948	<i>Caldivirga maquilingensis</i>	Crenarchaeota	Archaea
Crenarchaeota_399549	399549	<i>Metallosphaera sedula</i>	Crenarchaeota	Archaea
Crenarchaeota_399550	399550	<i>Staphylothermus marinus</i>	Crenarchaeota	Archaea
Crenarchaeota_415426	415426	<i>Hyperthermus butylicus</i>	Crenarchaeota	Archaea
Crenarchaeota_419942	419942	<i>Sulfolobus islandicus</i>	Crenarchaeota	Archaea
Crenarchaeota_444157	444157	<i>Pyrobaculum neutrophilum</i>	Crenarchaeota	Archaea
Crenarchaeota_453591	453591	<i>Ignicoccus hospitalis</i>	Crenarchaeota	Archaea
Crenarchaeota_490899	490899	<i>Desulfurococcus amylolyticus</i>	Crenarchaeota	Archaea
Crenarchaeota_633148	633148	<i>Thermosphaera aggregans</i>	Crenarchaeota	Archaea
Crenarchaeota_666510	666510	<i>Acidilobus saccharovorans</i>	Crenarchaeota	Archaea
Crenarchaeota_765177	765177	<i>Desulfurococcus mucosus</i>	Crenarchaeota	Archaea
DPANN_662762	662762	<i>Candidatus Parvarchaeum acidophilum ARMAN-5</i>	DPANN	Archaea
DPANN_1294122	1294122	<i>Candidatus Nanobsidianus stetteri</i>	DPANN	Archaea
DPANN_1801881	1801881	<i>Candidatus Pacearchaeota archaeon RBG_13_36_9</i>	DPANN	Archaea
DPANN_1805293	1805293	<i>Candidatus Pacearchaeota archaeon CG1_02_30_18</i>	DPANN	Archaea
DPANN_1805297	1805297	<i>Candidatus Pacearchaeota archaeon CG1_02_35_32</i>	DPANN	Archaea
DPANN_1805298	1805298	<i>Candidatus Pacearchaeota archaeon CG1_02_39_14</i>	DPANN	Archaea
DPANN_1912863	1912863	<i>Candidatus Micrarchaeum acidiphilum ARMAN-1</i>	DPANN	Archaea
Euryarchaeota_79929	79929	<i>Methanothermobacter marburgensis</i>	Euryarchaeota	Archaea
Euryarchaeota_190192	190192	<i>Methanopyrus kandleri</i>	Euryarchaeota	Archaea
Euryarchaeota_259564	259564	<i>Methanococcoides burtonii</i>	Euryarchaeota	Archaea
Euryarchaeota_269797	269797	<i>Methanosarcina barkeri</i>	Euryarchaeota	Archaea
Euryarchaeota_272569	272569	<i>Haloarcula marismortui</i>	Euryarchaeota	Archaea
Euryarchaeota_304371	304371	<i>Methanocella paludicola</i>	Euryarchaeota	Archaea
Euryarchaeota_309800	309800	<i>Haloferax volcanii</i>	Euryarchaeota	Archaea
Euryarchaeota_323259	323259	<i>Methanospirillum hungatei JF-1</i>	Euryarchaeota	Archaea
Euryarchaeota_339860	339860	<i>Methanosphaera stadtmanae</i>	Euryarchaeota	Archaea
Euryarchaeota_348780	348780	<i>Natronomonas pharaonis</i>	Euryarchaeota	Archaea
Euryarchaeota_349307	349307	<i>Methanosaeta thermophila</i>	Euryarchaeota	Archaea
Euryarchaeota_362976	362976	<i>Haloquadratum walsbyi</i>	Euryarchaeota	Archaea
Euryarchaeota_368407	368407	<i>Methanoculleus marisnigri</i>	Euryarchaeota	Archaea
Euryarchaeota_410358	410358	<i>Methanocorpusculum labreanum</i>	Euryarchaeota	Archaea
Euryarchaeota_416348	416348	<i>Halorubrum lacusprofundi</i>	Euryarchaeota	Archaea
Euryarchaeota_419665	419665	<i>Methanococcus aeolicus</i>	Euryarchaeota	Archaea
Euryarchaeota_420247	420247	<i>Methanobrevibacter smithii</i>	Euryarchaeota	Archaea
Euryarchaeota_478009	478009	<i>Halobacterium salinarum</i>	Euryarchaeota	Archaea
Euryarchaeota_485914	485914	<i>Halomicrobium mukohataei</i>	Euryarchaeota	Archaea
Euryarchaeota_521011	521011	<i>Methanosphaerula palustris</i>	Euryarchaeota	Archaea
Euryarchaeota_543526	543526	<i>Haloterrigena turkmenica</i>	Euryarchaeota	Archaea
Euryarchaeota_547558	547558	<i>Methanohalophilus mahii</i>	Euryarchaeota	Archaea
Euryarchaeota_547559	547559	<i>Natrialba magadii</i>	Euryarchaeota	Archaea
Euryarchaeota_572546	572546	<i>Archaeoglobus profundus</i>	Euryarchaeota	Archaea
Euryarchaeota_573064	573064	<i>Methanocaldococcus fervens</i>	Euryarchaeota	Archaea
Euryarchaeota_589924	589924	<i>Ferroglobus placidus</i>	Euryarchaeota	Archaea
Euryarchaeota_644295	644295	<i>Methanohalobium evestigatum</i>	Euryarchaeota	Archaea
Euryarchaeota_795797	795797	<i>Halalkalicoccus jeotgali</i>	Euryarchaeota	Archaea
Euryarchaeota_1609968	1609968	<i>Methanobrevibacter</i> sp. YE315	Euryarchaeota	Archaea
Euryarchaeota_1641383	1641383	<i>Methanobacterium</i> sp. 42_16	Euryarchaeota	Archaea
Euryarchaeota_1860099	1860099	<i>Methanobrevibacter</i> sp. A27	Euryarchaeota	Archaea
Euryarchaeota_387957	387957	<i>Methanobrevibacter</i> sp. 87.7	Euryarchaeota	Archaea
Korarchaeota_374847	374847	<i>Korarchaeum cryptofilum</i>	Korarchaeota	Archaea
Nanoarchaeota_228908	228908	<i>Nanoarchaeum equitans</i>	Nanoarchaeota	Archaea
Thaumarchaeota_311458	311458	<i>Candidatus Caldiarchaeum subterraneum</i>	Thaumarchaeota	Archaea
Thaumarchaeota_414004	414004	<i>Cenarchaeum symbiosum</i>	Thaumarchaeota	Archaea
Thaumarchaeota_436308	436308	<i>Nitrosopumilus maritimus</i>	Thaumarchaeota	Archaea
Thaumarchaeota_886738	886738	<i>Candidatus Nitrosoarchaeum limnia SFB1</i>	Thaumarchaeota	Archaea
Thaumarchaeota_926571	926571	<i>Nitrososphaera viennensis EN76</i>	Thaumarchaeota	Archaea
Thaumarchaeota_1229908	1229908	<i>Candidatus Nitrosopumilus koreensis AR1</i>	Thaumarchaeota	Archaea
Thaumarchaeota_1229909	1229909	<i>Candidatus Nitrosopumilus sediminis</i>	Thaumarchaeota	Archaea
Thaumarchaeota_1237085	1237085	<i>Nitrososphaera gargensis</i>	Thaumarchaeota	Archaea
Thaumarchaeota_1410606	1410606	<i>Candidatus Nitrosopelagicus brevis</i>	Thaumarchaeota	Archaea
Thaumarchaeota_1459636	1459636	<i>Candidatus Nitrososphaera evergladensis SR1</i>	Thaumarchaeota	Archaea
Thaumarchaeota_1580092	1580092	<i>Candidatus Nitrosopumilus adriaticus</i>	Thaumarchaeota	Archaea

Unclassified_1579367	1579367	archaeon GW2011_AR5	Unclassified	Archaea
Unclassified_1579370	1579370	archaeon GW2011_AR10	Unclassified	Archaea
Unclassified_1579373	1579373	archaeon GW2011_AR15	Unclassified	Archaea
Unclassified_1579378	1579378	archaeon GW2011_AR20	Unclassified	Archaea
Acidobacteri_234267	234267	<i>Solibacter usitatus</i>	Acidobacteri	Bacteria
Acidobacteria_240015	240015	<i>Acidobacterium capsulatum</i>	Acidobacteria	Bacteria
Actinobacteria_469371	469371	<i>Thermobispora bispora</i>	Actinobacteria	Bacteria
Actinobacteria_469378	469378	<i>Cryptobacterium curtum</i>	Actinobacteria	Bacteria
Aquificae_123214	123214	<i>Persephonella marina</i>	Aquificae	Bacteria
Aquificae_204536	204536	<i>Sulfurihydrogenibium azorense</i>	Aquificae	Bacteria
Aquificae_608538	608538	<i>Hydrogenobacter thermophilus</i>	Aquificae	Bacteria
Bacteroides_818	818	<i>Bacteroides thetaiotaomicron</i>	Bacteroides	Bacteria
Bacteroidetes_216432	216432	<i>Croceibacter atlanticus</i>	Bacteroidetes	Bacteria
Bacteroidetes_376686	376686	<i>Flavobacterium johnsoniae</i>	Bacteroidetes	Bacteria
Bacteroidetes_521097	521097	<i>Capnocytophaga ochracea</i>	Bacteroidetes	Bacteria
Chlorobi_290318	290318	<i>Chlorobium phaeovibrioides</i>	Chlorobi	Bacteria
Chlorobi_290512	290512	<i>Prosthecochloris aestuarii</i>	Chlorobi	Bacteria
Chlorobi_517418	517418	<i>Chloroherpeton thalassium</i>	Chlorobi	Bacteria
Chloroflexi_316274	316274	<i>Herpetosiphon aurantiacus</i>	Chloroflexi	Bacteria
Chloroflexi_479434	479434	<i>Sphaerobacter thermophilus</i>	Chloroflexi	Bacteria
Chloroflexi_552811	552811	<i>Dehalogenimonas lykanthroporepellens</i>	Chloroflexi	Bacteria
Cyanobacteria_197221	197221	<i>Thermosynechococcus elongatus</i>	Cyanobacteria	Bacteria
Cyanobacteria_251221	251221	<i>Gloeobacter violaceus</i>	Cyanobacteria	Bacteria
Deinococcus_Thermus_262724	262724	<i>Thermus thermophilus</i>	Deinococcus_Thermus	Bacteria
Deinococcus_Thermus_546414	546414	<i>Deinococcus deserti</i>	Deinococcus_Thermus	Bacteria
Deinococcus_Thermus_649638	649638	<i>Truepera radiovictrix</i>	Deinococcus_Thermus	Bacteria
Firmicutes_290402	290402	<i>Clostridium beijerinckii</i>	Firmicutes	Bacteria
Firmicutes_498761	498761	<i>Helio bacterium modesticaldum</i>	Firmicutes	Bacteria
Firmicutes_515620	515620	<i>Eubacterium eligens</i>	Firmicutes	Bacteria
Fusobacteria_190304	190304	<i>Fusobacterium nucleatum subsp. nucleatum</i>	Fusobacteria	Bacteria
Fusobacteria_519441	519441	<i>Streptobacillus moniliformis</i>	Fusobacteria	Bacteria
Fusobacteria_523794	523794	<i>Leptotrichia buccalis</i>	Fusobacteria	Bacteria
Planctomycetes_243090	243090	<i>Rhodopirellula baltica</i>	Planctomycetes	Bacteria
Planctomycetes_521674	521674	<i>Planctopirus limnophila</i>	Planctomycetes	Bacteria
Planctomycetes_530564	530564	<i>Pirellula staleyi</i>	Planctomycetes	Bacteria
Proteobacteria_265072	265072	<i>Methylobacillus flagellatus</i>	Proteobacteria	Bacteria
Proteobacteria_365044	365044	<i>Polaromonas naphthalenivorans</i>	Proteobacteria	Bacteria
Proteobacteria_557598	557598	<i>Laribacter hongkongensis</i>	Proteobacteria	Bacteria
Spirochaetes_243275	243275	<i>Treponema denticola</i>	Spirochaetes	Bacteria
Spirochaetes_573413	573413	<i>Sediminispirochaeta smaragdinae</i>	Spirochaetes	Bacteria
Thermotogae_381764	381764	<i>Fervidobacterium nodosum</i>	Thermotogae	Bacteria
Thermotogae_390874	390874	<i>Thermotoga petrophila</i>	Thermotogae	Bacteria
Thermotogae_391009	391009	<i>Thermosiphon melanesiensis</i>	Thermotogae	Bacteria
Unclassified_671143	671143	<i>Candidatus Methyloirabilis oxyfera</i>	Unclassified	Bacteria
Unclassified_1635277	1635277	candidate division TA06 bacterium 34_109	Unclassified	Bacteria
Unclassified_1640508	1640508	<i>Candidatus Dadabacteria bacterium CSP1-2</i>	Unclassified	Bacteria
Unclassified_1640516	1640516	candidate division NC10 bacterium CSP1-5	Unclassified	Bacteria
Unclassified_1703775	1703775	candidate division WOR_1 bacterium DG_54_3	Unclassified	Bacteria
Unclassified_1703779	1703779	candidate division WOR_3 bacterium SM23_42	Unclassified	Bacteria
Unclassified_1703780	1703780	candidate division WOR_3 bacterium SM23_60	Unclassified	Bacteria
Unclassified_1797270	1797270	<i>Candidatus Aminicenantes bacterium RBG_13_63_10</i>	Unclassified	Bacteria
Unclassified_1797273	1797273	<i>Candidatus Aminicenantes bacterium RBG_16_63_16</i>	Unclassified	Bacteria
Unclassified_1797275	1797275	<i>Candidatus Aminicenantes bacterium RBG_19FT_COMBO_58_17</i>	Unclassified	Bacteria
Unclassified_1797291	1797291	<i>Candidatus Atribacteria bacterium RBG_19FT_COMBO_35_14</i>	Unclassified	Bacteria
Unclassified_1798559	1798559	candidate division KSB1 bacterium RBG_16_48_16	Unclassified	Bacteria
Unclassified_1801658	1801658	candidate division NC10 bacterium RIFCSPLOWO2_02_FULL_66_22	Unclassified	Bacteria
Unclassified_1802102	1802102	<i>Candidatus Rokubacteria bacterium RIFCSPHIGHO2_02_FULL_73_26</i>	Unclassified	Bacteria
Unclassified_1805370	1805370	<i>Candidatus Rokubacteria bacterium 13_2_20CM_2_70_11</i>	Unclassified	Bacteria
Unclassified_1817851	1817851	<i>Candidatus Edwardsbacteria bacterium GWF2_54_11</i>	Unclassified	Bacteria
Unclassified_1817856	1817856	<i>Candidatus Eisenbacteria bacterium RBG_16_71_46</i>	Unclassified	Bacteria
Unclassified_1817859	1817859	<i>Candidatus Firestonebacteria bacterium RIFOXYA2_FULL_40_8</i>	Unclassified	Bacteria
Unclassified_1817861	1817861	<i>Candidatus Firestonebacteria bacterium RIFOXYC2_FULL_39_67</i>	Unclassified	Bacteria
Unclassified_1817863	1817863	<i>Candidatus Fischerbacteria bacterium RBG_13_37_8</i>	Unclassified	Bacteria
Unclassified_1817867	1817867	<i>Candidatus Glassbacteria bacterium RIFCSPLOWO2_12_FULL_58_11</i>	Unclassified	Bacteria
Unclassified_1817872	1817872	<i>Candidatus Margulisbacteria bacterium GWE2_39_32</i>	Unclassified	Bacteria
Unclassified_1817873	1817873	<i>Candidatus Margulisbacteria bacterium GWF2_35_9</i>	Unclassified	Bacteria
Unclassified_1817890	1817890	<i>Candidatus Raymondbacteria bacterium RIFOXYD12_FULL_49_13</i>	Unclassified	Bacteria
Unclassified_1968529	1968529	<i>Candidatus Aminicenantes bacterium 4484_214</i>	Unclassified	Bacteria
Unclassified_1970772	1970772	candidate division KSB1 bacterium 4484_87	Unclassified	Bacteria
Verrucomicrobia_349741	349741	<i>Akkermansia muciniphila</i>	Verrucomicrobia	Bacteria
Verrucomicrobia_481448	481448	<i>Methylacidiphilum inferorum</i>	Verrucomicrobia	Bacteria
Verrucomicrobia_583355	583355	<i>Coralimargarita akajimensis</i>	Verrucomicrobia	Bacteria
Alveolata_5823	5823	<i>Plasmodium berghei</i>	Alveolata	Eukarya
Alveolata_5825	5825	<i>Plasmodium chabaudi</i>	Alveolata	Eukarya

Alveolata_5851	5851	<i>Plasmodium knowlesi</i>	Alveolata	Eukarya
Alveolata_36329	36329	<i>Plasmodium falciparum</i>	Alveolata	Eukarya
Alveolata_126793	126793	<i>Plasmodium vivax</i>	Alveolata	Eukarya
Alveolata_137071	137071	<i>Plasmodium falciparum</i>	Alveolata	Eukarya
Alveolata_432359	432359	<i>Toxoplasma gondii</i>	Alveolata	Eukarya
Alveolata_1202447	1202447	<i>Symbiodinium minutum</i>	Alveolata	Eukarya
Choanoflagellida_81824	81824	<i>Monosiga brevicollis</i>	Choanoflagellida	Eukarya
Choanoflagellida_946362	946362	<i>Salpingoeca rosetta</i>	Choanoflagellida	Eukarya
Cryptophyta_905079	905079	<i>Guillardia theta CCMP2712</i>	Cryptophyta	Eukarya
Euglenozoa_5665	5665	<i>Leishmania mexicana</i>	Euglenozoa	Eukarya
Euglenozoa_5692	5692	<i>Trypanosoma congolense</i>	Euglenozoa	Eukarya
Euglenozoa_5699	5699	<i>Trypanosoma vivax</i>	Euglenozoa	Eukarya
Euglenozoa_5702	5702	<i>Trypanosoma brucei brucei</i>	Euglenozoa	Eukarya
Euglenozoa_31285	31285	<i>Trypanosoma brucei gambiense</i>	Euglenozoa	Eukarya
Euglenozoa_185431	185431	<i>Trypanosoma brucei brucei</i>	Euglenozoa	Eukarya
Euglenozoa_347515	347515	<i>Leishmania major strain Friedlin</i>	Euglenozoa	Eukarya
Euglenozoa_353153	353153	<i>Trypanosoma cruzi</i>	Euglenozoa	Eukarya
Euglenozoa_420245	420245	<i>Leishmania braziliensis</i>	Euglenozoa	Eukarya
Euglenozoa_435258	435258	<i>Leishmania infantum</i>	Euglenozoa	Eukarya
Fungi_5270	5270	<i>Ustilago maydis</i>	Fungi	Eukarya
Fungi_36080	36080	<i>Mucor circinelloides</i>	Fungi	Eukarya
Fungi_39416	39416	<i>Tuber melanosporum</i>	Fungi	Eukarya
Fungi_192523	192523	<i>Agaricus bisporus var. bisporus</i>	Fungi	Eukarya
Fungi_240176	240176	<i>Coprinopsis cinerea</i>	Fungi	Eukarya
Fungi_246409	246409	<i>Rhizopus delemar</i>	Fungi	Eukarya
Fungi_284590	284590	<i>Kluyveromyces lactis</i>	Fungi	Eukarya
Fungi_284591	284591	<i>Yarrowia lipolytica</i>	Fungi	Eukarya
Fungi_284811	284811	<i>Ashbya gossypii</i>	Fungi	Eukarya
Fungi_285006	285006	<i>Saccharomyces cerevisiae</i>	Fungi	Eukarya
Fungi_486041	486041	<i>Laccaria bicolor</i>	Fungi	Eukarya
Fungi_510953	510953	<i>Neurospora discreta</i>	Fungi	Eukarya
Fungi_602072	602072	<i>Aspergillus carbonarius</i>	Fungi	Eukarya
Fungi_644223	644223	<i>Komagataella phaffii</i>	Fungi	Eukarya
Glaucocestophyceae_2762	2762	<i>Cyanophora paradoxa</i>	Glaucocestophyceae	Eukarya
Haptophyceae_2903	2903	<i>Emiliana huxleyi</i>	Haptophyceae	Eukarya
Ichthyosporea_192875	192875	<i>Capsaspora owczarzaki</i>	Ichthyosporea	Eukarya
Ichthyosporea_667725	667725	<i>Sphaeroforma arctica JP610</i>	Ichthyosporea	Eukarya
Metazoa_6085	6085	<i>Hydra vulgaris</i>	Metazoa	Eukarya
Metazoa_7227	7227	<i>Drosophila melanogaster</i>	Metazoa	Eukarya
Metazoa_8355	8355	<i>Xenopus laevis</i>	Metazoa	Eukarya
Metazoa_9544	9544	<i>Macaca mulatta</i>	Metazoa	Eukarya
Metazoa_9600	9600	<i>Pongo pygmaeus</i>	Metazoa	Eukarya
Metazoa_37347	37347	<i>Tupaia belangeri</i>	Metazoa	Eukarya
Metazoa_132908	132908	<i>Pteropus vampyrus</i>	Metazoa	Eukarya
Rhizaria_753081	753081	<i>Bigeloviella natans</i>	Rhizaria	Eukarya
Rhodophyta_35688	35688	<i>Porphyridium purpureum</i>	Rhodophyta	Eukarya
Rhodophyta_45157	45157	<i>Cyanidioschyzon merolae</i>	Rhodophyta	Eukarya
Rhodophyta_130081	130081	<i>Galdieria sulphuraria</i>	Rhodophyta	Eukarya
Rhodophyta_280699	280699	<i>Cyanidioschyzon merolae</i>	Rhodophyta	Eukarya
stramenopiles_186039	186039	<i>Fragilariaopsis cylindrus</i>	stramenopiles	Eukarya
stramenopiles_296543	296543	<i>Thalassiosira pseudonana CCMP1335</i>	stramenopiles	Eukarya
stramenopiles_403677	403677	<i>Phytophthora infestans</i>	stramenopiles	Eukarya
stramenopiles_556484	556484	<i>Phaeodactylum tricornutum</i>	stramenopiles	Eukarya
Viridiplantae_3055	3055	<i>Chlamydomonas reinhardtii</i>	Viridiplantae	Eukarya
Viridiplantae_3067	3067	<i>Volvox carteri</i>	Viridiplantae	Eukarya
Viridiplantae_3068	3068	<i>Volvox carteri f. nagariensis</i>	Viridiplantae	Eukarya
Viridiplantae_3077	3077	<i>Chlorella vulgaris</i>	Viridiplantae	Eukarya
Viridiplantae_3659	3659	<i>Cucumis sativus</i>	Viridiplantae	Eukarya
Viridiplantae_3702	3702	<i>Arabidopsis thaliana</i>	Viridiplantae	Eukarya
Viridiplantae_4558	4558	<i>Sorghum bicolor</i>	Viridiplantae	Eukarya
Viridiplantae_39946	39946	<i>Oryza sativa subsp. indica</i>	Viridiplantae	Eukarya
Viridiplantae_41874	41874	<i>Bathycoccus prasinos</i>	Viridiplantae	Eukarya
Viridiplantae_70448	70448	<i>Ostreococcus tauri</i>	Viridiplantae	Eukarya
Viridiplantae_296587	296587	<i>Micromonas commoda</i>	Viridiplantae	Eukarya
Viridiplantae_381124	381124	<i>Zea mays subsp. mays</i>	Viridiplantae	Eukarya
Viridiplantae_385169	385169	<i>Ostreococcus sp</i>	Viridiplantae	Eukarya
Viridiplantae_436017	436017	<i>Ostreococcus lucimarinus</i>	Viridiplantae	Eukarya
Viridiplantae_554065	554065	<i>Chlorella variabilis</i>	Viridiplantae	Eukarya
Viridiplantae_564608	564608	<i>Micromonas pusilla</i>	Viridiplantae	Eukarya
Viridiplantae_574566	574566	<i>Coccomyxa subellipsoidea</i>	Viridiplantae	Eukarya
Viridiplantae_763042	763042	<i>Asterochloris sp</i>	Viridiplantae	Eukarya

**SI Table 2. List of sequence evolution models evaluated for 29 core-genes dataset**

Base Model	Heterogeneity	K	Likelihood	LLR	AIC	BIC
LG	+G+I	87	-492883	0	985940	986553
LG	+G+I+F	106	-493145	262	986502	987250
LG	+G	86	-493206	324	986585	987191
LG	+G+F	105	-493420	538	987051	987791
RtREV	+G+I+F	106	-494212	1330	988637	989385
WAG	+G+I+F	106	-496011	3128	992234	992982
WAG	+G+I	87	-496884	4001	993942	994555
VT	+G+I+F	106	-496978	4096	994169	994916
Blosum62	+G+I	87	-497004	4121	994182	994796
VT	+G+I	87	-497243	4360	994660	995274
RtREV	+G+I	87	-497421	4538	995016	995630
Blosum62	+G+I+F	106	-497536	4654	995285	996033
CpREV	+G+I	87	-498666	5783	997506	998120
CpREV	+G+I+F	106	-498692	5809	997596	998344
JTT	+G+I	87	-501937	9055	1004049	1004662
JTT	+G+I+F	106	-501986	9103	1004183	1004931
MtZoa	+G+I+F	106	-502729	9846	1005670	1006418
DCMut	+G+I+F	106	-503422	10539	1007055	1007803
Dayhoff	+G+I+F	106	-503425	10543	1007063	1007811
MtREV	+G+I+F	106	-505701	12818	1011614	1012362
MtArt	+G+I+F	106	-507424	14541	1015060	1015807
Flu	+G+I+F	106	-510433	17550	1021078	1021826
HIVb	+G+I+F	106	-514167	21284	1028546	1029294
HIVb	+G+I	87	-514445	21563	1029065	1029679
AB	+G+I+F	106	-515188	22305	1030588	1031336
MtMam	+G+I+F	106	-520429	27546	1041069	1041817
HIVw	+G+I+F	106	-532211	39328	1064634	1065381

Base Model = Generalized empirical amino acid exchange rate (probability)

Heterogeneity = Parameter for approximating site-specific variation

+GX; where G is discrete Gamma model, X is no. of categories (default X=4)

+I; proportion of invariant sites

+F; empirical amino acid frequencies estimated from the alignment

K = Number of parameters

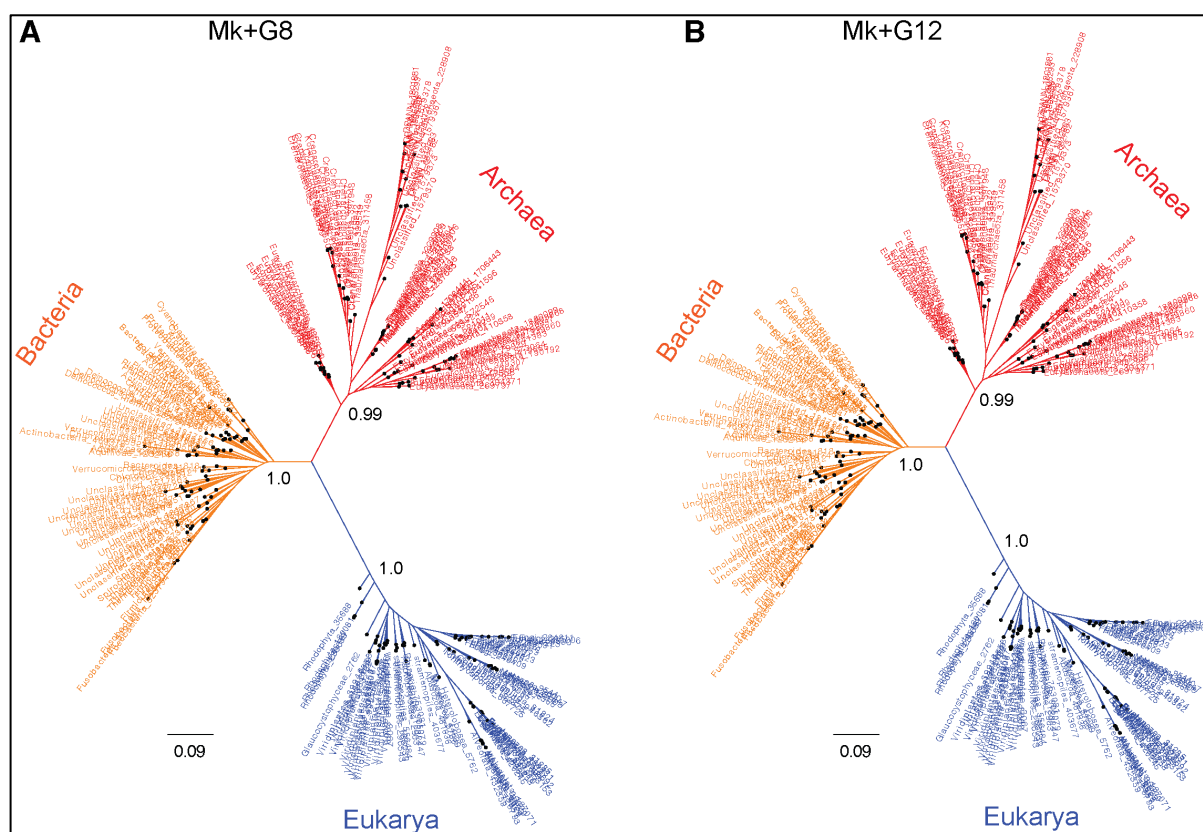
Likelihood = Raw likelihood score

LLR = Log likelihood ratio given as the difference from the best fitting model

AIC = Akaike information criterion

BIC = Bayesian information criterion





**SI Fig. 1** Unrooted genome trees derived from rate-heterogeneous versions of the Mk model. (A) Unrooted tree estimated from Mk+G8 model and (B) from Mk+G12 model. Scale bars represent expected number of changes per character. Branch support (posterior probability) is shown only for the major branches.