

1 **The asgardarchaeal-unique contribution to protein families**
2 **of the eukaryotic common ancestor was 0.3%**

3
4

5 Michael Knopp¹, Simon Stockhorst¹, Mark van der Giezen², Sriram G. Garg¹, Sven B. Gould¹

6

7 ¹Institute for Molecular Evolution, Heinrich-Heine-University Düsseldorf, 40225 Düsseldorf,
8 Germany

9 ²Centre for Organelle Research, University of Stavanger, 4021 Stavanger, Norway

10

11

12

13 **Keywords:** eukaryote origin, eukaryogenesis, asgardarchaea, compartmentalization

14

15

16 **Significance Statement**

17 Ever since the first report of a new archaeal lineage, the asgardarchaea, their metagenome
18 analyses have encouraged continued speculations on a type of cell biology ranging between
19 that of prokaryotes and eukaryotes. While it appears a tempting notion, recent microscopic
20 images of an asgardarchaeon suggest otherwise. We inspected the origin of eukaryotic protein
21 families with respect to their distribution across bacteria and archaea. This reveals that the
22 protein families shared exclusively between asgardarchaea and eukaryotes amounts to only
23 0.3% of the protein families conserved across all eukaryotes. Asgardarchaeal diversity is
24 likely unrivaled across archaea, but their cell biology remains prokaryotic in nature and lends
25 support for the importance of endosymbiosis in evolving eukaryotic traits.

26 **Summary**

27 The difference between pro- and eukaryotic biology is evident in their genomes, cell biology,
28 and evolution of complex and macroscopic body plans. The lack of intermediates between the
29 two types of cells places the endosymbiotic acquisition of the mitochondrion through an
30 archaeal host at the event horizon of eukaryote origin. The identification of eukaryote specific
31 proteins in a new archaeal phylum, the asgardarchaea, has fueled speculations about their
32 cellular complexity, suggesting they could be eukaryote-like. Here we analyzed the coding
33 capacity of 150 eukaryotes, 1000 bacteria, and 226 archaea, including the only cultured
34 member of the asgardarchaea, Candidatus *Promethoarchaeon syntrophicum* MK-D1.
35 Established clustering methods that recover endosymbiotic contributions to eukaryotic
36 genomes, recover an asgardarchaeal-unique contribution of a mere 0.3% to protein families
37 present in the last eukaryotic common ancestor, while simultaneously suggesting that
38 asgardarchaeal diversity rivals that of all other archaea combined. Furthermore, we show that
39 the number of homologs shared exclusively between asgardarchaea and eukaryotes is only 27
40 on average. Genomic and in particular cellular complexity remains a eukaryote-specific
41 feature and, we conclude, is best understood as the archaeal host's solution to housing an
42 endosymbiont and not as a preparation for obtaining one.

43

44

45 **Introduction**

46 Four billion years of prokaryotic evolution has only once resulted in the emergence of highly
47 compartmentalized cells and eventually macroscopic body plans: following the origin of
48 eukaryotes through endosymbiosis. The analysis of core eukaryotic features such as the
49 nucleus, mitochondria, sex and meiosis, compartmentalization and dynamic membrane
50 trafficking, and virtually all of the associated protein families, consistently point to their
51 presence in the last eukaryotic common ancestor (LECA) (Fritz-Laylin et al. 2010,
52 Koumandou et al. 2013, Koonin et al. 2013, Garg and Martin 2016). We possess a reasonable
53 understanding of the basic cellular features and coding capacity of LECA, owing to the
54 growing number of genome sequences spanning all of eukaryotic diversity. All eukaryotes
55 stem from a single ancestor that in terms of cellular and genomic complexity rivaled those of
56 extant eukaryotic supergroups (Fritz-Laylin et al. 2010, Koumandou et al. 2013, Koonin et al.
57 2013). It is certain that LECA was a product of the integration of an alphaproteobacterium
58 into an archaeal host following endosymbiosis (Lane 2011, Blackstone 2013, Martin et al.
59 2015, Dacks et al. 2016, Spang et al. 2019).

60

61 Through the description of the asgardarchaea, current debates once again concern the cellular
62 complexity of the host that came to house the endosymbiont and what contribution the
63 mitochondrion could have played in establishing the eukaryotic cell (Dacks et al. 2016,
64 Martin et al. 2017). Asgardarchaea, a novel phylum assembled from metagenome data, are
65 viewed as bridging the gap between pro- and eukaryotic cells, because they encode proteins
66 homologous to eukaryotic ones that are e.g. involved in intracellular vesicle trafficking and
67 the regulation of actin cytoskeleton dynamics (Spang et al. 2015, Zaremba-Niedzwiedzka et
68 al. 2017, Neveu et al. 2020). The cellular complexity of the host cell that acquired the
69 alphaproteobacterial endosymbiont has been a matter of speculation ever since the realization
70 that endosymbiosis was pivotal in the transition to eukaryotic life. Modern models of
71 eukaryogenesis differ regarding the timing of mitochondrial acquisition, the extent of the
72 cellular complexity of the host, and the selective reasons provided for explaining the
73 presence, function, and emergence of eukaryotic traits prior or ensuing endosymbiosis
74 (O'Malley 2010, Martin et al. 2015, Gould et al. 2016, Tria et al. 2019, Vosseberg et al.
75 2020).

76

77 Understanding the steps of eukaryogenesis is a demanding intellectual challenge that explores
78 the past of life and one of its most radical transitions. It holds the key to understanding the
79 steps towards cellular complexity, the timing of mitochondrial entry, and what limits
80 prokaryotes to frequently evolve eukaryote-like complexity. Was eukaryogenesis really a
81 matter of luck (Booth and Doolittle 2015) and how important was the energetic superiority
82 provided by the mitochondrion to the host cell (Lane and Martin 2010, Lynch and Marinov
83 2015, Lane and Martin 2016)? Any model that views endosymbiosis as some kind of terminal
84 coincidence on the evolutionary roadmap to the eukaryotic domain of life needs to explain the
85 singularity that is eukaryogenesis and the lack of comparable complexity among prokaryotes.

86

87 A consistent motivation for speculating on the archaeal host cell's grade of complexity is
88 trying to understand whether the host cell was phagocytotic or not (Cavalier-Smith 1987,
89 Yutin et al. 2009, Martijn and Ettema 2013, Martin et al. 2017). This is complicated by the
90 description of a phagocytosis-like process in a planctomycete (Shiratori et al. 2019) and the
91 conflicting evidence for intracellular prokaryotic endosymbionts in the absence of
92 phagocytosis (Fenchel and Bernard 1993, Emblay and Finlay 1993, Schmid 2003, Zientz et
93 al. 2004, Duplessis et al. 2004, Thacker 2005, Woyke et al. 2006, Husnik et al. 2013). It has

94 also been concluded that asgardarchaea are not phagocytotic (Burns et al. 2018), although
95 they encode actin-regulating profilins (Akil and Robinson 2018), small Rab-like GTPases
96 (Surkont and Pereira-Leal 2016), and prototypic SNARE proteins (Neveu et al. 2020).
97 Phagocytosis might have evolved multiple times independently (Yutin et al. 2009, Mills
98 2020) and is a mode of feeding, which is incompatible with the syntrophic foundation that
99 underpins eukaryogenesis (Martin et al. 2015, Spang et al. 2019, Martin and Müller 1998,
100 Vellai et al. 1998, Imachi et al. 2020). A sole focus on this single eukaryotic trait might
101 distract and furthermore discounts the complexity of the transition that was involved. What is
102 certain is that images of an asgardarchaeon, Candidatus *Promethoarchaeum syntrophicum*
103 MK-D1, reveal cells with typical archaeal morphology, half a micron in diameter, with
104 obligate syntrophy, and devoid of intracellular complexity (Imachi et al. 2020).

105

106 Here we clustered the available genomes of 150 eukaryotes, 1000 bacteria and 226 archaea
107 (including asgardarchaea metagenomic assemblies, and for comparison the complete genome
108 of the cultured Candidatus *P. syntrophicum* strain MK-D1) in order to evaluate the
109 asgardarchaeal-unique contribution to eukaryogenesis that is understood as support for early
110 cellular complexity in asgardarchaea.

111

112

113 **Results**

114

115 In order to evaluate to what degree asgardarchaea bridge the prokaryotic and eukaryotic
116 protein families, we performed a global comparison of clustered gene families across eleven
117 asgardarchaeal metagenome-assembled genomes (MAGs), the closed MAG of the cultured
118 asgardarchaeon Candidatus *P. syntrophicum* MK-D1, 214 other archaea, 1,000 bacteria, and
119 150 eukaryotes. Protein families for 150 eukaryotes were taken from Brueckner and Martin
120 (Brueckner and Martin 2020), which included 239,012 clusters. We further clustered proteins
121 from the prokaryotes, resulting in 352,384 bacterial clusters and 49,855 archaeal clusters.
122 Subsequently, the eukaryote and prokaryote clusters were merged in a reciprocal best cluster
123 approach previously described in Ku et al. 2015 (Ku et al. 2015), yielding Eukaryote-
124 Prokaryote-Clusters (EPCs). These EPCs contained proteins from eukaryotes and proteins
125 from either archaeal (Eukaryote-Archaea clusters, EA) or bacterial (Eukaryote-Bacteria
126 clusters, EB), or both (Eukaryote-Archaea-Bacteria clusters, EAB). This approach yielded
127 2,590 EPCs, of which 867 or 33.5% (330 EA clusters + 537 EAB clusters; Suppl. Table 1)

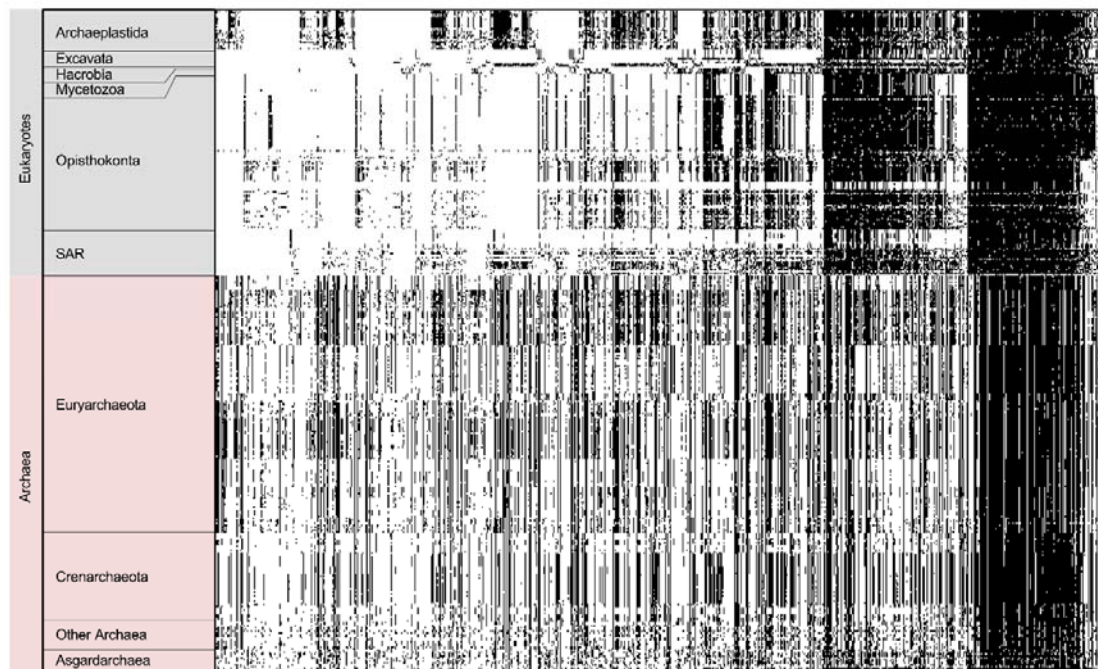
128 contained an archaeal component. Among these 867 EPCs, asgardarchaeal protein sequences,
129 including those of *Candidatus P. syntrophicum*, were present in about 75% of the protein
130 families (75% of EAB clusters and 75.2% of EA clusters).

131

132 A presence-absence pattern (PAP) of all 867 protein families with archaeal contribution to the
133 EPCs, including the 537 protein families present in all domains is shown in Fig. 1 (and Suppl.
134 Table 1). While gene distributions among eury- and crenarchaeota is highly similar, those of
135 the asgardarchaea are patchier and more diverse. Among all of our 239,012 eukaryotic
136 clusters, we could identify only six EA clusters with asgardarchaeal-unique contributions to
137 eukaryotes (Suppl. Table 2), representing 0.0025% of all extant eukaryotic diversity. To
138 calculate the asgardarchaeal-unique contribution to LECA, we filtered the eukaryotic clusters
139 for those that include at least one representative of each of the six eukaryotic supergroups,
140 resulting in 1880 LECA clusters and consequently an asgardarchaeal-unique contribution of
141 0.3191%.

142

FIG 1



143

144

145 **Fig. 1: Presence-absence pattern (PAP) of all 867 eukaryote-prokaryote clusters (EPCs)**
146 **with archaeal contribution among the investigated 150 eukaryotes, 212 archaea and twelve**
147 **asgardarchaea. The protein families were sorted by their distribution among six eukaryotic**
148 **supergroups (SAR; Stramenopila, Alveolata, Rhizaria), their presence is indicated in black**
149 **along the X-axis. The group of “Other Archaea” is comprised of all archaea within our**

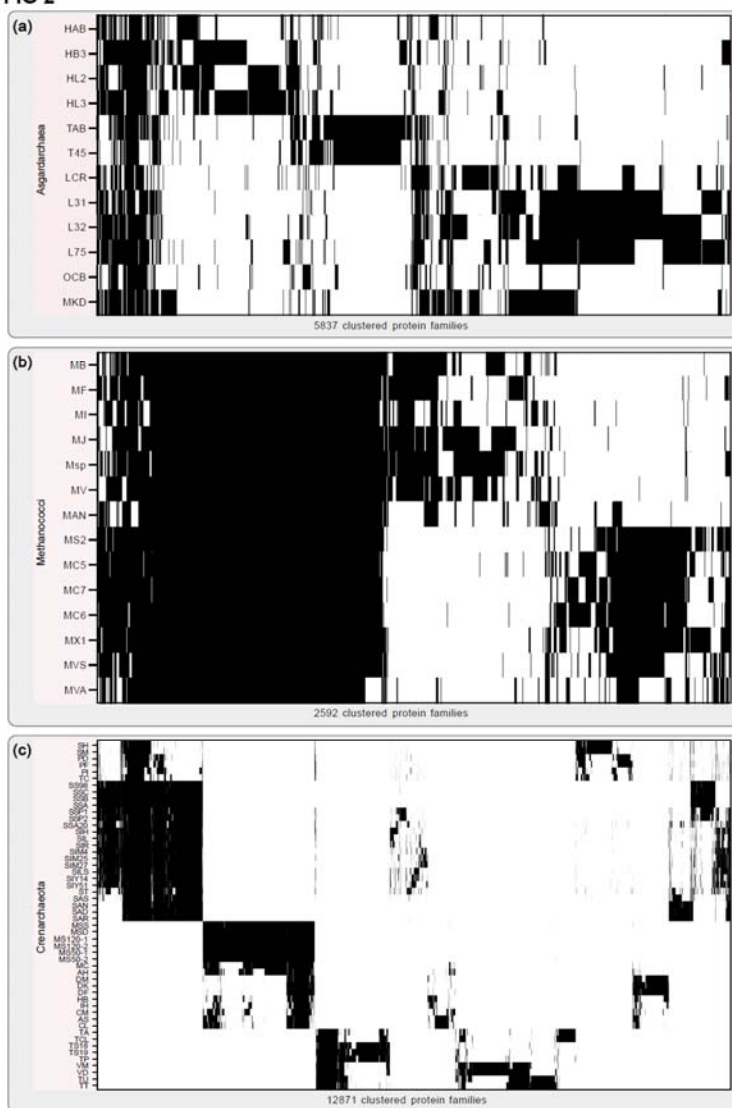
150 *database that have less than 15 members. The distribution of these EPCs among the*
151 *asgardarchaea does not reveal a particular pattern such as that seen for the other archaeal*
152 *groups, demonstrating their large genetic differences.*

153

154 But how closely are the asgardarchaea related to each other with respect to other archaea? To
155 investigate their kinship, we constructed protein families of only the asgardarchaeal protein
156 sequences, which generated a set of 5,837 protein families. The distribution of these protein
157 families among the asgardarchaea reveals a pronounced diversity, with only a small
158 proportion of the protein families being shared across all of them, which suggests one is
159 dealing with a kind of superphylum (Fig. 2a). This is evident from a comparison wherein, we
160 clustered all protein sequences of 14 members of the genus *Methanococci* and 52
161 crenarchaeote members from the TACK superphylum (Fig. 2b and Fig. 2c), revealing 2,592
162 and 12,871 protein families, respectively.

163

FIG 2



164

165

166 **Fig. 2: Archaeal protein family distributions.** (a) Distribution of all 5,837 calculated
167 asgardarchaeal protein families among the investigated asgardarchaea. The protein families
168 were obtained by globally comparing all asgardarchaeal protein sequences in a pairwise all-
169 vs-all Diamond BLASTp approach including subsequent clustering via MCL. The result was
170 sorted via hierarchical clustering along the X-axis. The vast majority of protein families is not
171 shared among all asgardarchaea but they are rather specific to the individual taxon.
172 *Candidatus Prometheoarchaeon syntrophicum MK-D1* shares the highest number of its
173 protein families with members of the Lokiarchaeota. (b) For comparison, we calculated
174 protein families for 14 members of the class Methanococci in the same manner, generating a
175 total of 2,592 protein families. Hierarchical clustering revealed a striking difference in
176 protein families shared between members within the two clusterings, underlining the
177 investigated asgardarchaea's diversity among each other. (c) On the contrary, clustering of
178 52 members of the crenarchaeota, a highly diverse taxonomic group, results in 12,871 protein
179 families. The matrix was sorted along the X- and Y-axis via hierarchical clustering. All
180 Abbreviations used are listed in Suppl. Table 6.

181

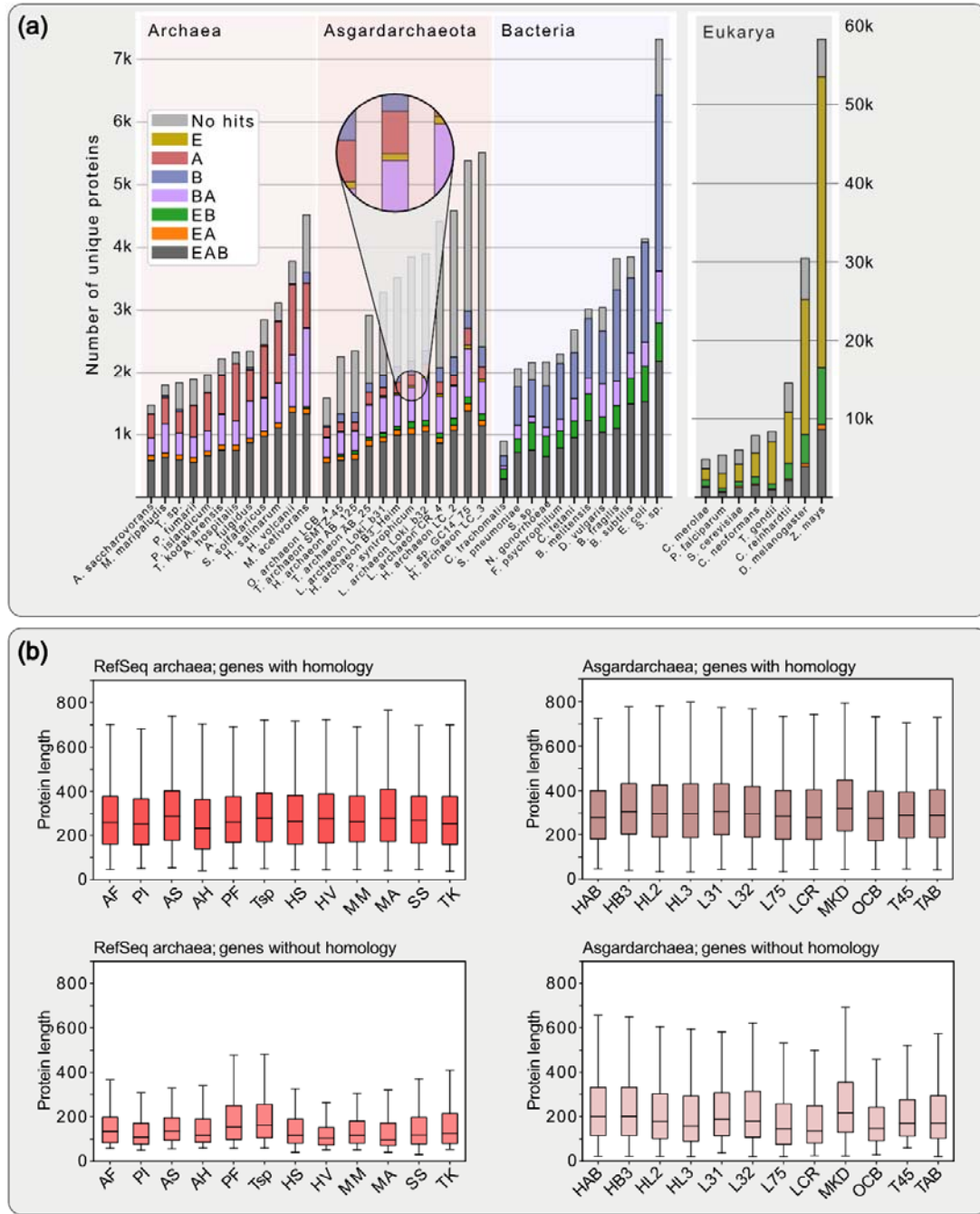
182 A key difference between archaea and eukaryotes is the difference in the number of protein
183 families. In a previous study, 239,012 protein families were generated on the basis of all
184 pairwise protein sequence comparisons of 150 eukaryotic proteomes (Brueckner and Martin
185 2020). The same approach yields a total of 49,855 protein families on the basis of all pairwise
186 protein sequence comparisons between 226 archaea, including 11 asgardarchaea and
187 *Candidatus P. syntrophicum* MK-D1. The patchy distribution of the asgardarchaeal EPCs
188 (Fig. 1) and the high diversity between each other (Fig. 2a) prompts a closer look, but the vast
189 difference in the number of protein families reflects the sudden inflation of protein family
190 emergence at the origin of eukaryotes.

191

192 Since many asgardarchaeal protein sequences could not be clustered due to a lack of
193 homology, we conducted an analysis of BLASTp hits against a bacterial database consisting
194 of 5,443 proteomes, an archaeal database consisting of 212 proteomes (excluding
195 asgardarchaea) and a eukaryotic database consisting of 150 proteomes. For all investigated
196 asgardarchaea and a selection of twelve well known and diverse archaea and bacteria each, as
197 well as eight eukaryotes, we quantified the amount of sequences with homologs in one of
198 these databases, any combination of these, or no significant homologs at all (Fig. 3a), while
199 ignoring hits from the same genus to counter database composition biases. For asgardarchaea,
200 only 27 protein sequences on average have homologs unique to eukaryotes (Fig 3a, magnified
201 area), supporting the initial EPC analysis which only uncovered six protein families that were
202 unique to eukaryotes and asgardarchaea. Furthermore, this test reveals that high proportions
203 of the asgardarchaeal proteomes do not retrieve significant hits in our three tested databases

204 (Fig. 3a, grey bars). In a few cases, such as for example for Candidatus *Heimdallarchaeota*
 205 *LC_3* or Candidatus *Lokiarchaeota archaeon CR_4*, the number of proteins for which no
 206 homology was detected in any other species, exceeds half of the respective genome's coding
 207 capacity (Fig. 3a).
 208

FIG 3



209

210 **Fig. 3: Distribution of protein homologs and length comparison of archaeal genes with and**
211 **without homology. (a) BLAST hit analysis comparing the proteomes of known archaea and**
212 **bacteria as well as the twelve investigated asgardarchaea. Sequences were blasted against a**
213 **prokaryotic database comprising 5,655 prokaryotic genomes (including the 1000 bacteria**
214 **and 212 archaea which were used for protein family calculation) and a eukaryotic database**
215 **comprised of 150 genomes. For each organism we quantified the amount of sequences with**
216 **subjects among eukaryotes, bacteria or archaea only, or any combination of these three. To**
217 **counter possible biases in database composition, hits from the same genus were excluded. E,**
218 **Eukaryotes; A, Archaea; B, Bacteria; AB, Archaea-Bacteria; et cetera. (b) For proteins with-**
219 **or without database hits, protein length distributions are shown as box and whiskers plots.**
220 **This reveals a noticeable difference in protein lengths of sequences without database hits**
221 **from only the asgardarchaea compared to their sequences with database hits, and in**
222 **comparison, to those of Refseq archaea. p-values of FDR-corrected, pairwise double-sided**
223 **Kolmogorov-Smirnov tests in Suppl. Table 7. All used abbreviations are listed in Suppl. Table**
224 **6.**
225

226 We plotted the protein sequence length distributions for each proteome, separating hits with
227 and hits without significant database hits (Fig. 3b). While the length distributions of
228 sequences with significant database hits were comparable between the asgard- and the RefSeq
229 archaea, the distributions of hits without significant database hits show a major difference.
230 (Fig. 3b). Furthermore, simulating missing data by ignoring not only hits from the same genus
231 but also the same phylum produced a more similar result for the reference organisms and
232 asgardarchaea (Suppl. Fig.1), recapitulating the extensive diversity within the asgardarchaea.
233

234 **Discussion**

235

236 The identification of asgardarchaea from deep-sediment metagenome data (Seitz et al. 2016;
237 Zaremba-Niedzwiedzka et al. 2017) provides valuable new information from which to re-
238 evaluate key issues surrounding the tree of life and the emergence of its eukaryotic branch. To
239 begin with, the iconic tree that introduced three aboriginal lineages (Woese and Fox 1977)
240 might require a revision. Phylogenomic analysis of asgardarchaea provides evidence for a
241 two-domains tree and the emergence of the host cell lineage of eukaryogenesis from within
242 the archaeal domain (Cox et al. 2008, McInerney et al. 2014, Hug et al. 2016, Williams et al.
243 2020) with some skepticism, however, remaining (Forterre 2015, Da Cunha et al. 2018, Liu et
244 al. 2020). Parallel to the discovery of the asgardarchaea were immediate speculations
245 regarding their cellular complexity (Dacks et al. 2016, Pittis and Gabaldón 2016, Rout and
246 Field 2017, Akil and Robinson 2018, Zachar et al. 2018, Neveu et al. 2020) and a faith of
247 having identified the missing link between pro- and eukaryotic biology based on the
248 identification of a few eukaryote signature proteins (ESPs), that we find to be only 27 on

249 average. In light of these numbers, the potential of these archaea to display eukaryote-like cell
250 complexity is hard to maintain.

251

252 Our analysis confirms a patchy distribution of ESPs among asgardarchaea (Dacks et al. 2016,
253 Zaremba-Niedzwiedzka et al. 2017, Imachi et al. 2020, Liu et al. 2020, Klinger et al. 2016,
254 Inoue et al. 2020, Bulzu et al. 2019) (Fig. 1). While of course the archaeal host brought in
255 1000s of genes, the unique contribution of this new archaeal lineage to the eukaryotic protein
256 families is substantially less than what one might infer from the original metagenome reports
257 and subsequent interpretations (Dacks et al. 2016, Pittis and Gabaldón 2016, Rout and Field
258 2017, Zachar et al. 2018, Akil and Robinson 2018, Neveu et al. 2020, Vosseberg et al. 2020).
259 The irregular gene distribution (Fig. 2a) might reflect differential gene loss upon the
260 segregation of the common ancestor of asgardarchaea and the archaeal host cell lineage (Eme
261 et al. 2017). Considering the role of pangenomes in the transformation of prokaryotic lineages
262 and the conquering of ecological niches (McInerney et al. 2017), pangenomes offer a
263 complementary explanation to the differential loss of genes. The archaeal ancestor of
264 eukaryotes might have tapped a shared gene pool more extensively than the sister lineages
265 leading to extant asgardarchaea.

266

267 The notion that asgardarchaeal contributions to eukaryotes were higher only to be eventually
268 replaced by bacterial (endosymbiotic) contributions might be brought up (Pittis and Gabaldon
269 2016; Eme et al. 2017; Vosseberg et al. 2020). The absence of extant archaeal relatives with
270 similarly higher ESPs, however, indicates that archaea neither have the necessity nor the
271 selective pressure for maintaining ESPs in the absence of an endosymbiont. Any theory that
272 hinges upon a larger presence of ESPs in archaea or bacteria prior to mitochondrial
273 endosymbiosis ignores the complete lack of “accumulated ESPs” in extant prokaryotes to a
274 degree that even remotely matches that of any given eukaryotic lineage. Extinction and
275 absence of geological conditions that promoted eukaryote origins considered for ESPs, must
276 also be considered for any genes that are currently thought to be recent independent gene
277 transfer events from prokaryotes to eukaryotes and not of endosymbiotic (mitochondrial)
278 origin.

279

280 Our protein family clustering method, which readily detects the mitochondrial contribution
281 (and the cyanobacterial contributions in the case of the Archaeplastida; Suppl. Fig. 2) failed to
282 detect a comparable asgardarchaeal-unique contribution. A stacked bar diagram puts gene

283 family cluster contribution in each lineage into a global perspective (Fig. 3a; Supp. Fig 1).
284 There is a small proportion of eukaryotic homologs (E, mustard yellow) visible, e.g. in the
285 Candidatus *P. syntrophicum* MK-D1, but it is substantially smaller in comparison to the
286 eukaryote-bacteria (EB, green) specific homologs evident in eukaryotes (and bacteria *vice*
287 *versa*) that reflects the mitochondrial contribution to eukaryogenesis (Brueckner and Martin
288 2020). Neglecting the surprisingly low number of asgardarchaeal-unique homologs to
289 eukaryotic genomes, our analysis demonstrates that asgardarchaea are among the most
290 genetically diverse group of archaea when comparing it to the genus Methanococci and the
291 phylum Crenarchaeota (Fig. 2). The odd length distributions of the asgardarchaeal proteins
292 with no homology (Fig. 3b) are strange as well, since protein length across pro- and
293 eukaryotes is usually well conserved (Xu et al. 2006). This could hint at an assembly and/or
294 binning issue, which was also observed regarding the anomalous phylogenetic behavior of
295 their ribosomal proteins and concatenated gene trees (Da Cunha et al. 2018; Garg et al. 2021).
296 If not, it is a biological phenomenon absent in other sequenced prokaryotes.

297

298 Considering the amount of data gathered in just the few years (Imachi et al. 2020, Liu et al.
299 2020, Klinger et al. 2016, Inoue et al. 2020, Bulzu et al. 2019, Villanueva et al. 2016), it is
300 surprising asgardarchaea have escaped identification for so long. Their habitats had been
301 sampled before, so it is likely that the method used and maybe an obligate dependency on
302 syntrophy, hindered culturing except for one imposing exception (Imachi et al. 2020).
303 Dedicated phylogenomic efforts are necessary to resolve their taxonomic classification, while
304 only culturing can picture their cell morphology.

305

306 Analyses of asgardarchaeal ESPs shows they have the potential to function similar to their
307 eukaryotic homologs in a eukaryotic system (Neveu et al. 2020, Akıl and Robinson 2018,
308 Rout and Field 2017, Klinger et al. 2016), but cross-kingdom inferences have their limits
309 (Dey et al. 2016). The analysis of archaeal small GTPases (Surkont and Pereira-Leal 2016)
310 and homologs of ESCRT proteins, the CDVs (Lindås et al. 2008, Lindås and Bernander 2013,
311 Caspi and Dekker 2018), serve as examples. One needs to interpret asgardarchaeal ESPs in
312 their prokaryotic context and in cells lacking an endosymbiont. The first images of an
313 asgardarchaeon, those of Candidatus *P. syntrophicum* MK-D1, and its dependency on a
314 bacterial partner (Imachi et al. 2020) define the current standard from which to plot
315 eukaryogenesis and the steps leading to eukaryotic cell- and genome complexity.

316

317 The identification of the asgardarchaea and the culturing of one representative represent an
318 important milestone in micro- and evolutionary biology. Their phylogenetic analysis echoes
319 two previously predicted outcomes: (i) eukaryotes to branch from within archaea, solidifying
320 the two-domains tree of life, and (ii) that the closer we zoom in on the two prokaryotic
321 partners from which eukaryotes evolved, the higher the number of otherwise eukaryote-
322 typical genes we identify in prokaryotes. The description of Candidatus *P. syntrophicum* MK-
323 D1 (Imachi et al. 2020) reminds us to not conflate genotypic with phenotypic complexity.
324 This predicts that future asgardarchaea we see cultured will lack eukaryotic traits, too, and
325 most, if not all, will depend on syntrophy. Our analysis also predicts that much of
326 asgardarchaeal diversity remains to be described, but that the gap between the pro- and
327 eukaryotic protein families will remain decisive and to change little. Placing the
328 endosymbiotic event and the energetic benefit offered by mitochondria to fuel the transition
329 early in eukaryogenesis, explains the lack of physical evidence for eukaryote-like complexity
330 in asgardarchaea despite them encoding ESPs. It offers a comprehensive full-service theory
331 for the singularity that is the origin of the eukaryotic cell that mitochondria-late models fail to
332 provide.

333

334 **Acknowledgements:** SBG would like to thank the German research council (267205415 –
335 SFB 1208) and the VolkswagenStiftung (Life) for funding. SGG was furthermore supported
336 by the Moore–Simons Project on the Origin of the Eukaryotic Cell GBMF9743.
337 Computational infrastructure and support were provided by the Centre for Information and
338 Media Technology at HHU Düsseldorf.

339

340 **Author Contributions:** SGG, MK and SBG conceived the analysis, which was carried out
341 predominantly by MK but also SS. SGG, MK, MvdG and SBG wrote the manuscript, whose
342 final version was approved all authors.

343

344 **Methods**

345 Calculation of protein families

346 Prokaryotic gene families were calculated from complete genomes of 1000 bacteria and 226
347 archaea of the Refseq database (Version September 2016) including eleven representatives of
348 the asgardarchaea¹³ and Candidatus *P. syntrophicum* MK-D1 (Imachi et al. 2020), separately.
349 Bacterial and archaeal protein families were calculated via MCL (van Dongen 2000, Enright
350 2002) (--abc -scheme 7) from all reciprocal best BLAST hits with pairwise global identities
351 (Rice et al. 2000) of at least 25% identity and a maximum e-value of 1×10^{-10} among all
352 investigated bacteria and archaea, respectively. Eukaryotic gene families on the basis of 150
353 eukaryotic genomes were calculated as part of a previous study (Brueckner et al. 2020).
354 Prokaryotic and eukaryotic cluster were combined into EPCs if at least 50% of all sequences
355 of a eukaryotic cluster had their best hit in a prokaryotic cluster and vice versa according to
356 the ‘reciprocal best cluster approach’ described in Ku et al. 2015 (Ku et al. 2015). All
357 proteomes used in this study are listed in Supplementary Table 5. Protein families of
358 Methanococci and Crenarchaeota were calculated from all pairwise global identities of all
359 protein sequences, using the above-mentioned identity and e-value cutoffs, via MCL (van
360 Dongen 2000, Enright 2002).

361

362 BLAST hit analysis

363 We performed a Diamond BLASTp (Buchfink et al. 2015) hit analysis on all protein
364 sequences of the eleven asgardarchaeal proteomes including MK-D1. We compared the
365 results to twelve bacterial and archaeal proteomes each, plus eight eukaryotic ones. We
366 blasted all protein sequences of the chosen proteomes against our database of 5655
367 prokaryotic and 150 eukaryotic proteomes. For each proteome, we quantified the number of
368 sequences that showed significant hits (at least 25% identity and a minimal e-value of 1×10^{-5})
369 within bacteria, archaea, eukaryota or any combination of these three. To counter
370 overrepresentation of some genera within our database, we excluded hits of the same genus
371 for all tested protein sequences.

372

373 InterProScan ESP analysis

374 InterProScan version 5.39-77.0 (Quevillon et al. 2005) with standard parameters was used to
375 annotate all asgardarchaeal proteomes, the MK-D1 proteome and all 212 archaeal proteomes
376 within our prokaryotic database. As in Zaremba et al. 2017 (Zaremba-Niedzwiedzka et al.
377 2017), we searched for InterPro-Identifiers that correspond with Eukaryote-specific-proteins

378 and plotted the results of all investigated asgard archaea together with the results of 14 model
379 Refseq archaea for comparison. (see Suppl. Figure 3)

380

381 LECA cluster filtering

382 Since eukaryotic inheritance is strictly vertical the eukaryotic protein families were filtered
383 for families that contained at least one protein sequence from one member of each of the six
384 supergroups included in our dataset of 150 eukaryotes, being *Archaeplastida*, *Opisthokonta*,
385 *SAR*, *Hacrobia*, *Excavata* and *Mycetozoa* resulting in 1880 protein families passing this
386 criterion.

387

388 **Suppl. material titles and legends**

389

390 **Suppl. Fig. 1: BLAST hit analysis.** The BLAST hit analysis underlying Fig. 3a was redone,
391 this time ignoring all hits from the same phylum. This was done in order to test for the
392 asgardarchaeas' taxonomic rank, since the amount of protein sequences without any
393 significant database hits, could either stem from poor sequence quality or just from the fact
394 that asgardarchaea are inherently different from all archaea we know so far. As expected, the
395 amount of Refseq archaeal protein sequences without significant database hits did increase
396 but the effect was not sufficient to recover the same proportion of “unknown” protein
397 sequences as we see in asgardarchaea.

398

399 **Suppl. Fig. 2: PAPs of sets of Eukaryote-Prokaryote-Clusters.** (a) Presence-Absence
400 pattern of all 537 EPCs containing members from eukaryotes, bacteria and archaea. (b)
401 Presence-Absence pattern of all 1723 EPCs containing only proteins from eukaryotes and
402 bacteria. The cyanobacterial contribution to the Archaeplastida is marked with red boxes. In
403 both cases, protein families were sorted by their distribution among the six eukaryotic
404 supergroups Archaeplastida, Excavata, Hacrobia, Mycetozoa, Opisthokonta and the SAR
405 group.

406

407 **Suppl. Fig. 3: InterProScan results for all investigated asgardarchaea as well as 14**
408 **reference archaea.** InterProScan was used to annotate the proteomes of all 12 investigated
409 asgardarchaea and 14 reference archaea. InterProScan results were used to search for ESP
410 candidate proteins within all analyzed proteomes and their presence indicated by a black dot.
411 InterPro identifiers pointing towards the presence of ESPs were obtained from Zaremba et al.

412 2017 and in some cases revised to be more strict. For most ESPs results show a clear divide
413 between asgardarchaea and the chosen reference archaea, while the DNA-directed RNA-
414 Polymerases I/III and II, Cyclins or Cyclin-like proteins and the STT3 subunit of the OST
415 complex could also be detected in the majority of the analyzed reference archaea.

416

417 **Suppl. Table 1.** EPC overview including protein family counts, protein sequence listing for
418 each protein family and relative frequency of occurrence within protein families of all
419 taxonomic groups present in the clustering. One protein family per line, one member each
420 cell.

421

422 **Suppl. Table 2.** Six identified Eukaryote-Asgardarchaea unique protein families annotated
423 via InterProScan using standard parameters and showing results from all subject databases.

424

425 **Suppl. Table 3.** Archaeal and bacterial clusterings. One protein family per line, one member
426 each cell.

427

428 **Suppl. Table 4.** InterProScan results for all investigated Asgardarchaea, underlying Fig. 3. in
429 addition to the list of InterPro accessions used to identify ESPs within the investigated
430 asgardarchaeal MAGs.

431

432 **Suppl. Table 5.** Listing of all proteome assemblies used in cluster creation and BLAST hit
433 analysis.

434

435 **Suppl. Table 6.** Species abbreviations used in Figure 2 and Figure 3.

436

437 **Suppl. Table 7.** FDR-corrected (Benjamini-Hochberg) p-values of pairwise double-sided KS
438 tests of gene length distributions.

439

440 **References**

441 Akil C, & Robinson RC (2018). Genomes of Asgard archaea encode profilins that regulate
442 actin. *Nature*, **562**, 439–443.

443 Blackstone NW (2013). Why did eukaryotes evolve only once? Genetic and energetic aspects
444 of conflict and conflict mediation. *Phil trans R Soc B*, **368**, 20120266.

- 445 Booth A & Doolittle WF (2015). Eukaryogenesis, how special really? *Proc Natl Acad Sci*
446 *USA*, **112**, 10278–10285.
- 447 Brueckner J, & Martin WF (2020). Bacterial genes outnumber archaeal genes in eukaryotic
448 genomes *Genome Biol Evol*, **12**, 282-292.
- 449 Buchfink B, Xie C, & Huson DH (2015). Fast and sensitive protein alignment using
450 DIAMOND. *Nature Methods*, **12**, 59–60.
- 451 Bulzu P-A, Andrei A-S, Salcher MM, Mehrshad M, Inoue K, Kandori H, Beja O, Ghai R, &
452 Banciu HL (2019). Casting light on Asgardarchaeota metabolism in a sunlit microoxic niche.
453 *Nat Microbiol*, **4**, 1129–1137.
- 454 Burns JA, Pittis AA & Kim E (2018). Gene-based predictive models of trophic modes suggest
455 Asgard archaea are not phagocytotic. *Nat Ecol Evol*, **2**, 697–704.
- 456 Caspi Y & Dekker C (2018). Dividing the archaeal way: the ancient cdv cell-division
457 machinery. *Front Microbiol*, **9**, 174.
- 458 Cavalier-Smith T (1987) The origin of eukaryote and archaebacterial cells. *Ann N Y Acad Sci*,
459 **503**, 17–54.
- 460 Cox CJ, Foster PG, Hirt RP, Harris SR &, Embley TM (2008). The archaebacterial origin of
461 eukaryotes. *Proc Natl ional Acad Sci USA*, **105**, 20356–20361.
- 462 Da Cunha V, Gaia M, Nasir A &, Forterre P (2018). Asgard archaea do not close the debate
463 about the universal tree of life topology. *PLoS Genet*, **14**, e1007215.
- 464 Dacks JB, Field MC, Buick R, Eme L, Gribaldo S, Roger AJ, Brochier-Armanet C &, Devos
465 DP. (2016). The changing view of eukaryogenesis – fossils, cells, lineages and how they all
466 come together. *Journal of Cell Science Sci*, **129**, 3695-3703.
- 467 Dey G, Thattai M & Baum B (2016). On the archaeal origins of eukaryotes and the challenges
468 of inferring phenotype from genotype. *Trends Cell Biol.*, **26**, 476–485.
- 469 Duplessis MR, Ziebis W, Gros O, Caro A, Robidart J &, Felbeck H. 2004. Respiration
470 strategies utilized by the gill endosymbiont from the host lucinid *Codakia orbicularis*
471 (*Bivalvia*: *Lucinidae*). *Appl Environ Microbiol*, **70**, 4144–4150.
- 472 Embley TM & Finlay BJ. 1993. Systematic and morphological diversity of endosymbiotic
473 methanogens in anaerobic ciliates. *Antonie Van Leeuwenhoek*, **64**, 261–271.
- 474 Eme L, Spang A, Lombard J, Stairs CW & Ettema TJG. 2017. Archaea and the origin of
475 eukaryotes. *Nat Rev Microbiol*, **15**, 711–723.
- 476 Enright, AJ. 2002. An efficient algorithm for large-scale detection of protein families. *Nucleic*
477 *Acids Res*, **30**, 1575–1584.

- 478 Fenchel T & Bernard C. 1993. A purple protist. *Nature*, **362**, 300–300.
- 479 Forterre P. 2015. The universal tree of life: an update. *Front Microbiol*, **6**, 717.
- 480 Fritz-Laylin LK, Prochnik SE, Ginger ML, Dacks JB, Carpenter ML, et al..., Dawson SC.
481 (2010). The genome of *Naegleria gruberi* illuminates early eukaryotic versatility. *Cell*. **140**,
482 631–642.
- 483 Garg SG, Kapust N, Lin W, Knopp M, Tria F, Nelson-Sathi S, Gould SB, Fan L, Zhu R,
484 Zhang C, Martin WF. 2021. Anomalous phylogenetic behavior of ribosomal proteins in
485 metagenome assembled genomes. *Genome Biol Evol* , **13**, evaa238.
- 486 Garg SG & Martin WF. 2016. Mitochondria, the cell cycle, and the origin of sex via a
487 syncytial eukaryote common ancestor. *Genome Biol Evol.*, **8**, 1950-1970.
- 488 Gould SB, Garg SG. & Martin WF. 2016. Bacterial vesicle secretion and the evolutionary
489 origin of the eukaryotic endomembrane system. *Trends Microbiol.*, **24**, 525–534.
- 490 Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, ..., Banfield JF. 2016. A new
491 view of the tree of life. *Nat. Microbiol.*, **1**, 16048.
- 492 Husnik F, Nikoh N, Koga R, Ross L, Duncan RP, ..., McCutcheon JP. 2013. Horizontal gene
493 transfer from diverse bacteria to an insect genome enables a tripartite nested mealybug
494 symbiosis. *Cell*, **153**, 1567.
- 495 Imachi H, Nobu MK, Nakahara N, Morono Y, Ogawara M, Takaki Y, Takano Y, Uematsu K,
496 Ikuta T, Ito M, Matsui Y, Miyazaki M, Murata K, Saito Y, Sakai S, Song C, Tasumi E,
497 Yamanaka Y, Yamaguchi T, ..., Takai K. 2020. Isolation of an archaeon at the prokaryote–
498 eukaryote interface. *Nature*, **577**, 519–525.
- 499 Inoue K, Tsunoda SP, Singh M, Tomida S, Hososhima S, ..., Kandori H. 2020.
500 Schizorhodopsins: A family of rhodopsins from Asgard archaea that function as light-driven
501 inward H⁺ pumps. *Sci Adv*, **6**, eaaz2441.
- 502 Klinger CM, Spang A, Dacks JB & Ettema TJG. 2016. Tracing the archaeal origins of
503 eukaryotic membrane-trafficking system building blocks. *Mol Biol Evol*, *Mol. Biol. Evol.*, **33**,
504 1528–1541.
- 505 Koonin EV, Csuros M & Rogozin IB. 2013. Whence genes in pieces: reconstruction of the
506 exon-intron gene structures of the last eukaryotic common ancestor and other ancestral
507 eukaryotes. *Wiley Interdiscip Rev RNA*, **4**, 93–105.
- 508 Koumandou VL, Wickstead B, Ginger ML, van der Giezen M, Dacks JB & Field MC. 2013.
509 Molecular paleontology and complexity in the last eukaryotic common ancestor. *Crit Rev*
510 *Biochem Mol* , **48**, 373–396.

- 511 Ku C, Nelson-Sathi S, Roettger M, Sousa FL, Lockhart PJ, Bryant D, Hazkani-covo E,
512 McInerney JO, Landan G, & Martin WF. 2015. Endosymbiotic origin and differential loss of
513 eukaryotic genes. *Nature*, **524**, 427–437.
- 514 Lane N & Martin WF. 2010. The energetics of genome complexity. *Nature*, **467**, 929–934.
- 515 Lane N & Martin WF. 2016. Mitochondria, complexity, and evolutionary deficit spending.
516 *Proc National Acad Sci*, **113**, E666–E666.
- 517 Lane N. 2011. Energetics and genetics across the prokaryote-eukaryote divide. *Biol. Direct*, **6**,
518 35.
- 519 Lindås A-C & Bernander R. 2013. The cell cycle of archaea. *Nat Rev Microbiol*, **11**, 627–38.
- 520 Lindås A-C, Karlsson EA, Lindgren MT, Ettema TJG & Bernander R. 2008. A unique cell
521 division machinery in the Archaea. *Proc National Acad Sci* **105**, 18942–18946.
- 522 Liu Y, Makarova KS, Huang WC, Wolf YI, Nikolskaya A, Zhang X, Cai M, Zhang CJ, Xu
523 W, Luo Z, Cheng L, Koonin EV & Li M. 2020. Expanding diversity of Asgard archaea and
524 the elusive ancestry of eukaryotes. bioRxiv
- 525 Lynch M & Marinov GK. 2015. The bioenergetic costs of a gene. *Proc Natl Acad Sci USA*,
526 **112**, 15690–15695.
- 527 Martijn J & Ettema TJ. 2013. From archaeon to eukaryote: the evolutionary dark ages of the
528 eukaryotic cell. *Biochem Soc Trans*, **41**, 451-457.
- 529 Martin W & Müller M. 1998. The hydrogen hypothesis for the first eukaryote. *Nature*, **392**,
530 37–41.
- 531 Martin WF, Garg SG & Zimorski V. 2015. Endosymbiotic theories for eukaryote origin. *Phil*
532 *Trans R Soc B*, **370**, 20140330.
- 533 Martin WF, Tielens AGM, Mentel M, Garg SG & Gould SB. 2017. The physiology of
534 phagocytosis in the context of mitochondrial origin. *Microbiol Mol Biol Rev*, **81**, 1–36.
- 535 McInerney JO, O’Connell MJ & Pisani D. 2014. The hybrid nature of the Eukaryota and a
536 consilient view of life on Earth. *Nat Rev Microbiol*, **12**, 449–455.
- 537 McInerney J, McNally A & O’Connell M. 2017. Why prokaryotes have pangenomes. *Nat*
538 *Microbiol*, **2**, 17040.
- 539 Mills DB. 2020. The origin of phagocytosis in earth history. *Interface Focus* **10**, 20200019.
- 540 Neveu E, Khalifeh D, Salamin N, & Fasshauer D. 2020. Prototypic SNARE proteins are
541 encoded in the genomes of Heimdallarchaeota, potentially bridging the gap between the
542 prokaryotes and eukaryotes. *Curr Biol*, **30**, 2468-2480.

- 543 O'Malley MA. 2010. The first eukaryote cell: an unfinished history of contestation. *Stud Hist*
544 *Philos Biol Biomed Sci*, **41**, 212–224.
- 545 Pittis AA & Gabaldón T. 2016. Late acquisition of mitochondria by a host with chimaeric
546 prokaryotic ancestry. *Nature*, **531**, 101–104.
- 547 Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, & Lopez R. 2005.
548 InterProScan: protein domains identifier. *Nucleic Acids Research*, **33**(Web Server), W116–
549 W120.
- 550 Rice P, Longden L, & Bleasby A. 2000. EMBOSS: The European Molecular Biology Open
551 Software Suite. *Trends Genet*, **16**, 276–277.
- 552 Rout MP & Field MC. 2017. The evolution of organellar coat complexes and organization of
553 the eukaryotic cell. *Annu Rev Biochem*, **86**, 637–657.
- 554 Schmid AMM. 2003. Endobacteria in the diatom *Pinnularia* (Bacillariophyceae). I. ‘Scattered
555 ct-nucleoids’ explained: DAPI-DNA complexes stem from exoplastidial bacteria boring into
556 the chloroplasts. *J Phycol*, **39**, 122–138.
- 557 Seitz KW, Lazar CS, Hinrichs K-U, Teske AP & Baker BJ. 2016. Genomic reconstruction of
558 a novel, deeply branched sediment archaeal phylum with pathways for acetogenesis and
559 sulfur reduction. *ISME J* **10**, 1696–1705.
- 560 Shiratori T, Suzuki S, Kakizawa Y & Ishida K-I. 2019. Phagocytosis-like cell engulfment by
561 a planctomycete bacterium. *Nat Commun*, **10**, 5529.
- 562 Spang A, Saw JH, Jørgensen SL, Zaremba-Niedzwiedzka K, Martijn J, ..., Ettema TJG.
563 2015. Complex archaea that bridge the gap between prokaryotes and eukaryotes. *Nature*, **521**,
564 173–179.
- 565 Spang A, Stairs CW, Dombrowski N, Eme L, Lombard J, ..., Ettema TJG. 2019. Proposal of
566 the reverse flow model for the origin of the eukaryotic cell based on comparative analyses of
567 Asgard archaeal metabolism. *Nat Microbiol*, **4**, 1138–1148.
- 568 Surkont J & Pereira-Leal JB. 2016. Are there Rab GTPases in archaea? *Mol Biol Evol*, **33**,
569 1833–1842.
- 570 Thacker RW. 2005. Impacts of shading on sponge-cyanobacteria symbioses: A comparison
571 between host-specific and generalist associations. *Integr Comp Biol*, **45**, 369–376.
- 572 Tria FDK, Brückner J, Skejo J, Xavier JC, Zimorski V, Gould SB., Garg SG, & Martin WF.
573 2019. Gene duplications trace mitochondria to the onset of eukaryote complexity. *BioRxiv*,
574 **781211**. <https://doi.org/10.1101/781211>
- 575 van Dongen S. 2000. A cluster algorithm for graphs. *Inf Syst*, **R0010**, 1–40.

- 576 Vellai T, Takács K & Vida G. 1998. A new aspect to the origin and evolution of eukaryotes.
577 *J Mol Evol*, **46**, 499-507.
- 578 Villanueva L, Schouten S & Damsté JSS. 2016. Phylogenomic analysis of lipid biosynthetic
579 genes of Archaea shed light on the “lipid divide”. *Environ Microbiol*, **19**, 54–69.
- 580 Vosseberg J, van Hooff JJE, Marcet-Houben M, van Vlimmeren A, van Wijk LM, Gabaldón
581 T & Snel B. 2020. Timing the origin of eukaryotic cellular complexity with ancient
582 duplications. *Nat Ecol Evol*, **5**, 92–100.
- 583 Williams TA, Cox CJ, Foster PG, Szöllösi GJ & Embley TM. 2020. Phylogenomics provides
584 robust support for a two-domains tree of life. *Nat Ecol Evol*, **4**, 138–147.
- 585 Woese CR. & Fox GE. 1977. Phylogenetic structure of the prokaryotic domain: The primary
586 kingdoms. *Proc Natl Acad Sci USA*, **74**, 5088–5090.
- 587 Woyke T, Teeling H, Ivanova NN, Huntemann M, Richter M, ..., Dubilier N. 2006.
588 Symbiosis insights through metagenomic analysis of a microbial consortium. *Nature*, **443**,
589 950–955.
- 590 Xu L, Chen H, Hu X, Zhang R, Zhang Z, & Luo ZW. 2006. Average gene length is highly
591 conserved in prokaryotes and eukaryotes and diverges only between the two kingdoms. *Mol*
592 *Biol Evol*, **23**, 1107–1108.
- 593 Yutin N, Wolf MY, Wolf YI & Koonin EV. 2009. The origins of phagocytosis and
594 eukaryogenesis. *Biol Direct*, **4**, 9.
- 595 Zachar I, Szilágyi A, Számadó S & Szathmáry E. 2018. Farming the mitochondrial ancestor
596 as a model of endosymbiotic establishment by natural selection. *Proc Nat Acad Sci USA*, **115**,
597 201718707.
- 598 Zaremba-Niedzwiedzka K, Caceres EF, Saw JH, Bäckström D, Juzokaite L, ..., & Ettema
599 TJG. 2017. Asgard archaea illuminate the origin of eukaryotic cellular complexity. *Nature*,
600 **541**, 353–358.
- 601 Zientz E, Dandekar T & Gross R. 2004. Metabolic Interdependence of Obligate Intracellular
602 Bacteria and Their Insect Hosts. *Microbiol Mol Biol Rev*, **68**(4), 745–770.