1 **Phylogenomics provides robust support for a two-domains tree of life**

2

3 Tom A. Williams[1*], Cymon J. Cox[2], Peter G. Foster[3], Gergely J. Szöllősi[4,5,6], T.

4 Martin Embley[7]

5

6   1. School of Biological Sciences, University of Bristol, Life Sciences Building, Tyndall

7      Avenue, Bristol BS8 1TQ, United Kingdom.

8   2. Centro de Ciências do Mar, Universidade do Algarve, Gambelas, 8005-319 Faro,

9      Portugal.

10   3. Department of Life Sciences, Natural History Museum, London SW7 5BD, United

11      Kingdom.

12   4. MTA-ELTE "Lendület" Evolutionary Genomics Research Group, 1117 Budapest,

13      Hungary

14   5. Dept. of Biological Physics, Eötvös Loránd University, 1117 Budapest, Hungary

15   6. Evolutionary Systems Research Group, Centre for Ecological Research, Hungarian

16      Academy of Sciences, H-8237 Tihany, Hungary

17   7. Institute for Cell and Molecular Biosciences, University of Newcastle, Newcastle upon

18      Tyne NE2 4HH, United Kingdom.

19

20 *Correspondence to tom.a.williams@bristol.ac.uk

21

22

23

24

25 **Hypotheses about the origin of eukaryotic cells are classically framed within the**

26 **context of a universal "tree of life" based upon conserved core genes. Vigorous**

27 **ongoing debate about eukaryote origins is based upon assertions that the topology of**

**28**     **the tree of life depends on the taxa included and the choice and quality of genomic**

**29**     **data analysed. Here we have reanalysed the evidence underpinning those claims and**

**30**     **bring more data to bear on the question by using supertree and coalescent methods**

**31**     **to interrogate >3000 gene families in Archaea and eukaryotes. We find that eukaryotes**

**32**     **consistently originate from within the Archaea in a two-domains tree when due**

**33**     **consideration is given to the fit between model and data. Our analyses support a**

**34**     **close relationship between eukaryotes and Asgard Archaea and identify the**

**35**     **Heimdallarchaeota as the current best candidate for the closest archaeal relatives of**

**36**     **the eukaryotic nuclear lineage.**

**37**

**38**     Current hypotheses about eukaryotic origins generally propose at least two partners in that

**39**     process: a bacterial endosymbiont that became the mitochondrion and a host cell for that

**40**     endosymbiosis[1–4]. The identity of the host has been informed by analyses of conserved

**41**     genes for the transcription and translation machinery that are considered essential for

**42**     cellular life[5]. Traditionally, the host was considered to be a eukaryote based upon ribosomal

**43**     RNA trees in either unrooted[6,7] or rooted form[8]. In these trees, Archaea, Bacteria and

**44**     Eukarya form three separate primary domains, with the rooted version suggesting that

**45**     Archaea and Eukarya are more closely related to each other than to Bacteria[8]. A criticism of

**46**     these three-domains (3D) trees is that they were constructed using overly simple

**47**     phylogenetic models[5,9,10]. Phylogenetic analyses using models that better fit features of the

**48**     data[10–12], coupled with an expanded sampling of prokaryotic diversity[13–15], have supported a

**49**     two domains (2D) tree consistent with the eocyte[16] hypothesis whereby the eukaryotic

**50**     nuclear lineage - that is, the host for the mitochondrial endosymbiont - originated from within

**51**     the Archaea (reviewed in[5,17]). The 2D tree has gained increasing traction in the field[18],

**52**     particularly with the discovery of the Asgard archaea[19,20]. The Asgard archaea branch

**53**     together with eukaryotes in phylogenetic trees, and their genomes encode homologues of

**54**     eukaryotic signature proteins - that is, proteins which underpin the defining cellular structures

**55**     of eukaryotes, and which were previously thought[7,21] to be unique to eukaryotes. However,

2

56 the discoveries and analyses that support the 2D tree have been criticised from a variety of

57 perspectives.

58

59 It has been suggested[22,23] that the close relationship between eukaryotes and Asgard

60 archaea in 2D trees[19,20] is due to eukaryotic contamination of Asgard metagenomes

61 combined with phylogenetic artifacts caused by the choice of genes analysed and the

62 inclusion of fast evolving Archaea in tree reconstructions[22–24]; see also the comment[25] and

63 response[24] to those analyses. The phenomenon of long branch attraction (LBA) due to the

64 presence of fast-evolving sequences (FES) is a well-known artifact in phylogenetic

65 analyses[26–28]. Indeed, it has previously been suggested that it is the 3D tree, rather than the

66 2D tree, that is an artifact of LBA[5,9–11], both because analyses under better-fitting models

67 have recovered a 2D tree, but also because the 3D topology is one in which the two longest

68 branches in the tree of life - the stems leading to bacteria and to eukaryotes - are grouped

69 together. Nevertheless, when putative FES were removed, Forterre and colleagues[22,24]

70 recovered a monophyletic Archaea within a three-domains tree, whether analysing 35 core

71 genes, a particular subset of 6 genes, or RNA polymerases alone. Claims that the 2D tree is

72 a product of unbalanced taxonomic sampling and inclusion of FES have also been made by

73 others[29].

74

75 In a more general criticism it has been suggested[30–33] that protein sequences do not harbour

76 sufficient signal to resolve the 2D/3D debate due to mutational saturation (but see[11,12]). One

77 suggested solution is to analyse conserved structural motifs (folds) in proteins rather than

78 primary sequence data[31,33,34]. Three-dimensional structures are thought to be more highly

79 conserved than primary sequences. It has therefore been suggested that they should

80 provide a more reliable indicator of ancient relationships, although it is not yet clear how best

81 to analyse fold data for this purpose. Published unrooted trees based upon analyses of

82 protein folds have recovered Archaea, Bacteria and Eukaryotes as separate groups[34,35], a

83 result that is consistent with the 3D, but not the 2D tree. Analyses of protein folds have

84    recently been extended to use non-stationary models to infer a rooted tree of life[31]. In these

85    analyses the inferred root separated cellular life into prokaryotes (Archaea plus Bacteria,

86    termed akaryotes) and eukaryotes[31,33]. This tree is incompatible with the idea that Archaea

87    and Eukaryotes share closer common ancestry, and recapitulates the hypothesis[36] that the

88    deepest division in cellular life is between prokaryotes and eukaryotes.

89

90    In this paper, we have evaluated the analyses and data that have led to conflicting

91    hypotheses of relationships between the major groups of cellular life, and for the position of

92    the eukaryotic nuclear lineage. We have also performed phylogenomic analyses using the

93    best-available supermatrix, supertree, and coalescent methods on an expanded sample of

94    genes and taxa, to further explore the deep structure of the tree of life and the relationship

95    between archaea and eukaryotes.

96

97    **Results and Discussion**

98

99    *Analysis of core genes consistently supports two primary domains, not three*

100

101   It has recently been argued[22–24] that the 2D tree is an artifact of data and taxon sampling,

102   and that resolution of those issues provides support for a 3D tree. The molecular data at the

103   core of this debate had first been used[19] to support a 2D tree in which eukaryotes clustered

104   within Archaea as the closest relatives of the Asgard Archaea. The original dataset[19]

105   comprised a concatenation of 36 "universal" genes for 104 taxa.  In the initial critique, it was

106   claimed that the close relationship reported[19] between Asgard archaea and eukaryotes was

107   caused by the inclusion in the data set of a contaminated Elongation Factor 2 (EF2) gene for

108   *Lokiarchaeum* sample Loki3[22] (now *Heimdallarchaeota*[20]), and by the inclusion of fast-

109   evolving archaeal lineages in the analysis. However, recent data suggest that the EF2 gene

110   of *Heimdallarchaeota* is not contaminated with eukaryotic sequences because similar  EF2

111  sequences have been found in additional *Heimdallarchaeota* metagenome-assembled

112  genomes (MAGs) prepared from different environmental DNA (eDNA) samples in different

113  laboratories[20,37].

114

115  The claim[22–24] that the presence of "fast evolving sequences" (FES) might be affecting the

116  topology recovered could be seen as a reasonable challenge, since LBA can influence the

117  tree topology recovered. A problem for this specific critique[22] however, is that no single, clear

118  and consistent criterion was used to identify the "fast evolving" sequences that were

119  removed from the original dataset[19] in order to recover the 3D tree. Long-branched archaea

120  might result from either a fast evolutionary rate or a long period of time, and these

121  possibilities are difficult to distinguish *a priori*. Moreover, the historical papers[38,39] cited[22] as

122  providing topological evidence that some sequences are "fast evolving" used site- and time-

123  homogenous phylogenetic models (that is, models in which the process of evolution is

124  constant over the sites of the alignment and branches of the tree) which often fit data

125  poorly[5]. To investigate further we ranked all of the taxa in the original dataset[19] according to

126  their root-to-tip distances for each species. This is equal to the summed branch length

127  (expressed as expected number of substitutions/site) from the root of the tree (rooted

128  between Bacteria and Archaea) to the relevant tip. We calculated distributions and 95%

129  credibility intervals (Table S1) for each of these root-to-tip distances from the samples drawn

130  during an MCMC analysis under the best-fitting (see below) CAT+GTR+G4 model in

131  PhyloBayes, in order to perform Bayesian relative rates tests (Table S1). The 23 taxa

132  previously identified as FES are not the 23 taxa with the longest root-to-tip distances; while

133  some of the taxa  chosen for exclusion (*Parvarchaeum*, *Micrarchaeum, Nanoarchaeum*

134  *Nst1*, *Nanosalinarum,* and *Korarchaeum*) are indeed relatively long-branching, others

135  (*Iainarchaeum*, *Nanoarchaeum G17* and *Aenigmaarchaeon*) are in the bottom half of the

136  branch length distribution, and many of the longest-branching Archaea (including the

137  Thaumarchaeota) were retained.   Nevertheless, analysis[22] of the reduced dataset did

138  recover a 3D tree, raising the question of why this result was obtained. In the following

139    analyses we have followed the recent renaming[20] of the 3 "Loki" MAGs originally analysed

140    as *Lokiarchaeum sp. GC14_75* (formerly Loki1), *Heimdallarchaeota archaeon LC_2* (Loki2),

141    and *Heimdallarchaeota archaeon LC_3* (Loki3).

142

143    The published 3D tree[22] was recovered from the 35-gene concatenated data set under the

144    LG+G4+F model[40] in PhyML 3.1[41], with moderate support (76% bootstrap) for monophyletic

145    Archaea (Figure 5(b) in [22]). In repeating this analysis, we noted that although PhyML

146    returned a three-domains tree, analysis of the same alignment under the same substitution

147    model (LG+G4+F) with IQ-Tree 1.6.2[42] and RAxML 8.2.4[43], two other maximum-likelihood

148    phylogeny packages, instead yielded a 2D tree where Heimdallarchaeota and *Lokiarchaeum*

149    were together the sister group to eukaryotes, with a better likelihood score (Figure S1, Table

150    S2). To investigate further, we computed the log likelihoods of the 2D and 3D trees in all

151    three packages, keeping the alignment and model constant (Table S2). All three

152    implementations accord the 2D tree a higher likelihood than the 3D tree (lnl ~= -684701.2,

153    compared to ~= -684716.1 for the 3D tree). It thus appears that the recovery of a 3D tree

154    reflects a failure of PhyML to find the more likely 2D tree, rather than to the removal of

155    problematic sequences. The differences between the likelihoods are not significant

156    according to an approximately-unbiased test (AU = 0.229 for the 3D tree, 0.771 for the 2D),

157    meaning that analysis of the 35-gene dataset under LG+G4+F is equivocal with respect to

158    the 2D and 3D trees; contrary to previous claims[22], analysis of the 35-gene concatenation

159    under the LG+G4+F model provides no unambiguous evidence to prefer the 3D tree.

160

161    A number of newer models accommodate particular features of empirical data better than

162    the LG+G4+F, so we investigated which trees were produced from the 35-gene dataset

163    using these models. We addressed three issues in particular: among-site compositional

164    heterogeneity due to site-specific biochemical constraints[44], changing composition in

165    different lineages over time[45], and variations in site- and lineage-specific evolutionary rates

166    (heterotachous evolution)[46].

168    The CAT+GTR+G4 model[44,47] is an extension to the standard GTR model that allows

169    compositions to vary across sites. Analysis of the 35-gene dataset using this model

170    produced a 2D tree where eukaryotes group with Heimdallarchaeota and *Lokiarchaeum* with

171    maximal support (Figure 1). It was previously reported[22] that convergence in Bayesian

172    analyses is a problem for this data set using the CAT+GTR+G4[22] model. In our analyses, we

173    achieved good convergence between chains as assessed both by comparison of split

174    frequencies and, for the continuous parameters of the model, means and effective sample

175    sizes (Table S4). As an additional check, we also carried out ML analyses using the

176    LG+C60+G4+F model, which improves on the LG+G4+F model by modelling site-specific

177    compositional heterogeneity using a mixture of 60 composition categories. This model fits

178    the data much better than the LG+G4+F according to the BIC (Table S3) and, like

179    CAT+GTR+G4, it recovered a 2D tree with high bootstrap support (Figure S1(c)). The 3D

180    tree (AU = 0.036) could also be rejected at P < 0.05 using an AU test, based on the

181    LG+C60+G4+F model and the 35-gene alignment.

182

183    Bayesian posterior predictive simulations[48] provide a tool for evaluating the adequacy of

184    models, by testing whether data simulated under a model is similar to the empirical data.

185    Figure 2 plots the 2D tree (inferred under CAT+GTR+G4) and the 3D tree (inferred under

186    LG+G4+F in PhyML) on the same scale (Figure 2(a)), revealing that --- from the same

187    alignment --- CAT+GTR+G4 infers that many more substitutions have occurred in the core

188    gene set during the evolutionary history of life. Model fit tests (Figure 2(b), Table S4) indicate

189    that LG+G4+F provided a much poorer fit to the data (larger Z-scores) than CAT+GTR+G4

190    in terms of across-site compositional heterogeneity (Z = 64.2 for LG+G4+F, Z = 6.9 for

191    CAT+GTR+G4), and therefore systematically under-estimated the probability of convergent

192    substitutions (Z = 19.7 for LG+G4+F; Z = 7.62 for CAT+GTR+G4). These differences arise

193    because LG+G4+F assumes that amino acid frequencies are the same at all sites, whereas

194    in empirical datasets different sites have different compositions, arising from distinct

195     biochemical and selective constraints. Since this means the effective number of amino acids

196     per site is in reality lower than that predicted by LG+G4+F, the probability of parallel

197     convergence to the same amino acid in independent lineages is higher (Table S5).

198     CAT+GTR+G4 accounts for this across-site variation by incorporating site-specific

199     compositions, and is therefore less prone to underestimating rates of convergent

200     substitution. This is important because the longest branches in both the 2D and 3D trees are

201     the lineages leading to the bacteria and eukaryotes. The lesser ability of LG+G4+F to detect

202     convergent substitutions along these branches may favour inference of a 3D tree. While

203     CAT+GTR+G4 provides a better fit than LG+G4+F, neither model completely fits the

204     composition of the data (P = 0 for all tests; Table S5). As a further data exploration step, we

205     recoded[49] the amino acid alignment into four categories of biochemically similar amino acids

206     (AGNPST, CHWY, DEKQR, FILMV). Recoding has been shown to ameliorate sequence

207     saturation and compositional heterogeneity[49,50], and in this case it improved model fit (as

208     judged by the magnitude of Z-scores; Table S5). Analysis of this SR4-recoded alignment

209     under CAT+GTR+G4 recovered a 2D tree where eukaryotes grouped with the

210     Heimdallarchaeota (PP = 0.98, Figure S2).

211

212     Variation in sequence composition across the branches of the tree is also a pervasive

213     feature of data that has been used to investigate the tree of life[10,11]. We tested each of the

214     genes in the 35-gene dataset (see Methods), and found that 23/35 showed significant

215     evidence of across-branch heterogeneity at P < 0.05 (Table S6). Analysis of the

216     concatenation of the 12 composition-homogeneous genes under CAT+GTR+G4 gave a 2D

217     tree with maximal posterior support (PP = 1, Figure S3), as did a partitioned analysis using

218     the best-fitting homogeneous model for each of the 12 gene partitions (LG+G4+F in all

219     cases; Figure S3; PP = 1). We also inferred a phylogeny from the entire 35-gene dataset

220     under the branch-heterogeneous node-discrete compositional heterogeneity (NDCH)2

221     model, which explicitly incorporates changing sequence compositions across the tree.

222     NDCH2 is an extension of the NDCH model[45]; it has a separate composition vector for each

223    tree node and is constrained via a sampled concentration parameter of a Dirichlet prior.

224    Thus, the model adjusts to the level of across-branch compositional heterogeneity in the

225    data during the MCMC analysis. For reasons of computational tractability, this analysis could

226    only be run on the SR4-recoded version of the 35-gene alignment. NDCH2 obtained

227    adequate model fit with respect to across-branch compositional heterogeneity (P = 0.7838),

228    and recovered a 2D tree with Heimdallarchaeota as the sister group to eukaryotes (PP =

229    0.85; Figure S2).

230

231    A failure to account for heterotachy, or rates of molecular evolution that are both site- and

232    branch-specific, has been posited as a potential issue for phylogenomic analyses of ancient

233    core genes[51,52]. We used the GHOST[53] model of IQ-Tree to analyze the 35-gene alignment.

234    GHOST is an edge-unlinked mixture model in which the sites of the alignment evolve along

235    a shared tree topology, but are fit by a finite mixture of GTR exchangeabilities, sequence

236    compositions and branch lengths. We fit a four component mixture model to both the original

237    amino acid alignment (LG+G4+F components) and the SR4-recoded version (GTR+F

238    components). The resulting trees were a weakly-supported (amino acids; 58% bootstrap

239    support for eukaryotes plus Heimdallarchaeota and *Lokiarchaeum*) or strongly-supported

240    (recoded data; 95% bootstrap support for eukaryotes plus *Heimdallarchaeota*) 2D  tree

241    (Figure S5).

242

243    In summary, all of our analyses of the 35-gene alignment using better models recovered a

244    2D tree in which eukaryotes are either the sister group of Heimdallarchaeota plus

245    *Lokiarchaeum* or Heimdallarchaeota alone, rather than the 3D tree which the data has

246    previously been claimed[22] to support.

247

248    *Do some core genes have different histories?*

249

9

250　Based upon AU tests under the LG+G4+F model for individual genes in the 35-gene dataset,

251　it was suggested[22] that the 35-gene dataset contains two subsets of genes with different

252　evolutionary histories: a larger set supporting the 2D tree and a smaller set supporting the

253　3D tree. We used the better-fitting CAT+GTR+G4 model to analyse a concatenated dataset

254　of the 6 genes which significantly favoured the 3D tree under LG+G4+F, and we also

255　analysed a four-state recoded version of the same alignment. Analysis of the original amino

256　acids recovered a moderately-supported 3D tree, while analysis of the recoded alignment

257　recovered a weakly-supported 2D tree (Figure S4); posterior predictive simulations indicated

258　that model fit was improved by SR4 recoding (Table S7), suggesting that support for the 3D

259　tree from these 6 genes under LG+G4+F may be due to model misspecification.

260

261　It has also been suggested that phylogenetic analyses of RNA polymerase subunits[22]

262　provide robust support for a 3D tree. By contrast, other[11] analyses of RNA polymerase

263　subunits have already suggested that better fitting models prefer a 2D tree. We evaluated

264　the fit of both models, LG+G4+F and CAT+GTR+G4, used[22] to recover a 3D tree from RNA

265　polymerase subunits, using posterior predictive simulations (Supplemental Text), and found

266　that both models provide an inadequate fit to the data (Table S8). Model fit was improved

267　following SR4 recoding (Table S8), and this analysis recovered a weakly-supported and

268　poorly-resolved 2D tree (Figure S6).

269

270　*Expanded gene and taxon sampling supports a clade of eukaryotes and Asgard archaea*

271

272　We took advantage of the recent dramatic improvements in genomic and transcriptomic

273　sampling of free-living bacteria, archaea, and microbial eukaryotes to assemble a dataset of

274　125 species, including 53 eukaryotes, 39 archaea (including an expanded set of Asgard

275　MAGs[20] representing two new groups, Odinarchaeota and Thorarchaeota), and 33 bacteria,

276　on the principle that improved sampling can sometimes help to resolve difficult phylogenetic

277 problems[54,55]. We used free-living representatives of eukaryotic groups to avoid the well-

278 documented problems for tree reconstruction caused by sequences from parasitic

279 eukaryotes[26]. Our sampling of archaea and bacteria was also expanded to include

280 representatives from the large number of uncultivated lineages that have recently been

281 identified by single cell-genomics and metagenomics[15,56,57].

282

283 To further investigate the claim[22] that the tree inferred depends on the choice of universal

284 marker genes, we used the Orthologous MAtrix (OMA[58]) algorithm to identify single-copy

285 orthologues *de novo* on the 125 genome set. Benchmarks[59] indicate that OMA is

286 conservative, in that it returns a relatively low number of orthologues, but that these

287 orthologues perform better than other methods at recovering the species tree. Combining

288 OMA analysis with manual filtering to remove EF2 and genes of endosymbiotic origin (see

289 Methods), we identified 21 broadly-conserved marker genes found in at least half of our set

290 of bacteria, archaea, and eukaryotes, and 43 genes encoded by at least half of the archaea

291 and eukaryotes (see Methods). We concatenated the 21 genes conserved in all three

292 domains and inferred a tree under CAT+GTR+G4 (Figure 3a). Rooting on the branch

293 separating bacteria and archaea resulted in a 2D tree, in which eukaryotes form a

294 maximally-supported clade with Asgard archaea (Figure 3a); within Asgards, the closest

295 relatives of eukaryotes was recovered as the Heimdallarchaeota, although with only modest

296 support (PP = 0.79).

297

298 We next analyzed the expanded set of genes conserved between archaea and eukaryotes,

299 placing the root outside the TACK/Asgard/eukaryote clade as suggested by the previous

300 analysis including bacteria. The consensus tree under CAT+GTR+G4 (Figure 3b) resolves a

301 clade of eukaryotes and Heimdallarchaeota with maximal posterior support; within that

302 clade, eukaryotes group with one *Heimdallarchaeota* metagenome bin (LC3) with high (PP =

303 0.95) support.

304

305    Given ongoing debates about the impact of even single genes within concatenated datasets,

306    we investigated in detail the overlap between the 35-gene set, the 21-genes selected by

307    OMA, and a 29-gene set used in some previous analyses[10,11,14,60,61] (Table S10). After

308    removing EF2, 7 genes are found in all three sets; 27 in at least two of the three, and 50

309    genes in total are present in at least one of the datasets. We obtained the orthologues for

310    the 50 gene families from the 125 species dataset, and inferred trees using the best-fit ML

311    model in IQ-Tree on the 7-, 27- and 50-gene concatenations (Figure S8). We also expanded

312    species sampling for the 35 genes to compare with the analyses described above. Analysis

313    under the best-fitting ML model for all four concatenates resulted in a 2D tree, with either all

314    Asgards (the 7- and 35-gene datasets) or Heimdallarchaeota (27 and 50 gene datasets) as

315    sister to eukaryotes with moderate (7-gene set) to high (the other sets) bootstrap support.

316    These results indicate that there is a congruent signal for a 2D tree, and a relationship

317    between eukaryotes and Asgard archaea, that is robust to moderate differences in the

318    choice of marker genes. The results of all our concatenation analyses are summarised in

319    Table S11.

320

321    *Supertree and multispecies coalescent methods support the two-domains tree*

322

323    Concatenation allows phylogenetic signal to be pooled and permits the use of complex,

324    parameter-rich substitution models, but its assumptions are problematic in the context of

325    microbial evolution. In particular, concatenation requires that all of the genes share a

326    common phylogeny[62,63], an assumption that is difficult to test because trees inferred from

327    individual genes are often poorly supported. Some incongruence between single gene trees

328    can be attributed to stochastic error or model misspecification[14], but genuinely different

329    evolutionary histories for different genes can arise from incomplete lineage sorting, gene

330    duplication and loss, and horizontal gene transfer. We therefore investigated alternative

331    methods for integrating phylogenetic signal from multigene datasets that account for gene

332    tree incongruence in different ways. The probabilistic supertree method of Steel and Rodrigo

333    (SR2008)[64], and the Split Presence-Absence (SPA) method[65], are supertree methods that

334    model differences between gene trees as stochastic noise; ASTRAL is a supertree method

335    that is consistent under the multispecies coalescent[66]. These methods have their own

336    assumptions and limitations[63], but these are distinct from --- and provide a useful contrast to

337    --- concatenation. As these methods do not require genes to be broadly conserved across

338    the species of interest, we analyzed a set of 3199 single-copy orthologues found in at least

339    four of the taxa in our dataset (of these 3199 gene families, 479 included at least one

340    archaeon and one eukaryote; see Table S12 for the taxonomic distribution and phylogenetic

341    relationships supported by the individual trees).

342

343    All of these analyses resolved a 2D tree including a clade of eukaryotes and Asgard archaea

344    with high to maximal support (Figures S9-S10). Supertrees inferred under the SPA method

345    and ASTRAL placed eukaryotes within the Asgard archaea as the sister lineage to the three

346    Heimdallarchaeota metagenome bins (Figures S9-10), while the SR2008 supertree

347    recovered eukaryotes and Asgard archaea as monophyletic sister lineages (Figure S10). To

348    compare these supertrees independently of their models and assumptions, we calculated

349    the summed quartet distances between the set of input trees and each supertree: that is, the

350    total number of quartets (subtrees of four leaves) that differ between the input trees and

351    each supertree (Table 1). The tree with the best score by this metric was the SPA supertree

352    which, like the model-based ASTRAL analysis, recovered Heimdallarchaeota and

353    eukaryotes as sister taxa. These results suggest that there is a congruent genome-wide

354    signal for a specific relationship between eukaryotes and the Heimdallarchaeota, and that

355    the 2D tree does not appear to be an artifact of concatenation.

356

357    *Is there support from protein folds for a root between prokaryotes and eukaryotes?*

358

359    Debates about the 2D and 3D trees have typically assumed that the root of the tree lies on

360    the branch separating bacteria and archaea[67–69] or within the bacteria[70–72]. Recently, a non-

361    stationary model of binary character evolution (the KVR[73] model) was used[31,33] to infer a

362    rooted tree of life from a matrix of protein fold presence/absence data. Fold presence and

363    absence were quantified by searching HMMs corresponding to Structural Classification of

364    Proteins (SCOP) families against a set of bacterial, archaeal and eukaryotic genomes. The

365    inferred trees are intrinsically rooted because the model is non-stationary: in this model there

366    is one composition (probability of protein fold presence) at the root of the tree, and a second

367    composition elsewhere. These analyses recovered a root between prokaryotes and

368    eukaryotes[31,33], suggesting this is the primary division within cellular life and rejecting both

369    the 2D and 3D trees.

370

371    We performed simulations to evaluate the ability of the KVR model to recover the root of the

372    tree from protein fold datasets.  When data were simulated under the KVR model, the

373    method recovered the true root of the simulation tree as might be expected. However, when

374    protein fold compositions were allowed to vary over the tree, something which is observed in

375    the empirical data[31,33], the model fails to find the true root.  Under these conditions, KVR

376    finds a root on one of the branches with atypical sequence composition (see Supplementary

377    Text).  In the empirical data matrix, the eukaryotes encode significantly more protein folds

378    than either bacteria or archaea (median of 871 folds per eukaryotic genome, compared to

379    521 for archaea and 615 for bacteria; $P < 10^{-8}$ for the eukaryote-archaea and eukaryote-

380    bacteria comparisons, $P = 0.000278$ comparing bacteria and archaea; $n = 47$ eukaryotes, 47

381    bacteria and 47 archaea, Wilcoxon rank-sum tests), but their higher compositions are in the

382    minority because the matrix contains an equal number of genomes from each of the three

383    domains. Thus, the inferred root between prokaryotes and eukaryotes may result from the

384    model's bias in placing the root on a branch with atypical composition; in simulations, the

385    root inference can be controlled by varying which composition among tips - high or low - is in

386    the majority (Supplementary Text). These results agree with recent work[72,74] in suggesting

387    that non-reversible models may provide reliable rooting information when the assumptions of

388    the model are met, but that root inferences are sensitive to model misspecification. The KVR

389    model is only one of the many possible non-stationary and non-homogeneous models, and

390    does not appear to be well-suited to these data. Models that better describe the process by

391    which fold (or sequence) compositions change through time, and across the tree --- or

392    indeed those that make use of other sources of time information[75,76] --- may perform better

393    for rooting deep phylogenies. How best to root ancient radiations remains an open question,

394    and method development is still at an early stage. A key challenge will be the development

395    of methods that account for the heterogeneity of the evolutionary process across the data

396    and through evolutionary time (that is, across the branches of the tree).

397

398    A potentially bigger problem than model misspecification for the published analyses[31,33] is

399    their assumption that the entire protein fold set evolves on a single underlying tree. This

400    assumption is unlikely to be realistic because of the different histories generated by

401    widespread horizontal gene transfer and, in eukaryotes, by endosymbiotic gene transfer

402    from the bacterial progenitors of mitochondria and plastids[77]. The assumption of a single

403    underlying tree to explain fold distributions also means that, despite claims to the contrary[31],

404    the published analyses cannot be used to reject the 2D tree because, as generally

405    formulated[5,16,78], it seeks to explain the inheritance of only a subset of the genes on cellular

406    genomes.

407

408    To evaluate whether the protein folds in the published matrix[31,33] share a common

409    evolutionary tree, we inferred single-gene phylogenies for each fold (Supplementary Text).

410    Although weakly supported, these trees are consistent with there being extensive

411    disagreement between single fold-based topologies: only 22 of the protein folds supported

412    the monophyly of eukaryotes, and none recovered all three domains as potentially

413    monophyletic groups, even though this was the consensus topology obtained from analysis

414    of the complete matrix. The trees contained signals for sister-group relationships between

415    eukaryotes and Alphaproteobacteria (the most frequent sister-group among the protein folds

416    shared between eukaryotes and bacteria) and for a relationship between eukaryotes and the

417    TACKL archaea. These analyses are consistent with endosymbiotic theory[2,79] and the ideas

418    that underpin the 2D tree, namely that eukaryotes contain a mixture of genes from the

419    archaeal host cell and the bacterial endosymbiont that became the mitochondrion[2,3,5]

420    (Supplemental Text).

421

422    *Conclusions*

423

424    Identifying the tree that best depicts the relationships between the major groups of life is

425    important for understanding eukaryotic origins and the evolution of the complexity that

426    distinguishes eukaryotic cells. It has recently been asserted that the tree recovered depends

427    upon the species investigated and the choice and quality of the molecular data analysed[22,23].

428    In the present study we have investigated the data sets used to underpin these claims and

429    find no compelling evidence to support them. Analyses using better-fitting phylogenetic

430    models consistently recovered a 2D tree[5,10,12,16,17,19,20] wherein eukaryotes are most closely

431    related to members of the recently discovered Asgard archaea. These results are also

432    supported by additional analyses of expanded concatenations and increased species

433    sampling, and from large-scale genome-wide data sets analysed using supertree and

434    coalescence methods.

435

436    We also investigated support from analyses of whole-genome protein folds for a rooted

437    universal tree in which the deepest division is between prokaryotes and eukaryotes. Taken

438    at face-value this tree would reject the 2D and 3D trees that are the focus of robust

439    discussion in the current literature[24,25]. However, while protein structure is a useful guide to

440    identifying homology when primary sequence similarity is weak, how best to analyse fold

441    data to resolve deep phylogenetic relationships is still not clear. Published analyses[31] do not

442    account for the varied evolutionary histories of individual folds due to endosymbiosis and

443    gene transfer, and our simulations suggest that root inference under existing models is

444    unreliable and affected by variation in the abundance and distribution of folds across

445    genomes. At present, the best supported root is on the branch separating bacteria and

446    archaea[67,68,80,81] or among the bacteria[70,72], and the hypothesis that eukaryotes are younger

447    than prokaryotes is supported by a range of phylogenetic, cell biological[2,3] and

448    palaeontological[61,82–84] evidence.

449

450    Our analyses and published trees[5,10,20] imply that the eukaryotic nuclear lineage evolved

451    from within the Archaea. They provide robust phylogenomic support for a clade of

452    eukaryotes and Asgard archaea, and identify the Heimdallarchaeota as the best candidate

453    among sampled lineages[19,20,85] for a sister group to eukaryotes. This sister group relationship

454    will no doubt change with further sampling of the potentially vast archaeal diversity in nature

455    still to be discovered. The prize will be ever more reliable inferences of the features that

456    were in place in the last common ancestor of both groups and an improved evidence-based

457    understanding of the building blocks that underpinned the transition from prokaryotic to

458    eukaryotic cells.

459

460    **Methods**

461

462    *Sequences and alignment*

463    For the reanalyses of the Da Cunha et al. and Spang et al. datasets, alignments were

464    obtained from the supplementary material of Da Cunha et al.[22], and the EF2 gene removed

465    according to the coordinates provided; the alignments from Spang et al. (2015) were

466    generously provided by the authors. OMA 2.1.1[58] was used to identify putative single-copy

467    orthologues among a dataset of 92 eukaryotic, archaeal and bacterial genomes. For putative

468    orthologues present in at least half of the sampled species, single gene trees were inferred

469    for each candidate under the LG+G4+F model in IQ-Tree, and the trees were manually

470    inspected to filter out eukaryotic genes that were acquired from the mitochondrial or plastid

471    endosymbionts. We also performed a BLASTP screen to identify organellar genes that might

472    have been missed via the tree inspection approach. This procedure resulted in a set of 43

473    single-copy orthologues shared between archaea and eukaryotes, and 21 genes shared

474    among all three domains, that were used for concatenation-based phylogenomic analyses.

475    For all OMA gene families found in at least four species, we used a BLASTP-based screen

476    to identify and filter out eukaryotic gene families of bacterial origin, resulting in 3261 gene

477    families in four or more species that are either eukaryote-specific inventions, or shared

478    between eukaryotes and archaea. For the comparisons of core gene sets, an iterative

479    process of manual comparisons, similarity searches and tree building was used to identify

480    common and distinct markers in the published sets, identify seed sequences for each marker

481    in the genomes of *Dictyostelium discoideum, Sulfolobus solfataricus* and *Escherichia coli*

482    *K12*, and build HMMs for each marker using the existing datasets. We used domain-specific

483    HMM searches in HMMER3[86] followed by the reciprocal best hit criterion against our

484    domain-specific reference genomes to identify candidate orthologues, followed by gene tree

485    inference and manual curation to assemble final marker sets. Sequences were aligned using

486    the L-INS-i mode in Mafft 7[87], and poorly aligning regions identified and removed using the

487    BLOSUM30 model in BMGE 1.12[88].

488

489    *Phylogenetics*

490    Maximum likelihood analyses were performed using IQ-Tree 1.6.2[42], and bootstrap supports

491    were computed using UFBoot2[89], except where indicated in the main text. Model fitting was

492    carried out using the MFP mode in IQ-Tree, adding the empirical site profile models (C20-

493    C60) to the default candidate model set. Bayesian phylogenies were inferred under the

494    CAT+GTR+G4 model in PhyloBayes-MPI 1.8[47], using the bpcomp and tracecomp programs

495    to monitor convergence of two MCMC chains for each analysis. Posterior predictive

496    simulations were performed using readpb_mpi in PhyloBayes. Tests for across-branch

497 compositional heterogeneity were performed in p4[62]: we inferred maximum-likelihood gene

498 trees for each of the 35 genes in the concatenation, then simulated data for each gene under

499 the LG+G4+F model. A Chi-square statistic reflecting compositional heterogeneity was

500 calculated on the original and simulated datasets, and the values from the simulated data

501 were used as a null distribution with which to evaluate the test statistic from the original data.

502

503 *Supertrees*

504 Supertrees were inferred from the maximum likelihood phylogenies for each single gene,

505 with substitution models chosen as described above. MRP, SR2008 and SPA supertrees

506 were inferred using p4[65]. Multispecies coalescent trees were inferred using ASTRAL-III[66].

507

508 **Data availability:** The data associated with our analyses are available in the FigShare

509 repository[91] at https://doi.org/10.6084/m9.figshare.8950859.v2.

510

511 **References**

512 1.   Embley, T. M. & Martin, W. Eukaryotic evolution, changes and challenges. *Nature* **440**,

513      623–630 (2006).

514 2.   Martin, W. F., Garg, S. & Zimorski, V. Endosymbiotic theories for eukaryote origin.

515      *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **370**, (2015).

516 3.   Roger, A. J., Muñoz-Gómez, S. A. & Kamikawa, R. The Origin and Diversification of

517      Mitochondria. *Curr. Biol.* **27**, R1177–R1192 (2017).

518 4.   Martijn, J. & Ettema, T. J. G. From archaeon to eukaryote: the evolutionary dark ages of

519      the eukaryotic cell. *Biochem. Soc. Trans.* **41**, 451–457 (2013).

520 5.   Williams, T. a., Foster, P. G., Cox, C. J. & Embley, T. M. An archaeal origin of

521      eukaryotes supports only two primary domains of life. *Nature* **504**, 231–236 (2013).

522 6.   Woese, C. R. & Fox, G. E. Phylogenetic structure of the prokaryotic domain: the primary

523      kingdoms. *Proc. Natl. Acad. Sci. U. S. A.* **74**, 5088–5090 (1977).

524   7.   Kurland, C. G., Collins, L. J. & Penny, D. Genomics and the irreducible nature of

525        eukaryote cells. *Science* **312**, 1011–1014 (2006).

526   8.   Woese, C. R., Kandler, O. & Wheelis, M. L. Towards a natural system of organisms:

527        proposal for the domains Archaea, Bacteria, and Eucarya. *Proc. Natl. Acad. Sci. U. S.*

528        *A.* **87**, 4576–4579 (1990).

529   9.   Tourasse, N. J. & Gouy, M. Accounting for evolutionary rate variation among sequence

530        sites consistently changes universal phylogenies deduced from rRNA and protein-

531        coding genes. *Mol. Phylogenet. Evol.* **13**, 159–168 (1999).

532   10.  Cox, C. J., Foster, P. G., Hirt, R. P., Harris, S. R. & Embley, T. M. The archaebacterial

533        origin of eukaryotes. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 20356–20361 (2008).

534   11.  Foster, P. G., Cox, C. J. & Embley, T. M. The primary divisions of life: a phylogenomic

535        approach employing composition-heterogeneous methods. *Philos. Trans. R. Soc. Lond.*

536        *B Biol. Sci.* **364**, 2197–2207 (2009).

537   12.  Raymann, K., Brochier-Armanet, C. & Gribaldo, S. The two-domain tree of life is linked

538        to a new root for the Archaea. *Proceedings of the National Academy of Sciences*

539        201420858 (2015).

540   13.  Guy, L. & Ettema, T. J. G. The archaeal 'TACK' superphylum and the origin of

541        eukaryotes. *Trends Microbiol.* **19**, 580–587 (2011).

542   14.  Williams, T. a., Foster, P. G., Nye, T. M. W., Cox, C. J. & Embley, T. M. A congruent

543        phylogenomic signal places eukaryotes within the Archaea. *Proc. Biol. Sci.* **279**, 4870–

544        4879 (2012).

545   15.  Hug, L. A. *et al.* A new view of the tree of life. *Nat Microbiol* **1**, 16048 (2016).

546   16.  Lake, J. a., Henderson, E., Oakes, M. & Clark, M. W. Eocytes: a new ribosome structure

547        indicates a kingdom with a close relationship to eukaryotes. *Proc. Natl. Acad. Sci. U. S.*

548        *A.* **81**, 3786–3790 (1984).

549   17.  Eme, L., Spang, A., Lombard, J., Stairs, C. W. & Ettema, T. J. G. Archaea and the origin

550        of eukaryotes. *Nat. Rev. Microbiol.* **15**, nrmicro.2017.133 (2017).

551   18.  Williams, T. A., Embley, T. M., Williams, T. A. & Embley, T. M. Changing ideas about

552     eukaryotic origins. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* (2015).

553  19. Spang, A. *et al.* Complex archaea that bridge the gap between prokaryotes and

554     eukaryotes. *Nature* 173–179 (2015).

555  20. Zaremba-Niedzwiedzka, K. *et al.* Asgard archaea illuminate the origin of eukaryotic

556     cellular complexity. *Nature* **541**, 353 (2017).

557  21. Hartman, H. & Fedorov, A. The origin of the eukaryotic cell: a genomic investigation.

558     *Proc. Natl. Acad. Sci. U. S. A.* **99**, 1420–1425 (2002).

559  22. Da Cunha, V., Gaia, M., Gadelle, D., Nasir, A. & Forterre, P. Lokiarchaea are close

560     relatives of Euryarchaeota, not bridging the gap between prokaryotes and eukaryotes.

561     *PLoS Genet.* **13**, e1006810 (2017).

562  23. Gaia, M., Da Cunha, V. & Forterre, P. The Tree of Life. in *Molecular Mechanisms of*

563     *Microbial Evolution* (ed. Rampelotto, P. H.) 55–99 (Springer International Publishing,

564     2018).

565  24. Da Cunha, V., Gaia, M., Nasir, A. & Forterre, P. Asgard archaea do not close the debate

566     about the universal tree of life topology. *PLoS genetics* **14**, e1007215 (2018).

567  25. Spang, A. *et al.* Asgard archaea are the closest prokaryotic relatives of eukaryotes.

568     *PLoS genetics* **14**, e1007080 (2018).

569  26. Hirt, R. P. *et al.* Microsporidia are related to Fungi: evidence from the largest subunit of

570     RNA polymerase II and other proteins. *Proc. Natl. Acad. Sci. U. S. A.* **96**, 580–585

571     (1999).

572  27. Lartillot, N., Brinkmann, H. & Philippe, H. Suppression of long-branch attraction artefacts

573     in the animal phylogeny using a site-heterogeneous model. *BMC Evol. Biol.* **7 Suppl 1**,

574     S4 (2007).

575  28. Bergsten, J. A review of long-branch attraction. *Cladistics* **21**, 163–193 (2005).

576  29. Nasir, A., Kim, K. M., Da Cunha, V. & Caetano-Anollés, G. Arguments Reinforcing the

577     Three-Domain View of Diversified Cellular Life. *Archaea* **2016**, 1851865 (2016).

578  30. Penny, D., McComish, B. J., Charleston, M. a. & Hendy, M. D. Mathematical elegance

579     with biochemical realism: the covarion model of molecular evolution. *J. Mol. Evol.* **53**,

580   711–723 (2001).

581 31. Harish, A. & Kurland, C. G. Empirical genome evolution models root the tree of life.

582   *Biochimie* **138**, 137–155 (2017).

583 32. Philippe, H. & Forterre, P. The rooting of the universal tree of life is not reliable. *J. Mol.*

584   *Evol.* **49**, 509–523 (1999).

585 33. Harish, A. & Kurland, C. G. Akaryotes and Eukaryotes are independent descendants of

586   a universal common ancestor. *Biochimie* **138**, 168–183 (2017).

587 34. Yang, S., Doolittle, R. F. & Bourne, P. E. Phylogeny determined by protein domain

588   content. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 373–378 (2005).

589 35. Caetano-Anolles, G. An Evolutionarily Structured Universe of Protein Architecture.

590   *Genome Research* **13**, 1563–1571 (2003).

591 36. Mayr, E. Two empires or three? *Proc. Natl. Acad. Sci. U. S. A.* **95**, 9720–9723 (1998).

592 37. Narrowe, A. B. *et al.* Complex evolutionary history of translation Elongation Factor 2 and

593   diphthamide biosynthesis in Archaea and parabasalids. *bioRxiv* 262600 (2018).

594   doi:10.1101/262600

595 38. Brochier, C., Forterre, P. & Gribaldo, S. Archaeal phylogeny based on proteins of the

596   transcription and translation machineries: tackling the Methanopyrus kandleri paradox.

597   *Genome Biol.* **5**, R17 (2004).

598 39. Brochier, C., Gribaldo, S., Zivanovic, Y., Confalonieri, F. & Forterre, P. Nanoarchaea:

599   representatives of a novel archaeal phylum or a fast-evolving euryarchaeal lineage

600   related to Thermococcales? *Genome Biol.* **6**, R42 (2005).

601 40. Le, S. Q. & Gascuel, O. An improved general amino acid replacement matrix. *Mol. Biol.*

602   *Evol.* **25**, 1307–1320 (2008).

603 41. Guindon, S. *et al.* New Algorithms and Methods to Estimate Maximum-Likelihood

604   Phylogenies: Assessing the Performance of PhyML 3.0. *Syst. Biol.* **59**, 307–321 (2010).

605 42. Nguyen, L. T., Schmidt, H. A., Von Haeseler, A. & Minh, B. Q. IQ-TREE: A fast and

606   effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol.*

607   *Evol.* **32**, 268–274 (2015).

608    43. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of

609        large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).

610    44. Lartillot, N. & Philippe, H. A Bayesian mixture model for across-site heterogeneities in

611        the amino-acid replacement process. *Mol. Biol. Evol.* **21**, 1095–1109 (2004).

612    45. Foster, P. Modeling Compositional Heterogeneity. *Syst. Biol.* **53**, 485–495 (2004).

613    46. Zhou, Y., Brinkmann, H., Rodrigue, N., Lartillot, N. & Philippe, H. A dirichlet process

614        covarion mixture model and its assessments using posterior predictive discrepancy

615        tests. *Mol. Biol. Evol.* **27**, 371–384 (2010).

616    47. Lartillot, N. L., Odrigue, N. I. R., Tubbs, D. A. S. & Icher, J. A. R. PhyloBayes MPI :

617        Phylogenetic Reconstruction with Infinite Mixtures of Profiles in a Parallel Environment.

618        **62**, 611–615 (2013).

619    48. Bollback, J. P. Bayesian model adequacy and choice in phylogenetics. *Mol. Biol. Evol.*

620        **19**, 1171–1180 (2002).

621    49. Susko, E. & Roger, A. J. On reduced amino acid alphabets for phylogenetic inference.

622        *Mol. Biol. Evol.* **24**, 2139–2150 (2007).

623    50. Hrdy, I. *et al.* Trichomonas hydrogenosomes contain the NADH dehydrogenase module

624        of mitochondrial complex I. *Nature* **432**, 618–622 (2004).

625    51. Whelan, S. Spatial and temporal heterogeneity in nucleotide sequence evolution. *Mol.*

626        *Biol. Evol.* **25**, 1683–1694 (2008).

627    52. Gouy, R., Baurain, D. & Philippe, H. Rooting the tree of life: the phylogenetic jury is still

628        out. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **370**, 20140329 (2015).

629    53. Crotty, S. M. *et al.* GHOST: Recovering Historical Signal from Heterotachously-evolved

630        Sequence Alignments. *bioRxiv* 174789 (2019). doi:10.1101/174789

631    54. Graybeal, A. Is it better to add taxa or characters to a difficult phylogenetic problem?

632        *Syst. Biol.* **47**, 9–17 (1998).

633    55. Hedtke, S. M., Townsend, T. M. & Hillis, D. M. Resolution of phylogenetic conflict in

634        large data sets by increased taxon sampling. *Syst. Biol.* **55**, 522–529 (2006).

635    56. Castelle, C. J. & Banfield, J. F. Major New Microbial Groups Expand Diversity and Alter

636      our Understanding of the Tree of Life. *Cell* **172**, 1181–1197 (2018).

637   57. Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes

638      substantially expands the tree of life. *Nat Microbiol* **2**, 1533–1542 (2017).

639   58. Roth, A. C. J., Gonnet, G. H. & Dessimoz, C. Algorithm of OMA for large-scale orthology

640      inference. *BMC Bioinformatics* **9**, 518 (2008).

641   59. Altenhoff, A. M. *et al.* Standardized benchmarking in the quest for orthologs. *Nat.*

642      *Methods* **13**, 425–430 (2016).

643   60. Williams, T. a. & Embley, T. M. Archaeal 'dark matter' and the origin of eukaryotes.

644      *Genome Biol. Evol.* **6**, 474–481 (2014).

645   61. Betts, H. C. *et al.* Integrated genomic and fossil evidence illuminates life's early

646      evolution and eukaryote origin. *Nat Ecol Evol* **2**, 1556–1562 (2018).

647   62. Roch, S. & Steel, M. Likelihood-based tree reconstruction on a concatenation of aligned

648      sequence data sets can be statistically inconsistent. *Theor. Popul. Biol.* **100C**, 56–62

649      (2015).

650   63. Roch, S., Nute, M. & Warnow, T. Long-Branch Attraction in Species Tree Estimation:

651      Inconsistency of Partitioned Likelihood and Topology-Based Summary Methods. *Syst.*

652      *Biol.* **68**, 281–297 (2019).

653   64. Steel, M. & Rodrigo, A. Maximum likelihood supertrees. *Syst. Biol.* **57**, 243–250 (2008).

654   65. Akanni, W. A., Wilkinson, M., Creevey, C. J., Foster, P. G. & Pisani, D. Implementing

655      and testing Bayesian and maximum-likelihood supertree methods in phylogenetics. *R*

656      *Soc Open Sci* **2**, 140436 (2015).

657   66. Zhang, C., Sayyari, E. & Mirarab, S. ASTRAL-III: Increased Scalability and Impacts of

658      Contracting Low Support Branches. in *Comparative Genomics* 53–75 (Springer, Cham,

659      2017).

660   67. Iwabe, N., Kuma, K., Hasegawa, M., Osawa, S. & Miyata, T. Evolutionary relationship of

661      archaebacteria, eubacteria, and eukaryotes inferred from phylogenetic trees of

662      duplicated genes. *Proc. Natl. Acad. Sci. U. S. A.* **86**, 9355–9359 (1989).

663   68. Gogarten, J. P. *et al.* Evolution of the vacuolar H+-ATPase: implications for the origin of

664     eukaryotes. *Proc. Natl. Acad. Sci. U. S. A.* **86**, 6661–6665 (1989).

665   69. Fournier, G. P. & Gogarten, J. P. Rooting the ribosomal tree of life. *Mol. Biol. Evol.* **27**,

666     1792–1801 (2010).

667   70. Lake, J. a., Skophammer, R. G., Herbold, C. W. & Servin, J. a. Genome beginnings:

668     rooting the tree of life. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **364**, 2177–2185 (2009).

669   71. Cavalier-Smith, T. Rooting the tree of life by transition analyses. *Biol. Direct* **1**, 19

670     (2006).

671   72. Williams, T. A. *et al.* New substitution models for rooting phylogenetic trees. *Philos.*

672     *Trans. R. Soc. Lond. B Biol. Sci.* (2015).

673   73. Klopfstein, S., Vilhelmsen, L. & Ronquist, F. A Nonstationary Markov Model Detects

674     Directional Evolution in Hymenopteran Morphology. *Syst. Biol.* **64**, 1089–1103 (2015).

675   74. Cherlin, S. *et al.* The effect of non-reversibility on inferring rooted phylogenies.

676     *Molecular Biology and Evolution* **35**, 984–1002 (2018).

677   75. Tria, F. D. K., Landan, G. & Dagan, T. Phylogenetic rooting using minimal ancestor

678     deviation. *Nature Ecology & Evolution* **1**, s41559–017–0193 (2017).

679   76. Szöllõsi, G. J., Rosikiewicz, W., Boussau, B., Tannier, E. & Daubin, V. Efficient

680     exploration of the space of reconciled gene trees. *Syst. Biol.* **62**, 901–912 (2013).

681   77. Timmis, J. N., Ayliffe, M. a., Huang, C. Y. & Martin, W. Endosymbiotic gene transfer:

682     organelle genomes forge eukaryotic chromosomes. *Nat. Rev. Genet.* **5**, 123–135

683     (2004).

684   78. McInerney, J. O., O'Connell, M. J. & Pisani, D. The hybrid nature of the Eukaryota and a

685     consilient view of life on Earth. *Nat. Rev. Microbiol.* **12**, 449–455 (2014).

686   79. Gray, M. W. & Doolittle, W. F. Has the endosymbiont hypothesis been proven?

687     *Microbiol. Rev.* **46**, 1–42 (1982).

688   80. Brown, J. R. & Doolittle, W. F. Root of the universal tree of life based on ancient

689     aminoacyl-tRNA synthetase gene duplications. *Proc. Natl. Acad. Sci. U. S. A.* **92**, 2441–

690     2445 (1995).

691   81. Zhaxybayeva, O., Lapierre, P. & Gogarten, J. P. Ancient gene duplications and the

692      root(s) of the tree of life. *Protoplasma* **227**, 53–64 (2005).

693  82.  Knoll, A. H. Paleobiological perspectives on early eukaryotic evolution. *Cold Spring*

694      *Harb. Perspect. Biol.* **6**, (2014).

695  83.  Butterfield, N. J. Early evolution of the Eukaryota. *Palaeontology* **58**, 5–17 (2015).

696  84.  Parfrey, L. W., Lahr, D. J. G., Knoll, A. H. & Katz, L. a. Estimating the timing of early

697      eukaryotic diversification with multigene molecular clocks. *Proc. Natl. Acad. Sci. U. S. A.*

698      **108**, 13624–13629 (2011).

699  85.  Spang, A. *et al.* Proposal of the reverse flow model for the origin of the eukaryotic cell

700      based on comparative analyses of Asgard archaeal metabolism. *Nature Microbiology*

701      (2019). doi:10.1038/s41564-019-0406-9

702  86.  Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**, e1002195

703      (2011).

704  87.  Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:

705      improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).

706  88.  Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a new

707      software for selection of phylogenetic informative regions from multiple sequence

708      alignments. *BMC Evol. Biol.* **10**, 210 (2010).

709  89.  Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2:

710      Improving the Ultrafast Bootstrap Approximation. *Mol. Biol. Evol.* **35**, 518–522 (2018).

711  90.  Williams, T. A. *et al.* Integrative modeling of gene and genome evolution roots the

712      archaeal tree of life. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E4602–E4611 (2017).

713  91. Williams et al. Data from: Phylogenomics provides robust support for a two-domains

714      tree of life. Figshare fileset https://doi.org/10.6084/m9.figshare.8950859.v2 (2019)

715

716  **Acknowledgements**

725

726    **Author contributions:** All authors contributed to the conception and design of the project,

727    and interpretation of results. TAW, CJC, PGF and GJSz performed analyses. TAW and TME

728    wrote the manuscript, with input from all authors.

729

730    **Competing interests statement:** The authors declare they have no competing interests.

731

732    **Figure legends**

733

734    **Figure 1: The 35-gene matrix of Da Cunha et al. favours a two-domains tree using the**

735    **best-fitting models in both maximum likelihood and Bayesian analyses.** The

736    eukaryotes (green) group with the sampled Asgard archaea (orange) with maximum

737    posterior support. Bacteria are in grey, TACK Archaea in yellow, Euryarchaeota in blue. This

738    is a consensus tree inferred under the CAT+GTR+G4 model in PhyloBayes-MPI; branch

739    lengths are proportional to the expected number of substitutions per site, as indicated by the

740    scale bar. A 2D topology was obtained under a variety of other models in ML analyses

741    (LG+G4+F, LG+PMSF+G4, LG+C60+G4+F; Figure S1), and also with 4-state Susko-Roger

742    recoding under the CAT+GTR+G4 and NDCH2 models (Figure S2).

743

**Figure 2. Evidence that the three-domains tree is an artifact of long branch attraction.**

(a) Da Cunha et al. analysed a dataset of 35 core protein-coding genes under the LG+G4+F model and obtained a 3D tree; the better-fitting (Table S4) CAT+GTR+G4 model recovers a 2D tree. Bootstrap support (a) and Bayesian posterior probability (b) are indicated for the key nodes defining the 3D and 2D trees. "Asgard" refers to a clade of Heimdallarchaeota and Lokiarchaeum. Plotting these trees to the same scale (in terms of substitutions per site) illustrates major differences in these analyses. The 3D/LG+G4+F analysis suggests that, on average, 30.77 changes have taken place per site; the two-domains/CAT+GTR+G4 analysis suggests that 47.4 changes per site have occurred. This difference amounts to ~128,511 additional substitutions in total inferred under the CAT+GTR+G4 model. (b) Posterior predictive tests indicate that CAT+GTR+G4 performs significantly better than LG+G4+F in capturing the site-specific evolutionary constraints reflected by lower biochemical diversity approaching that of the empirical data). This results in more realistic estimates of substitutional saturation and convergence found in the data. The longest branches on both the 3D and 2D tree are the stems leading to the bacteria and eukaryotes (in blue and green, respectively). CAT+GTR+G4 identifies many more convergent substitutions on these branches than does LG+G4+F, as can be seen by comparing the branch lengths in (a). This failure to detect convergent substitutions under LG+G4+F has the effect of drawing the bacterial and eukaryotic branches together, because convergences are mistaken for homologies (synapomorphies), resulting in a 3D tree.

**Figure 3: An expanded sampling of microbial diversity supports a two-domains tree.**

(a) Bayesian phylogeny of 21 concatenated proteins conserved across bacteria, archaea and eukaryotes under the CAT+GTR+G4 model, rooted on the branch separating bacteria and archaea. Eukaryotes group with Asgard archaea with maximum posterior support. (b) Bayesian phylogeny of 43 genes conserved between Archaea and eukaryotes under CAT+GTR+G4. Eukaryotes group with, or within, Heimdallarchaeota. All support values are Bayesian posterior probabilities, and branch lengths are proportional to the expected number

772    of substitutions per site, as indicated by the scale bars. The Euryarchaeota are paraphyletic

773    in the consensus tree in (a), consistent with some recent analyses using bacterial

774    outgroups[11,12], although the relevant support values are low and the analysis does not

775    robustly exclude the alternative hypothesis[90] of a monophyletic Euryarchaeota. The tree in

776    (b) is formally unrooted because it does not include a bacterial outgroup. Based on (a) and

777    published analyses[12,90], the root may lie between the Euryarchaeota and the other taxa, or

778    within the Euryarchaeota. Amino acid data were recoded using the 4-state scheme of Susko

779    and Roger, which our posterior predictive simulations (Table S7) suggest improved model fit

780    by ameliorating substitutional saturation and compositional heterogeneity; phylogenies

781    inferred on the original amino acid data are provided in Figure S7.

782

783    **Tables**

784

| Supertree method | Summed quartet distance | Asgard-eukaryote relationship |
|---|---|---|
| SR2008 | 17287838 | Sister groups |
| MSC (ASTRAL) | 17213379 | Eukaryotes with Heimdallarchaeota (0.28 quadripartition support) |
| **SPA** | **17195042** | Eukaryotes with Heimdallarchaeota (BPP 1.0t) |

785    **Table 1: Summed quartet distances between the supertrees produced by several**

786    **methods and the set of 3199 input trees.** All trees recover a clade of eukaryotes and

787    Asgard archaea; in addition, the SPA and ASTRAL trees place eukaryotes within Asgard

788    archaea, as the sister group to the Heimdallarchaeota.

*Thermotoga maritima*
*Synechocystis* PCC 6803
*Bacillus subtilis*
*Rhodopirellula baltica*
*Chlamydia trachomatis*
*Rickettsia prowazekii*
*Escherichia coli*
*Bacteroides thetaiotaomicron*
*Campylobacter jejuni*
*Borrelia burgdorferi*

Bacteria

*Methanocaldococcus jannaschii*
*Methanotorris igneus*
*Methanococcus maripaludis*
*Methanothermus fervidus*
*Methanothermobacter thermoautotrophicus*
*Methanosphaera stadtmanae*
*Methanobacterium alcaliphilum*
*Methanomassiliicoccus luminyensis*
*Aciduliprofundum boonei*
*Thermoplasma acidophilum*
*Picrophilus torridus*
*Ferroplasma acidarmanus*
*Ferroglobus placidus*
*Archaeoglobus fulgidus*
*Methanocorpusculum sp.*
*Methanoplanus petrolearius*
*Methanoculleus marisnigri*
*Methanospirillum hungatei*
*Methanosphaerula palustris*
*Methanosaeta thermophila*
*Methanosarcina acetivorans*
*Methanohalobium evestigatum*
*Methanocella paludicola*
*Haloferax volcanii*
*Halalkalicoccus jeotgali*
*Halobacterium salinarum* NRC1
*Haloarcula marismortui*
*Thermococcus kodakarensis*
*Pyrococcus furiosus*

Euryarchaeota

*Caldiarchaeum subterraneum*
*Nitrososphaera sp.*
*Cenarchaeum symbiosum*
*Nitrosopumilus maritimus*
*Nitrosoarchaeum limnia*
*Nitrosoarchaeum koreensis*
*Thermofilum pendens*
*Vulcanisaeta distributa*
*Caldivirga maquilingensis*
*Thermoproteus uzoniensis*
*Pyrobaculum calidifontis*
*Pyrobaculum aerophilum*
*Ignicoccus hospitalis*
*Staphylothermus marinus*
*Thermosphaera aggregrans*
*Desulfurococcus kamchatkensis*
*Pyrolobus fumarii*
*Hyperthermus butylicus*
*Aeropyrum pernix*
*Acidilobus saccharovorans*
*Fervidicoccus fontis*
*Ignisphaera aggregans*
*Sulfolobus tokodaii*
*Sulfolobus acidocaldarius*
*Sulfolobus solfataricus*
*Sulfolobus islandicus*
*Acidianus hospitalis*
*Metallosphaera sedula*
*Metallosphaera cuprina*

TACK Archaea

*Lokiarchaeum sp.* GC14_75
*Heimdallarchaeota archaeon* LC_2
*Heimdallarchaeota archaeon* LC_3

Asgard

*Trichomonas vaginalis*
*Entamoeba histolytica*
*Tetrahymena thermophila*
*Leishmania infantum*
*Plasmodium falciparum*
*Thalassiosira pseudonana*
*Arabidopsis thaliana*
*Dictyostelium discoideum*
*Saccharomyces cerevisiae*
*Homo sapiens*

Eukaryotes

0.5

(a)

1 substitution/site

Archaea

Asgard

Eukaryotes

76

Bacteria

**Three-domains tree**

Archaea

Asgard

Eukaryotes

1

**Two-domains tree**

Bacteria

(b)

*Model performance*

Observed

**CAT+GTR** (Z = 6.9)

**LG** (Z = 64.2)

7.5        8        8.5        9

Site-specific biochemical diversity

**LG** (Z = 19.7)

Observed

**CAT+GTR** (Z = 7.62)

0.42        0.43        0.44        0.45        0.46

Empirical convergence probability