

1 Universal markers support a long inter-domain 2 branch between Archaea and Bacteria

3
4 Edmund R. R. Moody¹, Tara A. Mahendrarajah², Nina Dombrowski², James W. Clark¹, Celine
5 Petitjean¹, Pierre Offre², Gergely J. Szollosi^{3,4,5}, Anja Spang^{2,6}, Tom A. Williams^{1*}
6

- 7 1. School of Biological Sciences, University of Bristol, Bristol BS8 1TQ, UK.
- 8 2. NIOZ, Royal Netherlands Institute for Sea Research, Department of Marine
9 Microbiology and Biogeochemistry; AB Den Burg, The Netherlands
- 10 3. Dept. of Biological Physics, Eötvös Loránd University, 1117 Budapest, Hungary
- 11 4. MTA-ELTE “Lendület” Evolutionary Genomics Research Group, 1117 Budapest,
12 Hungary;
- 13 5. Evolutionary Systems Research Group, Centre for Ecological Research, Hungarian
14 Academy of Sciences, 8237 Tihany, Hungary
- 15 6. Department of Cell- and Molecular Biology, Science for Life Laboratory, Uppsala
16 University, SE-75123, Uppsala, Sweden

17
18 *To whom correspondence should be addressed: tom.a.williams@bristol.ac.uk

19 Abstract

20 The tree of life is generally estimated from a core set of 16-56 genes coding for proteins
21 predominantly involved in translation and other conserved informational and cellular
22 processes. These markers represent only a fraction of the genes that were likely present in
23 the last universal common ancestor (LUCA), but are useful for deep phylogenetic
24 reconstructions because their mode of inheritance appears to be mainly vertical, which
25 satisfies the assumptions of gene concatenation and supertree methods. Previous
26 phylogenetic analyses of these genes recovered a long branch between Archaea and
27 Bacteria. By contrast, a recent study made use of a greatly expanded set of 381 marker genes
28 and recovered a much shorter branch length between Archaea and Bacteria, comparable to
29 some divergences within the domains. These analyses suggest that the apparent deep split
30 between Archaea and Bacteria may be the result of accelerated evolution of ribosomal genes.
31 Here we re-evaluate the evolutionary history of the expanded marker gene set and show that
32 substitutional saturation, inter-domain gene transfer, hidden paralogy, and poor model fit
33 contribute to the inference of an artificially shortened inter-domain branch. Our results do not
34 exclude a moderately faster rate of ribosomal gene evolution during the divergence of Archaea
35 and Bacteria, but indicate that vertically-evolving marker genes across all functional categories
36 support a major genetic divergence between the two primary domains of life.
37
38
39
40

41 Introduction

42

43 Much remains unknown about the earliest period of cellular evolution and the deepest
44 divergences in the tree of life. Phylogenies encompassing both Archaea and Bacteria have
45 been inferred from a “universal core” set of 16-56 genes encoding proteins involved in
46 translation and other aspects of the genetic information processing machinery¹⁻¹⁰. These
47 genes are thought to evolve vertically and so more closely track an underlying tree of life
48 compared to other genes^{3,4,8,11}. In these analyses, the branch separating Archaea from
49 Bacteria (hereafter, the AB branch) is often the longest internal branch in the tree^{2,6,12-15}. In
50 molecular phylogenetics, branch lengths are usually measured in expected numbers of
51 substitutions per site, with a long branch corresponding to a greater degree of genetic change.
52 Long branches can therefore result from high evolutionary rates, long periods of time, or a
53 combination of the two. Molecular clock models can, in principle, disentangle the contributions
54 of these effects, but only very few fossil calibrations¹⁶ are currently available that are old
55 enough to calibrate early divergences¹⁷⁻²⁰, and as a result the ages and evolutionary rates of
56 the deepest branches of the tree remain highly uncertain.

57

58 Recently, Zhu et al.²¹ inferred a phylogeny from 381 genes distributed across Archaea and
59 Bacteria using the supertree method ASTRAL²². In addition to a large increase in the number
60 of genes compared to other universal marker sets, the functional profile of these markers
61 comprises not only proteins involved in information processing but also proteins affiliated with
62 most other functional COG categories, including metabolic processes (Table S1).
63 Subsequently, the genetic distance (branch length) between the domains²¹ was estimated
64 from a concatenation of the same marker genes, resulting in a much shorter AB branch length
65 than observed with the core universal markers^{2,6}. These analyses motivated the hypothesis²¹
66 that the apparent deep divergence of Archaea and Bacteria might be the result of an
67 accelerated evolutionary rate of genes encoding ribosomal proteins along the AB branch as
68 compared to other genes. Interestingly, the same observation was made previously using a
69 smaller set of 38 non-ribosomal marker proteins⁵, although the difference in AB branch length
70 between ribosomal and non-ribosomal markers in that analysis was reported to be
71 substantially lower (roughly two-fold, compared to roughly ten-fold for the 381 protein set of
72 Zhu et al.)^{5,21}.

73

74 A higher evolutionary rate of ribosomal genes might result from the accumulation of
75 compensatory substitutions at the interaction surfaces among the protein subunits of the
76 ribosome^{5,23}, or as a compensatory response to the addition or removal of ribosomal subunits
77 early in evolution⁵. Alternatively, differences in the inferred AB branch length might result from
78 varying rates or patterns of evolution between the traditional core genes^{2,24} and the expanded
79 set²¹. Substitutional saturation (multiple substitutions at the same site²⁵) and across-site
80 compositional heterogeneity can both impact the inference of tree topologies and branch
81 lengths²⁶⁻³⁰. These difficulties are particularly significant for ancient divergences³¹. Failure to
82 model site-specific amino acid preferences has previously been shown to lead to under-
83 estimation of the AB branch length due to a failure to detect convergent changes^{2,32}, although
84 the published analysis of the 381 marker set did not find evidence of a substantial impact of
85 these features on the tree as a whole²¹. Those analyses also identified phylogenetic
86 incongruence among the 381 markers, but did not determine the underlying cause²¹.

87
88
89
90
91
92
93
94
95
96
97
98

This recent work²¹ raises two important issues regarding the inference of the universal tree: first, that estimates of the genetic distance between Archaea and Bacteria from classic “core genes” may be unrepresentative of ancient genomes as a whole, and second, that there may be many more suitable genes to investigate early evolutionary history than generally recognized, providing an opportunity to improve the precision and accuracy of deep phylogenies. Here, we address these points by examining the evolutionary history of the 381 marker set (hereafter, the expanded marker gene set), and by evaluating the impact of orthology assignment, horizontal gene transfer (HGT), substitutional saturation and substitution model fit on inferences of the genetic divergence between Archaea and Bacteria based on gene concatenations.

99 Results and Discussion

100

101 ***Genes from the expanded marker set are not widely distributed in Archaea***

102

103 The 381-gene set was originally derived from a larger set of 400 genes used to estimate the
104 phylogenetic placement of new bacterial lineages as part of the PhyloPhlAn method³³.
105 Perhaps reflecting the focus on bacterial phylogeny in the original application, the phylogenetic
106 distribution of the 381 marker genes in the expanded set varies substantially (Table S1), with
107 many being poorly represented in Archaea. Indeed 25% of the published gene trees
108 (<https://biocore.github.io/wol/>²¹) contain less than 0.5% archaeal homologues, with 21 (5%)
109 and 69 (18%) of these trees containing no or less than 10 archaeal homologues, respectively.
110 For the remaining 75% of the gene trees, archaeal homologs comprise 0.5%-13.4% of the
111 dataset. While there are many more sequenced bacteria than archaea, 63% of the gene trees
112 possessed genes from less than half of the 669 archaeal genomes included in the analysis ,
113 whereas only 22% of the gene trees possessed fewer than half of the total number of 9906
114 sampled bacterial genomes. These distributions suggest that many of these genes are not
115 broadly present in both domains, and that some might be specific to Bacteria.

116

117 ***Conflicting evolutionary histories of individual marker genes and the inferred species*** 118 ***tree***

119

120 In the focal analysis of the 381 gene set, the tree topology was inferred using the supertree
121 method ASTRAL²², with branch lengths inferred on this fixed tree from a marker gene
122 concatenation²¹. The topology inferred from this expanded marker set²¹ is similar to published
123 trees^{6,34} and recovers Archaea and Bacteria as reciprocally monophyletic domains, albeit with
124 a shorter AB branch than in earlier analyses. However, the individual gene trees²¹ disagree
125 regarding domain monophyly: Archaea and Bacteria are recovered as reciprocally
126 monophyletic groups in only 24 of the 381 published²¹ (Table S1) maximum likelihood (ML)
127 gene trees for the expanded marker set.

128

129 Since single gene trees often fail to strongly resolve ancient relationships, we used
130 approximately-unbiased (AU) tests³⁵ to evaluate whether the failure to recover domain
131 monophyly in the published ML trees is statistically supported. For computational tractability,

132 we performed these analyses on a 1000-species subsample of the full dataset that was
133 compiled in the original study²¹. For 79 of the 381 genes, we could not perform the test
134 because the gene was not found on any of the 74 archaeal genomes present in the 1000-
135 species subsample. For the remaining 302 genes, domain monophyly was rejected ($p < 0.05$)
136 for 232 out of 302 (76.8%) genes. As a comparison, we performed the same test on several
137 smaller marker sets used previously to infer a tree of life^{2,5,36}; none of the markers in those
138 sets rejected reciprocal domain monophyly ($AU > 0.05$ for all genes, Figure 1(a)). In what
139 follows, we refer to four published marker gene sets: the Expanded set (381 genes²¹), the
140 Core set (49 genes², encoding ribosomal proteins and other conserved information-processing
141 functions), the Non-ribosomal set (38 genes, broadly distributed and explicitly selected to
142 avoid genes encoding ribosomal proteins⁵), and the Bacterial set (29 genes used in a recent
143 analysis of bacterial phylogeny³⁶).

144
145 To investigate why 232 of the marker genes rejected the reciprocal monophyly of Archaea and
146 Bacteria, we returned to the full 10,575-species dataset and annotated each sequence in each
147 marker gene family by assigning proteins to KOs, Pfams, and Interpro domains, among others
148 (Table S1, see Methods for details). We labelled the tips on each published marker gene
149 phylogeny using the corresponding KO and PFAM annotations and descriptions (See
150 Methods) and manually inspected the tree topologies for reciprocal domain monophyly of
151 Archaea and Bacteria and the presence of paralogues (Table S1). This revealed that the major
152 cause of domain polyphyly observed in gene trees was inter-domain gene transfer (in 357 out
153 of 381 gene trees (93.7%)) and mixing of sequences from distinct paralogous families (in 246
154 out of 381 gene trees (64.6%)). For instance, marker genes encoding ABC-type transporters
155 (p0131, p0151, p0159, p0174, p0181, p0287, p0306, p0364), tRNA synthetases (i.e. p0000,
156 p0011, p0020, p0091, p0094, p0202), aminotransferases and dehydratases (i.e. p0073/4-
157 aminobutyrate aminotransferase; p0093/3-isopropylmalate dehydratase) often comprised a
158 mixture of paralogues. For example, the phylogenetic tree comprising spermidine/putrescine
159 import ATP-binding protein PotA (p0131) contains several paralogues families, such as
160 different sugar or ion ATP-binding proteins. Further, dTDP-glucose 4,6-dehydratase (p0134)
161 is representative of a gene tree that includes several paralogous families with different
162 substrate specificities that include dehydrorhamnose and glucuronate.

163
164 Together, these analyses indicate that the evolutionary histories of the individual markers of
165 the expanded set differ from each other and from the species tree. Zhu et al. acknowledged²¹
166 the varying levels of congruence between the marker phylogenies and the species tree, but
167 did not investigate the underlying causes. Our analyses establish the basis for these
168 disagreements in terms of gene transfers and the mixing of orthologues and paralogues within
169 and between domains. In principle, concatenation is based on the assumption that all of the
170 genes in the supermatrix evolve on the same underlying tree; genes with different gene tree
171 topologies should not be concatenated because the topological differences among sites are
172 not modelled, and so the impact on inferred branch lengths is difficult to predict. In practice, it
173 is often difficult to be certain that all of the markers in a concatenate share the same gene tree
174 topology, and the analysis proceeds on the hypothesis that a small proportion of discordant
175 genes are not expected to seriously impact the inferred tree. However, the concatenated tree
176 inferred from the expanded marker set differs from previous trees in that the genetic distance
177 between Bacteria and Archaea is greatly reduced, such that the AB branch length appears
178 comparable to the distance between bacterial phyla²¹. An accurate estimate of the AB branch
179 length is important because it has a major bearing on unanswered questions regarding the

180 root of the universal tree³¹ and the deepest divisions among cellular life. We therefore
181 evaluated the impact of the conflicting gene histories within the expanded marker set on
182 inferred AB branch length.

183

184 ***The inferred branch length between Archaea and Bacteria is artifactually shortened by***
185 ***inter-domain gene transfer and hidden paralogy***

186

187 To investigate the impact of gene transfers and mixed paralogy on the AB branch length
188 inferred by gene concatenations²¹, we compared branch lengths estimated from markers that
189 rejected ($AU < 0.05$) or did not reject ($AU > 0.05$) the reciprocal monophyly of Bacteria and
190 Archaea in the 381 marker set (Figure 1(a)). To estimate AB branch lengths for genes in which
191 the domains were not monophyletic in the ML tree, we first performed a constrained ML search
192 to find the best gene tree that was consistent with domain monophyly for each family under
193 the LG+G4+F model in IQ-TREE 2³⁷. While it may seem strained to estimate the length of a
194 branch that does not appear in the ML tree, we reasoned that this approach would provide
195 insight into the contribution of these genes to the AB branch length in the concatenation, in
196 which they conflict with the overall topology. AB branch lengths were significantly ($P =$
197 2.159×10^{-12} , Wilcoxon rank sum test) shorter for markers that rejected domain monophyly
198 (Figure 1(a); < 0.05 : mean AB branch length in expected substitutions/site 0.0130, > 0.05 : mean
199 AB branch length 0.559). This result suggests that inter-domain gene transfers reduce the AB
200 branch length when included in a concatenation. This behaviour might result from marker gene
201 transfers reducing the number of fixed differences between the domains, so that the AB branch
202 length in a tree in which Archaea and Bacteria are constrained to be reciprocally monophyletic
203 will tend towards 0 as the number of transfers increases. Consistent with this hypothesis, we
204 observed that ΔLL , the difference in log likelihood between the ML gene tree and the
205 constrained ML tree consistent with domain monophyly (a simple proxy for marker gene
206 verticality) correlates negatively with AB branch length (Figure 1(b)). Furthermore, AB branch
207 length decreased as increasing numbers of low-verticality markers were added to the
208 concatenate (Figure 1(c)). Taken together, these results indicate that the inclusion of genes
209 that do not support the reciprocal monophyly of Archaea and Bacteria in the universal
210 concatenate reduces the estimated AB branch length by homogenizing the genetic diversity
211 of the two domains.

212

213 ***The inferred Archaea-Bacteria branch is artifactually shortened by unaccounted-for***
214 ***substitutional saturation***

215

216 The longer AB branch length observed in concatenations of traditional core marker sets^{2,6,24}
217 compared to the expanded marker set²¹ has been interpreted as evidence for ribosomal gene-
218 specific accelerated evolution between domains^{5,21}, because ribosomal proteins and proteins
219 that physically interact with the ribosome make up a large proportion (47.1%) of the core set².
220 An alternative hypothesis is that faster rates of evolution in the genes of the expanded set
221 following the divergence of Archaea and Bacteria has resulted in substitutional saturation,
222 overwriting some of the early changes and shortening the AB branch. To distinguish between
223 these hypotheses, we compared AB branch lengths estimated from fast- and slow-evolving
224 sites in the expanded concatenate, for the 1000 taxa subset²¹. We estimated site-specific rates
225 of evolution using IQ-TREE 2³⁷, and constructed two new concatenates comprising the 25%
226 fastest- and 25% slowest-evolving sites. As expected, the total tree length inferred from the
227 slowest sites was shorter (105 substitutions/site compared to 1216 substitutions/site in the

228 25% fastest site concatenation). However, the relative length of the AB branch (that is, AB
229 branch length as a proportion of total tree length) was 8.8-fold higher in the concatenation
230 inferred from the 25% slowest sites compared to the fastest sites (Figure 2(a)). Considering
231 only the top 5% of the expanded marker set in terms of Δ LL score, the relative difference
232 increases to 11.1-fold. This analysis suggests that substitutional saturation has overwritten
233 many of the changes that originally occurred during the divergence of Archaea and Bacteria
234 at fast-evolving sites in the expanded gene set. Total tree lengths estimated for the expanded
235 marker genes did not significantly differ ($P = 0.5716$) between genes that did or did not reject
236 domain monophyly.

237

238 An earlier study⁵ also recovered a longer AB branch from ribosomal compared to non-
239 ribosomal genes, and none of the markers used in that study rejected domain monophyly
240 (Figure 1, “Non-ribosomal”). While the non-ribosomal markers from that study⁵ have similar
241 total tree lengths to those in the core gene set (Figure 3(a)), they have ~2.4-fold shorter AB
242 branches, Figure 3(b)). These results are consistent with the hypothesis that ribosomal
243 proteins may have longer AB branches than other marker genes.

244

245 ***Failure to model across-site amino acid preferences reduces estimates of the AB*** 246 ***branch length***

247

248 Amino acid preferences vary across the sites of a sequence alignment, due to variation in the
249 underlying functional constraints^{26,29,30}. The consequence is that, at many alignment sites, only
250 a subset of the twenty possible amino acids are tolerated by selection. Standard substitution
251 models, such as LG+G+F, are site-homogeneous, and approximate the composition of all
252 sites using the average composition across the entire alignment. Such models underestimate
253 the rate of evolution at highly constrained sites because they cannot account for the high
254 number of multiple substitutions that occur at such sites. The effect is that site-homogeneous
255 models underestimate branch lengths when fit to site-heterogeneous data. The AB branch has
256 previously been shown to be particularly susceptible to this effect². Zhu et al.²¹ investigated
257 the impact of site heterogeneity on tree topology and branch lengths using the posterior mean
258 site frequency (PMSF) model³⁸, and found that inferences under site-homogeneous and site-
259 heterogeneous models were similar over the dataset as a whole²¹. To evaluate the effect of
260 modelling site-heterogeneity on the AB branch length in particular, we inferred phylogenies
261 from the concatenate constructed from the top 5% of genes scored by Δ LL under a series of
262 models that account for site-heterogeneity using increasingly rich mixtures of site-specific
263 composition profiles (C10-C60)³⁰. Models that account for site heterogeneity support a branch
264 that is 1.6-2.1-fold longer than using a site-homogeneous model; the best-fitting model among
265 those we tested was LG+C60+G+F according to the BIC score, and this model inferred an AB
266 branch length of 2.4; that is, twice as long as that inferred under the site-homogeneous model.
267 Therefore, and consistent with previous work^{2,32}, substitution model fit has a significant effect
268 on the inferred genetic proximity of Archaea and Bacteria (Figure 4). We also find that trimming
269 poorly-aligned sites greatly increases the inferred length of the AB branch, likely because
270 trimmed sites tend to be fast-evolving.

271

272

273

274

275 ***The age of the last universal common ancestor (LUCA) inferred from the expanded***
276 ***gene set is driven by variation in the AB branch length***

277

278 Zhu et al.²¹ argued that the expanded marker set is useful for deep phylogeny because
279 estimates of the age of LUCA obtained by fitting molecular clocks to their dataset are in
280 agreement with the geological record: a root (LUCA) age of 3.6-4.2 Gyr was inferred from the
281 entire 381-gene dataset, consistent with the earliest fossil evidence for life, whereas estimates
282 from ribosomal markers alone supported a root age of 7 Gya. This age might be considered
283 unrealistic because it is much older than the age of the Earth (with the moon-forming impact
284 occurring ~4.51Gya^{42,43}). In the original analyses, the age of LUCA was estimated using a
285 maximum likelihood approach, as well as a Bayesian molecular clock with a strict clock
286 (assuming a constant evolutionary rate) or a relaxed clock with a single calibration. A strict
287 clock model does not permit changes in evolutionary rate through time or across branches,
288 and so a longer AB branch will lead to an older inferred LUCA age. Likewise, a relaxed clock
289 model with a single calibration may fail to differentiate molecular distances and geological
290 time. Given that the short AB branch in the expanded gene set results, in part, from
291 phylogenetic incongruence among markers, we evaluated the age of LUCA inferred from the
292 subset of the expanded gene set least affected by these issues. To do so, we analysed the
293 top 5% of gene families according to their Δ LL score (including only 1 ribosomal protein) under
294 the same clock model parameters as the original dataset. This analysis resulted in a
295 significantly more ancient age estimate for LUCA (5.5-6.5 Ga), suggesting that, under these
296 conditions, the inferred age for LUCA is driven by variation in the AB branch length, and is not
297 in itself a reliable indicator of marker quality. In principle, more reliable estimates of LUCA's
298 age might be obtained by using more calibrations. However, unambiguous calibrations remain
299 elusive, particularly for the root and other deep branches of the tree. Despite advances in
300 molecular clock methodology, such calibrations represent the only way to reliably capture the
301 relationship between genetic distance and divergence time.

302 **Conclusion**

303

304 The apparent genetic proximity of Archaea and Bacteria inferred from a concatenation of the
305 expanded gene set of Zhu et al.²¹ results from substitutional saturation and inter-domain HGT.
306 Saturation obscures evidence of substitutions along the AB branch, while undetected HGT
307 acts to artifactually homogenize the genetic diversity of the domains, and leads to a reduction
308 in the inferred AB branch when branch lengths are estimated on a fixed topology. Treatments
309 of the data that improve estimation of the AB branch length - the use of better-fitting
310 substitution models, or use of the markers least affected by inter-domain HGT - results in the
311 inference of a ~15-fold longer AB branch that is similar to estimates from previously published
312 datasets (Figure 6). These results emphasise the importance of the fit of the substitution model
313 to the data, and the need to evaluate whether the data meet the assumptions of the methods
314 used - in particular, the assumption that all genes evolve on the same underlying tree topology,
315 which underlies analyses of gene concatenations. The violation of that assumption in branch
316 length estimates from concatenated alignments of the expanded marker genes resulted in the
317 recovery of an artifactually short AB branch.

318

319 A distinct question is whether the long AB branch inferred from congruent marker gene sets
320 might be affected by accelerated ribosomal gene evolution. Comparison of AB branch length
321 for sets of congruent “core” and non-ribosomal marker sets (Figure 3) are consistent with a
322 higher relative rate of evolution along the AB branch for ribosomal genes, although the effect
323 is small by comparison with the impact of phylogenetic congruence, site trimming and
324 substitution model fit (Figure 6). In sum, all of our analyses support the conclusion that the
325 split separating Archaea from Bacteria is by far the longest internal branch on the tree of the
326 primary domains of life.

327
328 In the future, it may be useful to further evaluate the suitability of non-ribosomal markers for
329 concatenated gene phylogenies. Furthermore, methods that explicitly model HGTs^{44–46}, while
330 also calculating branch lengths in a manner that accounts for differences among the
331 underlying gene trees, and that incorporate alternative means of dating information, such as
332 gene transfers^{17,47–49}, may help to better constrain the divergence times of the deepest cellular
333 lineages.

334 Methods

335 **Data**

336 We downloaded the individual alignments from ²¹
337 (<https://github.com/biocore/wol/tree/master/data/>), along with the genome metadata and the
338 individual newick files. We checked each published tree for domain monophyly, and also
339 performed approximately unbiased (AU)³⁵ tests to assess support for domain monophyly on
340 the underlying sequence alignments using IQ-TREE 2³⁷. The phylogenetic analyses were
341 carried out using the ‘reduced’ subset of 1000 taxa outlined by the authors²¹, for computational
342 tractability. These markers were also trimmed according to the protocol in the original paper²¹,
343 i.e sites with >90% gaps were removed, followed by removal of sequences with >66% gaps.

344 We also downloaded the Williams et al.² (“core”), Petitjean et al.⁵ (“non-ribosomal”) and
345 Coleman et al.³⁶ (“bacterial”) datasets from their original publications.

346

347 **Annotations**

348 Proteins used for phylogenetic analyses by Zhu *et al.*²¹, were annotated to investigate the
349 selection of sequences comprising each of the marker gene families. To this end, we
350 downloaded the protein sequences provided by the authors from the following repository:
351 <https://github.com/biocore/wol/tree/master/data/alignments/genes>. To obtain reliable
352 annotations, we analysed all sequences per gene family using several published databases,
353 including the arCOGs (version from 2014)⁵⁰, KOs from the KEGG Automatic Annotation Server
354 (KAAS; downloaded April 2019)⁵¹, the Pfam database (Release 31.0)⁵², the TIGRFAM
355 database (Release 15.0)⁵³, the Carbohydrate-Active enZymes (CAZy) database (downloaded
356 from dbCAN2 in September 2019)⁵⁴, the MEROPs database (Release 12.0)^{55,56}, the
357 hydrogenase database (HydDB; downloaded in November 2018)⁵⁷, the NCBI- non-redundant
358 (nr) database (downloaded in November 2018), and the NCBI COGs database (version from
359 2020). Additionally, all proteins were scanned for protein domains using InterProScan (v5.29-
360 68.0; settings: --iprlookup --goterms)⁵⁸.

361

362 Individual database searches were conducted as follows: arCOGs were assigned using PSI-
363 BLAST v2.7.1+ (settings: -evaluate 1e-4 -show_gis -outfmt 6 -max_target_seqs 1000 -dbsize
364 100000000 -comp_based_stats F -seg no)⁵⁹. KOs (settings: -E 1e-5), PFAMs (settings: -E 1e-
365 10), TIGRFAMs (settings: -E 1e-20) and CAZymes (settings: -E 1e-20) were identified in all
366 archaeal genomes using hmmsearch v3.1b2⁶⁰. The MEROPs and HydDB databases were
367 searched using BLASTp v2.7.1 (settings: -outfmt 6, -evaluate 1e-20). Protein sequences were
368 searched against the NCBI_nr database using DIAMOND v0.9.22.123 (settings: -more-
369 sensitive -e-value 1e-5 -seq 100 -no-self-hits -taxonmap prot.accession2taxid.gz)⁶¹. For all
370 database searches the best hit for each protein was selected based on the highest e-value
371 and bitscore and all results are summarized in the Data Supplement,
372 Annotation_Tables/0_Annotation_tables_full/All_Zhu_marker_annotations_16-12-
373 2020.tsv.zip. For InterProScan we report multiple hits corresponding to the individual domains
374 of a protein using a custom script (parse_IPRdomains_vs2_GO_2.py).

375
376 Assigned sequence annotations were summarized and all distinct KOs and Pfams were
377 collected and counted for each marker gene. KOs and Pfams with their corresponding
378 descriptions were mapped to the marker gene file downloaded from the repository:
379 <https://github.com/biocore/wol/blob/master/data/markers/metadata.xlsx> and used in
380 summarization of the 381 marker gene protein trees (Table S1).

381
382 For manual inspection of single marker gene trees, KO and Pfam annotations were mapped
383 to the tips of the published marker protein trees, downloaded from the repository:
384 <https://github.com/biocore/wol/tree/master/data/trees/genes>. Briefly, the Genome ID, Pfam,
385 Pfam description, KO, KO description, and NCBI Taxonomy string were collected from each
386 marker gene annotation table and were used to generate mapping files unique to each marker
387 gene phylogeny, which links the Genome ID to the annotation information
388 (GenomeID|Domain|Pfam|Pfam Description|KO|KO Description). An in-house Perl script
389 `replace_tree_names.pl`
390 (https://github.com/ndombrowski/Phylogeny_tutorial/tree/main/Input_files/5_required_Scripts
391) was used to append the summarized protein annotations to the corresponding tips in each
392 marker gene tree. Annotated marker gene phylogenies were manually inspected using the
393 following criteria including: 1) retention of reciprocal domain monophyly (Archaea and
394 Bacteria) and 2) for the presence or absence of potential paralogous families. Paralogous
395 groups and misannotated families present in the gene trees were highlighted and violations of
396 search criteria were recorded in Table S1.

397 **Phylogenetic analyses**

398 *Constraint analysis*

399 We performed a maximum likelihood free topology search using IQ-TREE 2³⁷ under the
400 LG+G+F model, with 10,000 bootstrap replicates. We also performed a constrained analysis
401 with the same model, in order to find the maximum likelihood tree in which Archaea and
402 Bacteria were reciprocally monophyletic. We then compared both trees using the
403 approximately unbiased (AU)³⁵ test in IQ-TREE 2³⁷ with 10,000 RELL³⁵ bootstrap replicates.
404 To evaluate the relationship between marker gene verticality and AB branch length, we
405 calculated the difference in log-likelihood between the constrained and unconstrained trees in
406 order to rank the genes from the expanded marker set, made concatenates comprised of the

407 top 20-100 (intervals of 5) of these marker genes, and inferred the tree length under
408 LG+C10+G+F with 1000 bootstrap replicates.

409 *Site and gene evolutionary rates*

410 We inferred rates using the --rate option in IQ-TREE 2³⁷ for both the 381 marker concatenation
411 from Zhu²¹ and the top 5% of marker genes based on the results of difference in log-likelihood
412 between the constrained tree and free-tree search in the constraint analysis (above). We built
413 concatenates for the 25% slowest and 25% fastest sites, and inferred branch lengths from
414 each of these concatenates using the tree inferred from the complete dataset as a fixed
415 topology.

416 *Individual markers*

417 We inferred rates and trees (1000 bootstrap replicates) for each individual marker from
418 Petitjean *et al.*⁵, Williams *et al.*², Zhu *et al.*²¹ under the LG+G+F and LG+C20+G+F models
419 using IQ-TREE 2³⁷. Mean rates per sequence were then calculated using a python script (see
420 Data Supplement).

421 *Model complexity*

422 Model complexity tests were undertaken using the top 5% concatenate described above, with
423 the alignment being trimmed with BMGE 1.12⁴¹ with default settings (BLOSUM62, entropy 0.5)
424 for all of the analyses except the 'untrimmed' LG+G+F run, other models on the trimmed
425 alignment were LG+G+F, LG+R+F and LG+C10,20,30,40,50,60+G+F, with 1000 bootstrap
426 replicates. Model fitting was done using ModelFinder⁶² in IQ-TREE 2³⁷.

427 *Molecular clock analyses*

428 Molecular clock analyses were devised to test the effect of genetic distance on the inferred
429 age of LUCA. Following the approach of Zhu *et al.*²¹, we subsampled the alignment to 100
430 species. Five alternative alignments were analysed, representing conserved sites across the
431 entire alignment, randomly selected sites across the entire alignment, only ribosomal marker
432 genes, the top 5% of marker genes according to Δ LL and the top 5% of marker genes further
433 trimmed under default settings in BMGE 1.12⁴¹. Divergence time analyses were performed in
434 MCMCTree⁶³ under a strict clock model. We used the normal approximation approach, with
435 branch lengths estimated in codeml under the LG+G4 model. In each case, a fixed tree
436 topology was used alongside a single calibration on the Cyanobacteria-Melainabacteria split.
437 The calibration was modelled as a uniform prior distribution between 2.5 and 2.6 Ga, with a
438 2.5% probability that either bound could be exceeded. For each alignment, four independent
439 MCMC chains were run for 2,000,000 generations to achieve convergence.

440 *Plotting*

441 Statistical analyses were performed using R 3.6.3⁶⁴, and data were plotted with ggplot2⁶⁵.

442

443

444

445 **Figure legends**

446

447 **Figure 1: Expanded set genes in which Archaea and Bacteria are not monophyletic**

448 **support a shorter AB branch.** (A) Expanded set genes that reject domain monophyly (AU P

449 < 0.05) support significantly shorter AB branch lengths when constrained to follow a domain

450 monophyletic tree ($P = 2.159 \times 10^{-12}$, Wilcoxon rank-sum test). None of the marker genes from

451 several other published analyses reject domain monophyly (AU $p > 0.05$ for all genes tested).

452 (B) Marker gene verticality (ΔLL , see below) for the expanded gene set normalized by

453 alignment length correlates negatively with the length of the AB branch between Archaea and

454 Bacteria ($R^2=0.03998$, $P = 0.0004731$). (C) Concatenations of 20-100 markers of the

455 expanded set markers ranked by marker gene verticality (ΔLL) show the same trend, with a

456 reduction in AB branch length as markers with a greater ΔLL are added to the concatenate.

457 ΔLL is the difference between the log likelihood of the ML gene family tree under a free

458 topology search and the log likelihood of the best tree constrained to obey domain monophyly.

459 The trendline is estimated using a LOESS regression.

460

461 **Figure 2: Vertically-evolving genes and slow-evolving sites support a longer relative**

462 **AB branch length.** We estimated site-specific evolutionary rates for all marker genes in the

463 expanded dataset (A-B), as well as for the 20 genes with the smallest ΔLL (top 5%) in that

464 dataset (C-D). Concatenations based on the 25% slowest sites (A,C) and on the top 5%

465 vertical genes (C,D) support a longer AB branch. This suggests that the inference of a short

466 AB branch is impacted by both substitutional saturation and unmodelled inter-domain transfer

467 of marker genes. Phylogenies were inferred under the LG+G4+F model in IQ-TREE ²³⁷.

468 Branch lengths are the expected number of substitutions per site, as indicated by the scale

469 bars. Alignment lengths in amino acids: **A:** 36797, **B:** 67274, **C:** 2736, **D:** 3884.

470

471 **Figure 3. Evidence for a faster rate of inter-domain ribosomal protein evolution. (A)**

472 Tree lengths (total number of substitutions/site for each gene family) for the core and non-ribosomal

473 gene sets under the LG+G+F and LG+C20+G+F models. Modelling site heterogeneity

474 increases tree lengths (the number of inferred substitutions) due to improved modelling of the

475 site-specific features of the evolutionary process. Tree lengths are not significantly different

476 between the core and non-ribosomal marker sets $P = 0.4821/0.1651$ (for LG+G+F and

477 LG+C60+G+F respectively) (B) Non-ribosomal genes have moderately shorter AB branch

478 lengths than core genes, consistent with a moderately faster rate of ribosomal gene evolution

479 on the inter-domain branch. The difference is significant under both LG+G+F ($P = 8.78 * 10^{-9}$)

480 ⁹), and the better-fitting LG+C20+G+F model ($P = 2.237 * 10^{-7}$).

481 **Figure 4. The effect of modelling site heterogeneity on AB branch length.** Increasing the
482 number of protein mixture profiles, as well as trimming is associated with a change in AB
483 branch length on the expanded marker set²¹. All analyses used LG exchangeabilities, four rate
484 categories (Gamma-distributed or freely estimated), and included a general composition
485 vector containing the empirical amino acid frequencies (+F). Modelling of site heterogeneity
486 with the C10-C60 models increases the inferred AB branch length ~2-fold. Trimming poorly-
487 aligned sites slightly increases the AB branch estimation whereas relaxing the gamma rate
488 categories slightly decreases estimation of AB branch length. LG (LG substitution matrix), G
489 (four gamma rate categories), F (empirical site frequencies estimated from the data), C10-60
490 (number of protein mixture profiles used³⁰) R (four free rate categories which relax the
491 assumption of a gamma distribution for rates^{39,40}, BMGE (trimming using Block Mapping and
492 Gathering with Entropy⁴¹). The trendline is estimated using a LOESS regression.

493
494 **Figure 5. Divergence time estimation of the Archaea-Bacteria split.** Violins represent
495 posterior age estimates from Bayesian molecular clock analyses of 1) Conserved sites as
496 estimated previously²¹; 2) Random sites²¹ 3) Ribosomal genes²¹ 4) The top 5% of marker gene
497 families according to their Δ LL score (including only 1 ribosomal protein) and 5) The same top
498 5% of marker genes trimmed using BMGE⁴¹ to remove highly variable sites. In each case, a
499 strict molecular clock was applied, with the age of the Cyanobacteria-Melainabacteria split
500 constrained between 2.5 and 2.6 Ga.

501
502 **Figure 6. The impact of marker gene choice, phylogenetic congruence, alignment**
503 **trimming, and substitution model fit on estimates of the Archaea-Bacteria branch**
504 **length.** Analysis using a site-homogeneous model (LG+G+F) on the complete 381-gene
505 expanded set results in an AB branch substantially shorter than previous estimates. Removing
506 the genes most seriously affected by inter-domain gene transfer, trimming poorly-aligned sites
507 using BMGE⁴¹, and using the best-fitting site-heterogeneous model available (LG+C60+G+F)
508 substantially increase the estimated AB length, such that it is comparable with published
509 estimates from non-ribosomal and core gene sets. Branch lengths measured in expected
510 number of substitutions/site.

511

512 **Data and code availability**

513 All of the data, including sequence alignments, trees, annotation files, and scripts associated
514 with this manuscript have been deposited in the FigShare repository at DOI:
515 10.6084/m9.figshare.13395470.

516

517 **Acknowledgements**

518 ERRM was supported by a Royal Society Enhancement Award (RGF\EA\180199) to TAW.
519 CP was supported by NERC grant NE/P00251X/1 to TAW. TAW was supported by a Royal
520 Society University Research Fellowship (URF\R\201024). GJSz received funding from the
521 European Research Council under the European Union's Horizon 2020 research and
522 innovation program under Grant Agreement 714774 and Grant GINOP-2.3.2.-15-2016-
523 00057. AS was supported by the Swedish Research Council (VR starting grant 2016-03559

524 to AS) and the NWO-I foundation of the Netherlands Organisation for Scientific Research
525 (WISE fellowship to AS).
526
527
528
529

530 References

531

- 532 1. Harris, J. K., Kelley, S. T., Spiegelman, G. B. & Pace, N. R. The genetic core of the
533 universal ancestor. *Genome Res.* **13**, 407–412 (2003).
- 534 2. Williams, T. A., Cox, C. J., Foster, P. G., Szöllősi, G. J. & Embley, T. M. Phylogenomics
535 provides robust support for a two-domains tree of life. *Nat Ecol Evol* **4**, 138–147 (2020).
- 536 3. Ciccarelli, F. D. *et al.* Toward automatic reconstruction of a highly resolved tree of life.
537 *Science* **311**, 1283–1287 (2006).
- 538 4. Ramulu, H. G. *et al.* Ribosomal proteins: toward a next generation standard for
539 prokaryotic systematics? *Mol. Phylogenet. Evol.* **75**, 103–117 (2014).
- 540 5. Petitjean, C., Deschamps, P., López-García, P. & Moreira, D. Rooting the domain
541 archaea by phylogenomic analysis supports the foundation of the new kingdom
542 Proteoarchaeota. *Genome Biol. Evol.* **7**, 191–204 (2014).
- 543 6. Hug, L. A. *et al.* A new view of the tree of life. *Nature Microbiology* vol. 1 (2016).
- 544 7. Fournier, G. P. & Gogarten, J. P. Rooting the ribosomal tree of life. *Mol. Biol. Evol.* **27**,
545 1792–1801 (2010).
- 546 8. Theobald, D. L. A formal test of the theory of universal common ancestry. *Nature* **465**,
547 219–222 (2010).
- 548 9. Mukherjee, S. *et al.* 1,003 reference genomes of bacterial and archaeal isolates expand
549 coverage of the tree of life. *Nat. Biotechnol.* **35**, 676–683 (2017).
- 550 10. Raymann, K., Brochier-Armanet, C. & Gribaldo, S. The two-domain tree of life is linked
551 to a new root for the Archaea. *Proceedings of the National Academy of Sciences* **112**,
552 6670–6675 (2015).
- 553 11. Creevey, C. J., Doerks, T., Fitzpatrick, D. A., Raes, J. & Bork, P. Universally distributed
554 single-copy genes indicate a constant rate of horizontal transfer. *PLoS One* **6**, e22099
555 (2011).
- 556 12. Gogarten, J. P. *et al.* Evolution of the vacuolar H⁺-ATPase: implications for the origin of

- 557 eukaryotes. *Proc. Natl. Acad. Sci. U. S. A.* **86**, 6661–6665 (1989).
- 558 13. Iwabe, N., Kuma, K., Hasegawa, M., Osawa, S. & Miyata, T. Evolutionary relationship of
559 archaeobacteria, eubacteria, and eukaryotes inferred from phylogenetic trees of
560 duplicated genes. *Proc. Natl. Acad. Sci. U. S. A.* **86**, 9355–9359 (1989).
- 561 14. Pühler, G. *et al.* Archaeobacterial DNA-dependent RNA polymerases testify to the
562 evolution of the eukaryotic nuclear genome. *Proc. Natl. Acad. Sci. U. S. A.* **86**, 4569–
563 4573 (1989).
- 564 15. Cox, C. J., Foster, P. G., Hirt, R. P., Harris, S. R. & Embley, T. M. The archaeobacterial
565 origin of eukaryotes. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 20356–20361 (2008).
- 566 16. Sugitani, K., Mimura, K., Takeuchi, M., Lepot, K. & Ito, S. Early evolution of large micro-
567 organisms with cytological complexity revealed by microanalyses of 3.4 Ga organic-
568 walled microfossils. *geobiology* **13**, 507–521 (2015).
- 569 17. Betts, H. C. *et al.* Integrated genomic and fossil evidence illuminates life’s early
570 evolution and eukaryote origin. *Nat Ecol Evol* **2**, 1556–1562 (2018).
- 571 18. van Zuilen, M. A., Lepland, A. & Arrhenius, G. Reassessing the evidence for the earliest
572 traces of life. *Nature* **418**, 627–630 (2002).
- 573 19. Horita, J. & Berndt, M. E. Abiogenic methane formation and isotopic fractionation under
574 hydrothermal conditions. *Science* **285**, 1055–1057 (1999).
- 575 20. Lepland, A., Arrhenius, G. & Cornell, D. Apatite in early Archean Isua supracrustal
576 rocks, southern West Greenland: its origin, association with graphite and potential as a
577 biomarker. *Precambrian Res.* **118**, 221–241 (2002).
- 578 21. Zhu, Q. *et al.* Phylogenomics of 10,575 genomes reveals evolutionary proximity
579 between domains Bacteria and Archaea. *Nature Communications* **10**, 5477 (2019).
- 580 22. Mirarab, S. *et al.* ASTRAL: genome-scale coalescent-based species tree estimation.
581 *Bioinformatics* **30**, i541–8 (2014).
- 582 23. Valas, R. E. & Bourne, P. E. The origin of a derived superkingdom: how a gram-positive
583 bacterium crossed the desert to become an archaeon. *Biol. Direct* **6**, 16 (2011).
- 584 24. Spang, A. *et al.* Complex archaea that bridge the gap between prokaryotes and

- 585 eukaryotes. *Nature* **521**, 173–179 (2015).
- 586 25. Jeffroy, O., Brinkmann, H., Delsuc, F. & Philippe, H. Phylogenomics: the beginning of
587 incongruence? *Trends Genet.* **22**, 225–231 (2006).
- 588 26. Lartillot, N. & Philippe, H. A Bayesian mixture model for across-site heterogeneities in
589 the amino-acid replacement process. *Mol. Biol. Evol.* **21**, 1095–1109 (2004).
- 590 27. Foster, P. G. Modeling compositional heterogeneity. *Syst. Biol.* **53**, 485–495 (2004).
- 591 28. Lartillot, N., Brinkmann, H. & Philippe, H. Suppression of long-branch attraction artefacts
592 in the animal phylogeny using a site-heterogeneous model. *BMC Evol. Biol.* **7**, S4
593 (2007).
- 594 29. Wang, H.-C., Li, K., Susko, E. & Roger, A. J. A class frequency mixture model that
595 adjusts for site-specific amino acid frequencies and improves inference of protein
596 phylogeny. *BMC Evolutionary Biology* **8**, 331 (2008).
- 597 30. Quang, L. S., Gascuel, O. & Lartillot, N. Empirical profile mixture models for
598 phylogenetic reconstruction. *Bioinformatics* **24**, 2317–2323 (2008).
- 599 31. Gouy, R., Baurain, D. & Philippe, H. Rooting the tree of life: the phylogenetic jury is still
600 out. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **370**, 20140329 (2015).
- 601 32. Tourasse, N. J. & Gouy, M. Accounting for Evolutionary Rate Variation among
602 Sequence Sites Consistently Changes Universal Phylogenies Deduced from rRNA and
603 Protein-Coding Genes. *Molecular Phylogenetics and Evolution* **13**, 159–168 (1999).
- 604 33. Segata, N., Börnigen, D., Morgan, X. C. & Huttenhower, C. PhyloPhlAn is a new method
605 for improved phylogenetic and taxonomic placement of microbes. *Nat. Commun.* **4**,
606 2304 (2013).
- 607 34. Castelle, C. J. & Banfield, J. F. Major New Microbial Groups Expand Diversity and Alter
608 our Understanding of the Tree of Life. *Cell* **172**, 1181–1197 (2018).
- 609 35. Shimodaira, H. An approximately unbiased test of phylogenetic tree selection. *Syst.*
610 *Biol.* **51**, 492–508 (2002).
- 611 36. Coleman, G. A. *et al.* A rooted phylogeny resolves early bacterial evolution. *Cold Spring*
612 *Harbor Laboratory* 2020.07.15.205187 (2020) doi:10.1101/2020.07.15.205187.

- 613 37. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic
614 Inference in the Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 615 38. Wang, H.-C., Minh, B. Q., Susko, E. & Roger, A. J. Modeling Site Heterogeneity with
616 Posterior Mean Site Frequency Profiles Accelerates Accurate Phylogenomic Estimation.
617 *Syst. Biol.* **67**, 216–235 (2018).
- 618 39. Yang, Z. A space-time process model for the evolution of DNA sequences. *Genetics*
619 **139**, 993–1005 (1995).
- 620 40. Soubrier, J. *et al.* The influence of rate heterogeneity among sites on the time
621 dependence of molecular rates. *Mol. Biol. Evol.* **29**, 3345–3358 (2012).
- 622 41. Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a new
623 software for selection of phylogenetic informative regions from multiple sequence
624 alignments. *BMC Evol. Biol.* **10**, 210 (2010).
- 625 42. Hanan, B. B. & Tilton, G. R. 60025: relict of primitive lunar crust? *Earth Planet. Sci. Lett.*
626 **84**, 15–21 (1987).
- 627 43. Barboni, M. *et al.* Early formation of the Moon 4.51 billion years ago. *Sci Adv* **3**,
628 e1602365 (2017).
- 629 44. Szöllősi, G. J., Davín, A. A., Tannier, E., Daubin, V. & Boussau, B. Genome-scale
630 phylogenetic analysis finds extensive gene transfer among fungi. *Philos. Trans. R. Soc.*
631 *Lond. B Biol. Sci.* **370**, 20140335 (2015).
- 632 45. Morel, B., Kozlov, A. M., Stamatakis, A. & Szöllősi, G. J. GeneRax: A Tool for Species-
633 Tree-Aware Maximum Likelihood-Based Gene Family Tree Inference under Gene
634 Duplication, Transfer, and Loss. *Mol. Biol. Evol.* **37**, 2763–2774 (2020).
- 635 46. Bansal, M. S., Kellis, M., Kordi, M. & Kundu, S. RANGER-DTL 2.0: rigorous
636 reconstruction of gene-family evolution by duplication, transfer and loss. *Bioinformatics*
637 **34**, 3214–3216 (2018).
- 638 47. Davín, A. A. *et al.* Gene transfers can date the tree of life. *Nat Ecol Evol* **2**, 904–909
639 (2018).
- 640 48. Boussau, B. & Scornavacca, C. Reconciling gene trees with species trees. in

- 641 *Phylogenetics in the Genomic Era* (eds. Scornavacca, C., Delsuc, F. & Galtier, N.)
642 3.2:1–3.2:23 (No commercial publisher| Authors open access book, 2020).
- 643 49. Shih, P. M. & Matzke, N. J. Primary endosymbiosis events date to the later Proterozoic
644 994 with cross-calibrated phylogenetic dating of duplicated ATPase proteins. *Proc. Natl.*
645 *Acad. Sci. U. S. A.* **110**, 996 (2013).
- 646 50. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–
647 2069 (2014).
- 648 51. Aramaki, T. *et al.* KofamKOALA: KEGG Ortholog assignment based on profile HMM and
649 adaptive score threshold. *Bioinformatics* **36**, 2251–2252 (2020).
- 650 52. Bateman, A. *et al.* The Pfam protein families database. *Nucleic Acids Res.* **32**, D138–41
651 (2004).
- 652 53. Haft, D. H., Selengut, J. D. & White, O. The TIGRFAMs database of protein families.
653 *Nucleic Acids Res.* **31**, 371–373 (2003).
- 654 54. Cantarel, B. L. *et al.* The Carbohydrate-Active EnZymes database (CAZy): an expert
655 resource for Glycogenomics. *Nucleic Acids Res.* **37**, D233–8 (2009).
- 656 55. Rawlings, N. D., Barrett, A. J. & Finn, R. Twenty years of the MEROPS database of
657 proteolytic enzymes, their substrates and inhibitors. *Nucleic Acids Res.* **44**, D343–50
658 (2016).
- 659 56. Saier, M. H., Jr, Tran, C. V. & Barabote, R. D. TCDB: the Transporter Classification
660 Database for membrane transport protein analyses and information. *Nucleic Acids Res.*
661 **34**, D181–6 (2006).
- 662 57. Søndergaard, D., Pedersen, C. N. S. & Greening, C. HydDB: A web tool for
663 hydrogenase classification and analysis. *Sci. Rep.* **6**, 34212 (2016).
- 664 58. Jones, P. *et al.* InterProScan 5: genome-scale protein function classification.
665 *Bioinformatics* **30**, 1236–1240 (2014).
- 666 59. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein
667 database search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
- 668 60. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence

- 669 similarity searching. *Nucleic Acids Research* **39**, W29–W37 (2011).
- 670 61. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
671 DIAMOND. *Nat. Methods* **12**, 59–60 (2015).
- 672 62. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jeremiin, L. S.
673 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods*
674 **14**, 587–589 (2017).
- 675 63. Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**,
676 1586–1591 (2007).
- 677 64. R Core Team. *R: A language and environment for statistical computing*. [https://www.R-](https://www.R-project.org/)
678 [project.org/](https://www.R-project.org/) (2020).
- 679 65. Wickham, H. *ggplot2: Elegant graphics for data analysis*. Springer-Verlag, New York,
680 USA, (2016).
- 681

Figure 1

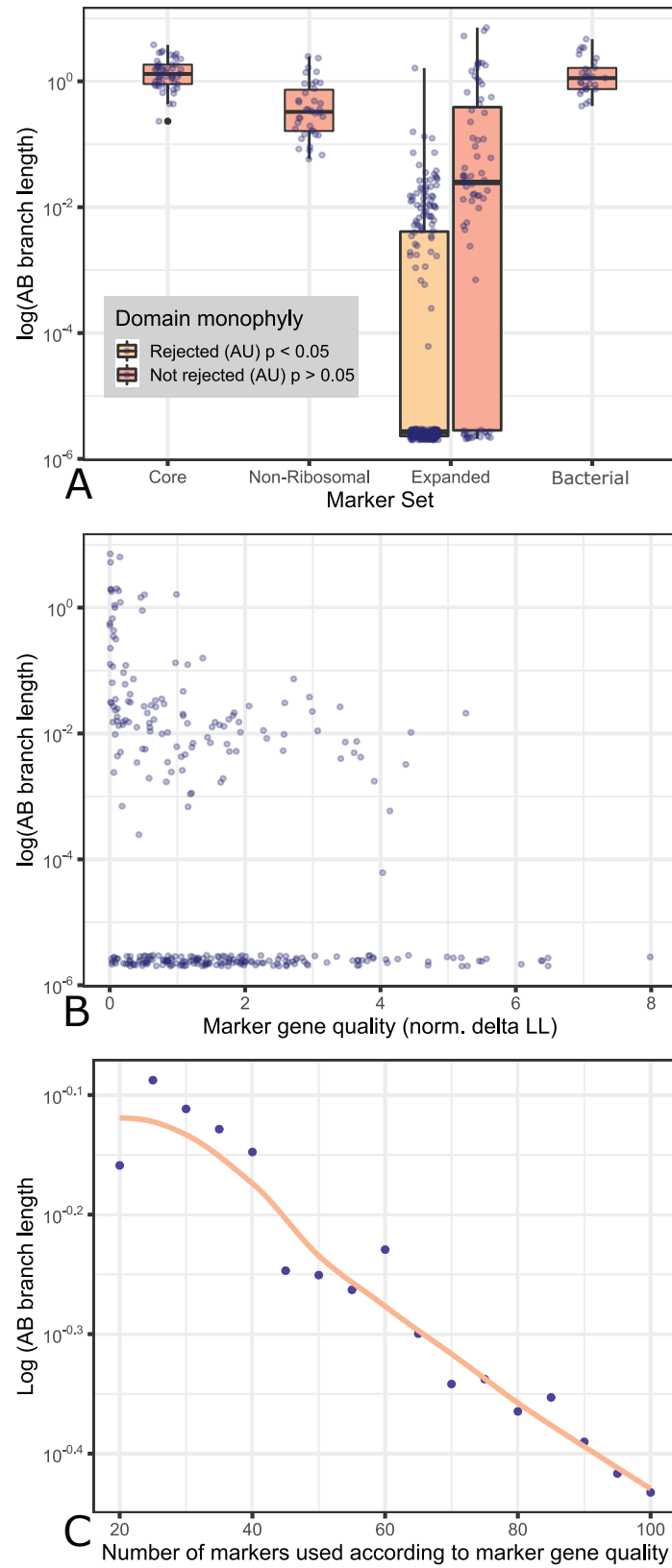


Figure 2

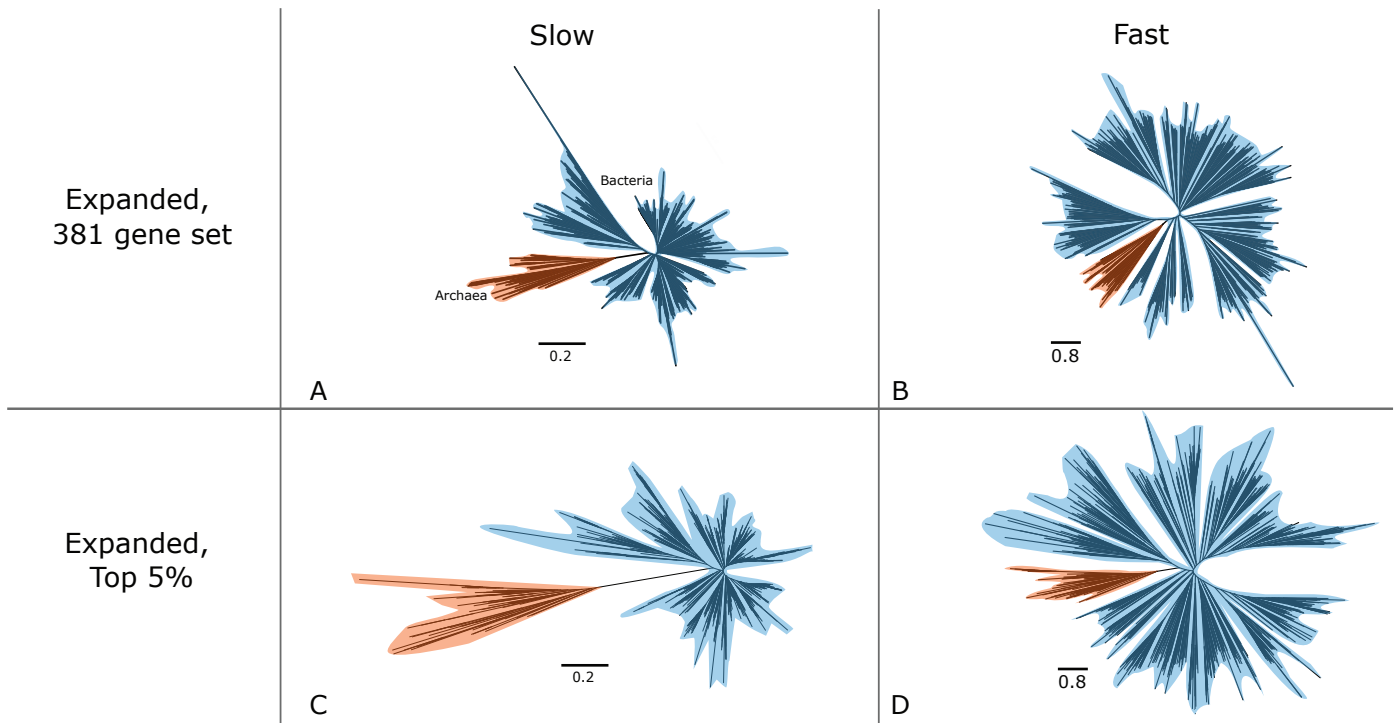


Figure 3

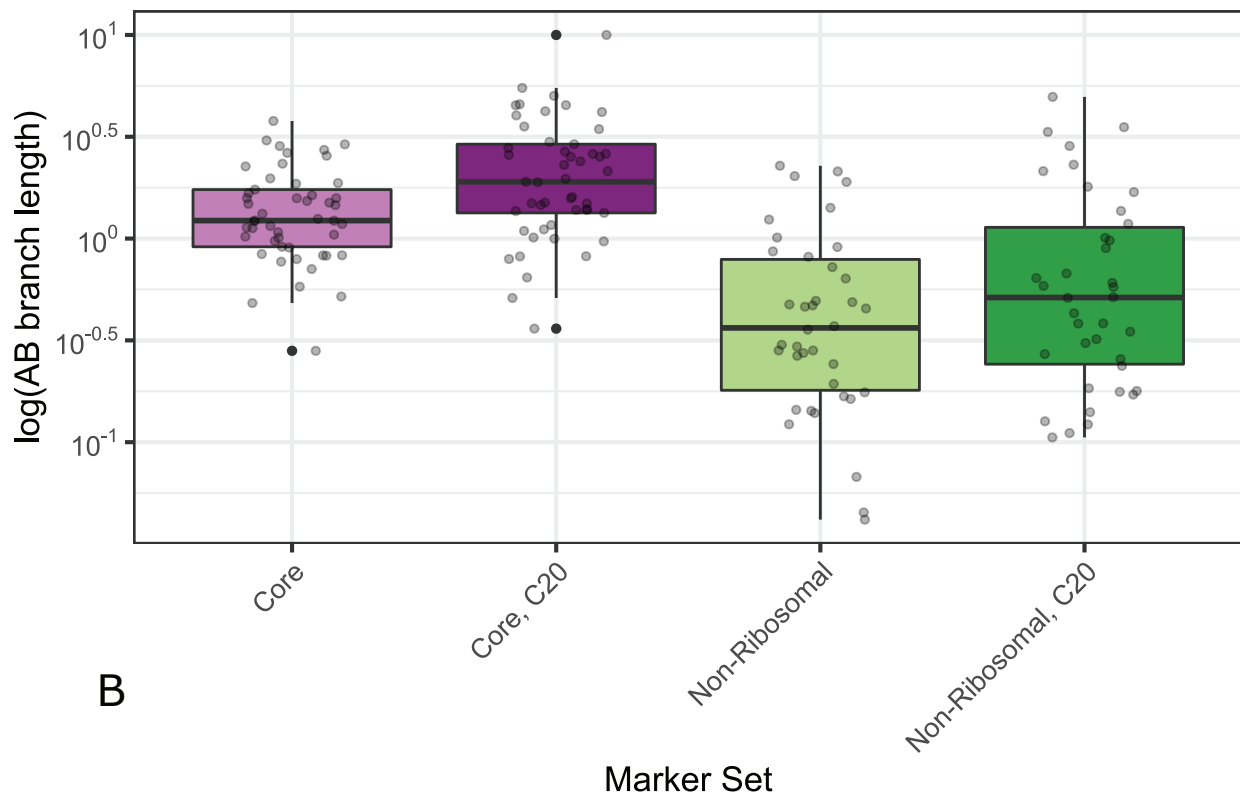
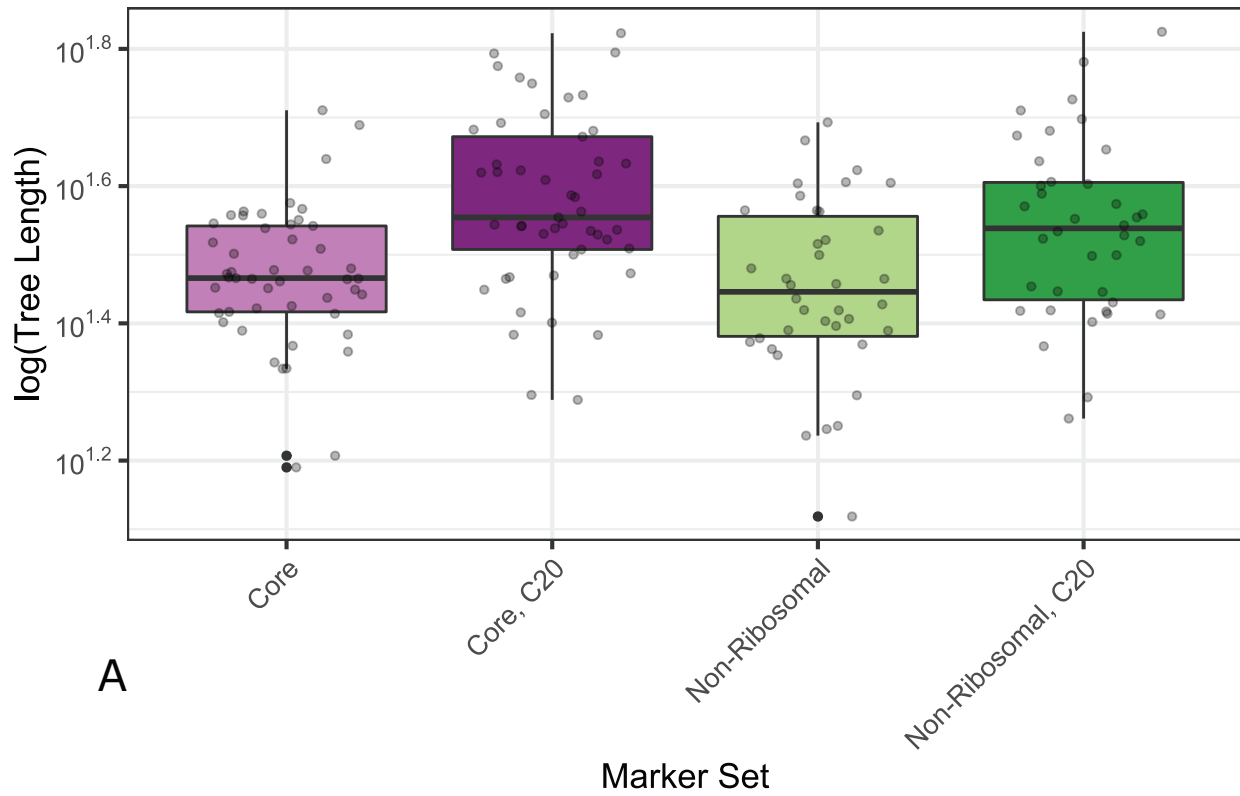


Figure 4

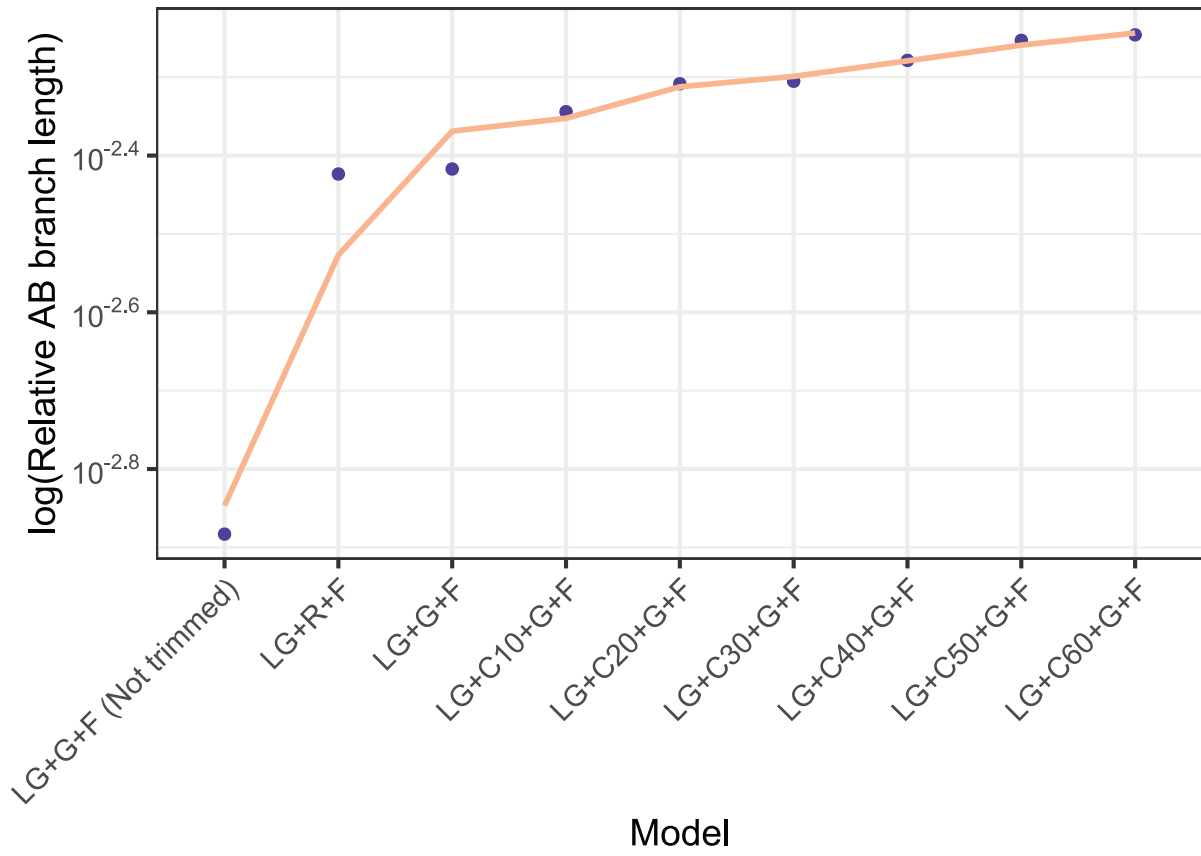


Figure 5

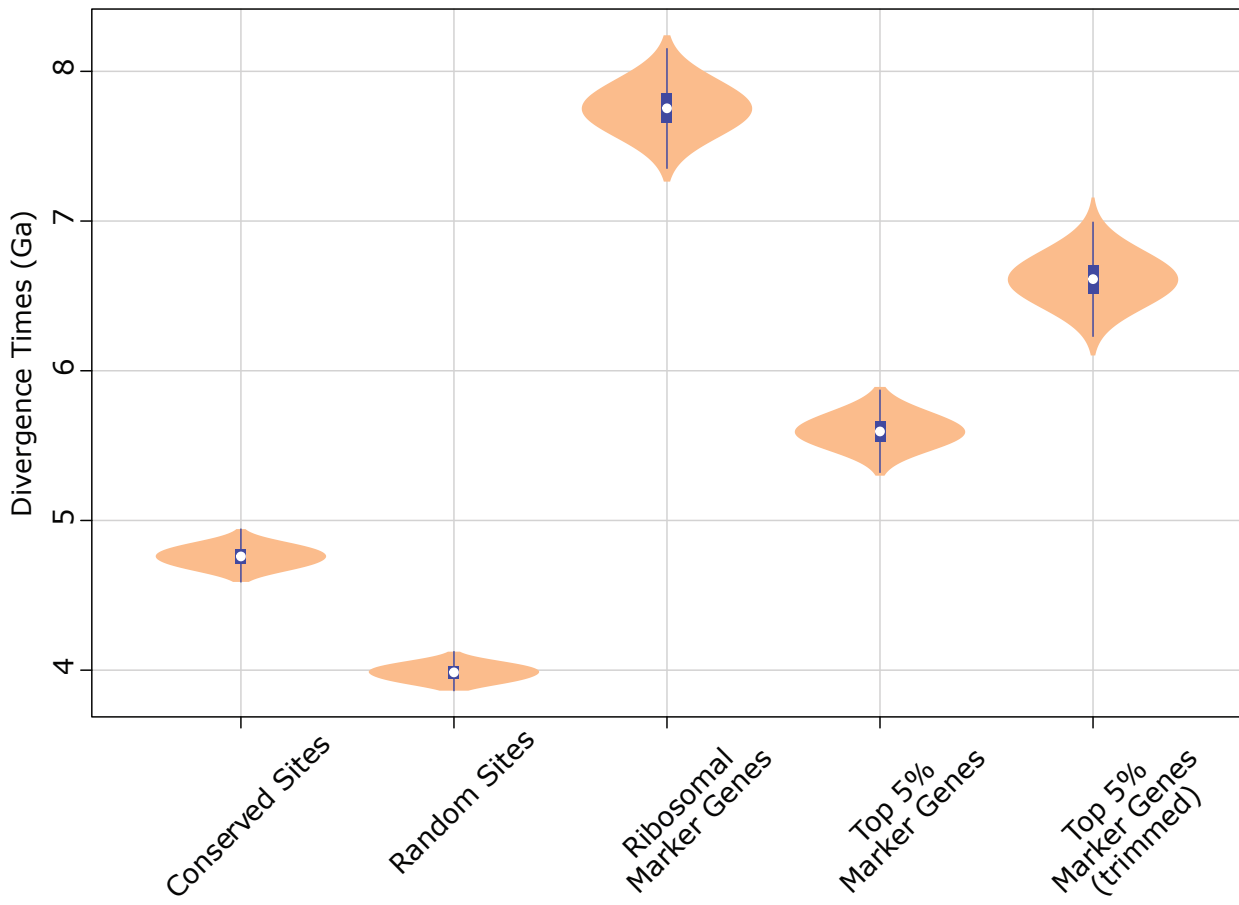


Figure 6

