

1 **Expanding diversity of Asgard archaea and the elusive ancestry of eukaryotes**

2

3 Yang Liu^{1†}, Kira S. Makarova^{2†}, Wen-Cong Huang^{1†}, Yuri I. Wolf², Anastasia Nikolskaya², Xinxu
4 Zhang¹, Mingwei Cai¹, Cui-Jing Zhang¹, Wei Xu³, Zhuhua Luo³, Lei Cheng⁴, Eugene V. Koonin^{2*}, Meng
5 Li^{1*}

6 1 Shenzhen Key Laboratory of Marine Microbiome Engineering, Institute for Advanced Study, Shenzhen
7 University, Shenzhen, Guangdong, 518060, P. R. China

8 2 National Center for Biotechnology Information, National Library of Medicine, National Institutes of
9 Health, Bethesda, Maryland 20894, USA

10 3 State Key Laboratory Breeding Base of Marine Genetic Resources, Key Laboratory of Marine Genetic
11 Resources, Fujian Key Laboratory of Marine Genetic Resources, Third Institute of Oceanography, State
12 Oceanic Administration, Xiamen 361005, P. R. China

13 4 Key Laboratory of Development and Application of Rural Renewable Energy, Biogas Institute of
14 Ministry of Agriculture, Chengdu 610041, P.R. China

15 † These authors contributed equally to this work.

16 *Authors for correspondence: koonin@ncbi.nlm.nih.gov or limeng848@szu.edu.cn

17

18

19 [Running title: Asgard archaea genomics](#)

20 [Keywords:](#)

21 **Abstract**

22 Comparative analysis of 162 (nearly) complete genomes of Asgard archaea, including 75 not reported
23 previously, substantially expands the phylogenetic and metabolic diversity of the Asgard superphylum,
24 with six additional phyla proposed. Phylogenetic analysis does not strongly support origin of eukaryotes
25 from within Asgard, leaning instead towards a three-domain topology, with eukaryotes branching outside
26 archaea. Comprehensive protein domain analysis in the 162 Asgard genomes results in a major expansion
27 of the set of eukaryote signature proteins (ESPs). The Asgard ESPs show variable phyletic distributions
28 and domain architectures, suggestive of dynamic evolution via horizontal gene transfer (HGT), gene loss,
29 gene duplication and domain shuffling. The results appear best compatible with the origin of the
30 conserved core of eukaryote genes from an unknown ancestral lineage deep within or outside the extant
31 archaeal diversity. Such hypothetical ancestors would accumulate components of the mobile archaeal
32 ‘eukaryome’ via extensive HGT, eventually, giving rise to eukaryote-like cells.

33 **Introduction**

34 The Asgard archaea are a recently discovered archaeal superphylum that is rapidly expanding, thanks to
35 metagenomic sequencing (1–5). The Asgard genomes encode a diverse repertoire of eukaryotic signature
36 proteins (ESPs) that far exceeds the diversity of ESPs in other archaea. The Asgard ESPs are particularly
37 enriched in proteins involved in membrane trafficking, vesicle formation and transport, cytoskeleton
38 formation and the ubiquitin network, suggesting that these archaea possess a eukaryote-type cytoskeleton
39 and an intracellular membrane system (2).

40 The discovery of the Asgard archaea rekindled the decades old but still unresolved fundamental debate on
41 the evolutionary relationship between eukaryotes and archaea that has shaped around the ‘2-domain (2D)
42 versus 3-domain (3D) tree of life’ theme (6–8). The central question is whether the eukaryotic nuclear
43 lineage evolved from a common ancestor shared with archaea, as in the 3D tree, or from within the
44 archaea, as in the 2D tree. The discovery and phylogenomic analysis of Asgard archaea yielded strong
45 evidence in support of the 2D tree, in which eukaryotes appeared to share common ancestry with one of
46 the Asgard lineages, Heimdallarchaeota (1, 2, 5). However, the debate is not over as arguments have been
47 made for the 3D topology, in particular, based on the phylogenetic analysis of RNA polymerases, some of
48 the most highly conserved, universal proteins (9, 10).

49 Molecular phylogenetic methods alone might be insufficient to resolve the ancient ancestral relationship
50 between archaea and eukarya. To arrive at a compelling solution, supporting biological evidence is crucial
51 (11), because, for example, the transition from archaeal, ether-linked membrane lipids to eukaryotic,
52 ester-linked lipids that constitute eukaryotic (and bacterial) membranes (12) and the apparent lack of the
53 phagocytosis capacity in Asgard archaea (13) are major problems for scenarios of the origin of eukaryotes
54 from Asgard or any other archaeal lineage. Some biochemical evidence indicates that Asgard archaea
55 possess an actin cytoskeleton regulated by accessory proteins, such as profilins and gelsolins, and the
56 endosomal sorting complex required for transport machinery (ESCRT) that can be predicted to function
57 similarly to the eukaryotic counterparts (14–16). Generally, however, the biology of Asgard archaea
58 remains poorly characterized, in large part, because of their recalcitrance to growth in culture (17). To
59 date, only one Asgard archaeon, *Candidatus* Prometheoarchaeum syntrophicum strain MK-D1, has been
60 isolated and grown in culture (17). This organism has been reported to form extracellular protrusions that
61 are involved in its interaction with syntrophic bacteria, but no visible organelle-like structure and,
62 apparently, little intracellular complexity.

63 The only complete, closed genome of an Asgard archaeon also comes from *Candidatus* P. syntrophicum
64 strain MK-D1 (17) whereas all other genome sequences were obtained by binning multiple metagenomics

65 contigs. Furthermore, these genome sequences represent but a small fraction of the Asgard diversity that
66 has been revealed by 16S rRNA sequencing (18, 19). An additional challenge to the study of the
67 relationship between archaea and eukaryotes is that identification and analysis of the archaeal ESPs are
68 non-trivial tasks due to the high sequence divergence of many if not most of these proteins. At present,
69 the most efficient, realistic approach to the study of Asgard archaea and their eukaryotic connections
70 involves obtaining high quality genome sequences and analyzing them using the most powerful and
71 robust of the available computational methods.

72 Here we describe metagenomic mining of the expanding diversity of the superphylum Asgard, including
73 the identification of six additional phylum-level lineages that thrive in a wide variety of ecosystems and
74 are inferred to possess versatile metabolic capacities. We show that these uncultivated Asgard groups
75 carry a broad repertoire of ESPs many of which have not been reported previously. Our in depth
76 phylogenomic analysis of these genomes provides insights into the evolution of Asgard archaea but calls
77 into question the origin of eukaryotes from within Asgard.

78

79 **Results**

80 **Reconstruction of Asgard archaeal genomes from metagenomics data**

81 We reconstructed 75 metagenome-assembled genomes (MAGs) from the Asgard superphylum that have
82 not been reported previously. These MAGs were recovered from various water depths of the Yap trench,
83 intertidal mangrove sediments of Mai Po Nature Reserve (Hong Kong, China) and Futian Mangrove
84 Nature Reserve (Shenzhen, China), seagrass sediments of Swan Lake Nature Reserve (Rongcheng,
85 China) and petroleum samples of Shengli oilfield (Shandong, China) (Supplementary Table 1,
86 Supplementary Figure 1). For all analyses described here, these 75 MAGs were combined with 87
87 publicly available genomes, resulting in a set of 162 Asgard genomes. The 75 genomes reconstructed here
88 were, on average, 82% complete and showed evidence of low contamination of about 3%, on average
89 (Supplementary Figure 2).

90

91 **Classification of Asgard genes into clusters of orthologs**

92 The previous analyses of Asgard genomes detected a large fraction of “dark matter” genes (20). For
93 example, in the recently published complete genome of *Candidatus Prometheoarchaeum syntrophicum*,
94 45% of the proteins are annotated as “hypothetical”. We made an effort to improve the annotation of
95 Asgard genomes by investigating this dark matter in greater depth, and developing a dedicated platform

96 for Asgard comparative genomics. To this end, we constructed Asgard Clusters of Orthologous Genes
97 (asCOGs) and used the most sensitive available methods of sequence analysis to annotate additional
98 Asgard proteins, attempting, in particular, to expand the catalogue of Asgard homologs of ESPs (see
99 Materials and Methods for details).

100 Preliminary clustering by sequence similarity and analysis of the protein cluster representation across the
101 genomes identified the set of 76 most complete Asgard MAGs (46 genomes available previously and 30
102 ones reported here) that cover most of the group diversity (Supplementary Table 1). The first version of
103 the asCOGs presented here consists of 14,704 orthologous protein families built for this 76-genome set.
104 The asCOGs cover from 72% to 98% (92% on average) of the proteins in these 76 genomes (additional
105 data file 1). Many asCOGs include individual domains of large, multidomain proteins.

106 The gene commonality plot for the asCOGs shows an abrupt drop at the right end, which reflects a
107 surprising deficit of nearly universal genes (Fig. 1). Such shape of the gene commonality curve appears
108 anomalous compared to other major groups of archaea or bacteria with many sequenced genomes (21).
109 For example, in the case of the TACK superphylum of archaea, for which the number of genomes
110 available is similar to that for Asgard, with a comparable level of diversity, the commonality plot shows
111 no drop at the right end, but instead, presents a clear uptick, which corresponds to the core of genes
112 represented in (almost) all genomes (Fig. 1). Apparently, most of the Asgard genomes remain incomplete,
113 such that conserved genes were missed randomly. Currently, there are only three gene families that are
114 present in all Asgard MAGs, namely, a Zn-ribbon domain, a Threonyl-tRNA synthetase and an
115 aminotransferase (additional data file 1).

116 We employed the asCOG profiles to annotate the remaining 86 Asgard MAGs, including those that were
117 sequenced in the later stages of this work (Supplementary Table 1). On average, 89% of the proteins
118 encoded in these genomes were covered by asCOGs (Supplementary Table 1). Thus, the asCOGs
119 database appears to be an efficient tool for annotation and comparative genomic analysis of Asgard
120 MAGs and complete genomes.

121

122 **Expanding the phylogenetic diversity of Asgard archaea**

123 Phylogenetic analysis of the Asgard MAGs based on a concatenated alignment of 209 core asCOGs (see
124 Methods and additional data file 2) placed many of the genomes reported here into the previously
125 delineated major Asgard lineages (Fig. 2a, Supplementary Table 1, additional data file 2), namely,
126 Thorarchaeota (n=20), Lokiarchaeota (n=18), Hermodarchaeota (n=9), Gerdarchaeota (n=3),

127 Helarchaeota (n=2), and Odinarchaeota (n=1). Additionally, we identified 6 previously unknown major
128 Asgard lineages that appear to be strong candidates to become additional phyla (Fig. 2a and b,
129 Supplementary Table 1; see also Taxonomic description of new taxa in the Supplementary Information
130 and additional data file 1). A clade formed by As_085 and As_075 is a deeply branching sister group to
131 the previously recognized Heimdallarchaeota(2). Furthermore, our phylogenetic analysis supported the
132 further split of “Heimdallarchaeota” into 4 phylum-level lineages according to the branch length in the
133 concatenated phylogeny (see Materials and Methods; see also Taxonomic description of new taxa in the
134 Supplementary Information and additional data file 1). The putative phyla within the old
135 Heimdallarchaeota included the previously defined Gerdarchaeota (4), and three additional phyla that
136 could be represented by As_002 (LC2) , As_003 (LC3) and AB_125 (As_001), respectively. Another 3
137 previously undescribed lineages were related, respectively, to Hel-, Loki-, Odin- and Thorarchaeota.
138 Specifically, As_181, As_178 and As_183 formed a clade that was deeply rooted at the Hel-Loki-Odin-
139 Thor clade; As_129 and As_130 formed a sister group to Odinarchaeota; and a lineage represented by
140 As_086 was a sister group to Thorarchaeota. These results were buttressed by the 16S rRNA gene
141 phylogeny, comparisons of the mean amino acid identity and 16S rRNA sequence identity (Fig. 2b,
142 Supplementary Figure 3, Supplementary Figure 4, Supplementary Table 2 and Supplementary Table 3).

143 We propose the name Wukongarchaeota after Wukong, a Chinese legendary figure who caused havoc in
144 the heavenly palace, for the putative phylum represented by MAGs As_085 and As_075 (*Candidatus*
145 Wukongarchaeum yapensis), and names of Asgard deities in the Norse mythology for the other 5
146 proposed phyla: (1) Hodarchaeota, after Hod, the god of darkness, for MAG As_027 (*Candidatus*
147 Hodarchaeum mangrove); (2) Kariarchaeota, after Kari, the god of the North wind, for MAG As_030
148 (*Candidatus* Kariarchaeum pelagius); (3) Borrarchaeota after Borr, the creator god and father of Odin, for
149 MAG As_133 (*Candidatus* Borrarchaeum yapensis); (4) Baldrarchaeota, after Baldr, the god of light and
150 brother of Thor, for MAG As_130 (*Candidatus* Baldrarchaeum yapensis); (5) and Hermodarchaeota after
151 Hermod, the messenger of the gods, son of Odin and brother of Baldr, for MAG As_086 (*Candidatus*
152 Hermodarchaeum yapensis) (Fig. 2a). For details, see Taxonomic Description of new taxa in the
153 **Supplementary Information.**

154 The gene content of Asgard MAGs agrees well with the phylogenetic structure of the group. The phyletic
155 patterns of the asCOG form clusters that generally correspond to the clades identified by phylogenetic
156 analysis (Fig. 3a), suggesting that gene gain and loss within Asgard archaea largely proceeded in a clock-
157 like manner and/or that horizontal gene exchange preferentially occurred between genomes within the
158 same clade.

159

160 **Phylogenomic analysis and positions of Asgard archaea and eukaryotes in the tree of life**

161 The outcome of phylogenetic reconstruction, especially, when deep branchings are involved, such as
162 those that are relevant for the 3D vs 2D conundrum, depends on the phylogenetic methods employed, the
163 selection of genes for phylogeny construction and, perhaps most dramatically, on the species sampling
164 (22–24). In many cases, initially uncertain positions of lineages in a tree settle over time once more
165 representatives of the groups in question and their relatives become available.

166 In our analysis of the universal phylogeny, we aimed to make the species set for phylogenetic
167 reconstruction as broadly representative as possible, while keeping its size manageable, to allow the use
168 of powerful phylogenetic methods. The tree was constructed from alignments of conserved proteins of
169 162 Asgard archaea, 286 other archaea, 98 bacteria and 72 eukaryotes (see Supplementary Material and
170 Methods for details of the procedure including the selection of a representative species set and
171 Supplementary Table 4). Members of 30 families of (nearly) universal proteins that appear to have
172 evolved without much HGT and have been previously employed for the reconstruction of the tree of life
173 (25) were used to generate a concatenated alignment of 7411 positions, after removing low information
174 content positions (Supplementary Table 4). For the phylogenetic reconstruction, we used the IQ-tree
175 program with several phylogenetic models (see Methods and Supplementary Table 5 for details).
176 Surprisingly, the resulting trees had the 3D topology, with high support values for all key bifurcations
177 (Fig. 2c, additional data file 2).

178 A full investigation of the effects of different factors, in particular, the marker gene selection, on the tree
179 topology is beyond the scope of this work. Nonetheless, we addressed the possibility that the 3D topology
180 resulted from the model used for the tree reconstruction and/or the species selection. To this end, we
181 constructed 100 trees from the same alignment by randomly sampling 5 representatives of Asgard
182 archaea, other archaea, bacteria and eukaryotes each. For these smaller sets of species, the best model
183 identified by Williams et al. (LG+C60+G4+F) could be employed, resulting in 50 3D and 50 2D trees
184 (Supplementary Table 5, additional data file 2). Because IQ-tree identified this model as over-specified
185 for such a small alignment, we also tested a more restricted model (LG+C20+G4+F), obtaining 58 3D and
186 42 2D trees for the same set of 100 samples (Supplementary Table 5, additional data file 2).

187 The results of our phylogenetic analysis indicate that: 1) species sampling substantially affects the tree
188 topology; 2) even the set of most highly conserved genes that appear to be minimally prone to HGT,
189 yields conflicting signals for different species sets. Additional markers, less highly conserved and more
190 prone to HGT, are unlikely to improve phylogenetic resolution and might cause systematic error. Notably,
191 the topology of our complete phylogenetic tree (Fig. 2a) within the archaeal clade is mostly consistent

192 with the tree obtained in a preliminary analysis of a larger set of archaeal genomes and a larger marker
193 gene set (26). Taking into account these observations and the fact that we used the largest set of Asgard
194 archaea and other archaea compared to all previous phylogenetic analyses, the appearance of the 3D
195 topology in our tree indicates that the origin of eukaryotes from within Asgard cannot be considered a
196 settled issue. Various factors affecting the tree topology, including further increased species
197 representation, particularly, of Asgard and the deep branches of the TACK superphylum, such as
198 Bathyarchaeota and Korarchaeota, remain to be explored in order to definitively resolve the evolutionary
199 relationship between archaea and eukaryotes.

200

201 **The core gene set of Asgard archaea**

202 We next analyzed the core set of conserved Asgard genes which we arbitrarily defined as all asCOGs that
203 are present at least in one third of the MAGs in each of the 12 phylum-level lineages, with the mean
204 representation across lineages >75%. Under these criteria, the Asgard core includes 378 asCOGs
205 (Supplementary Table 6). As expected, most of these protein families, 293 (77%), are universal (present
206 in bacteria, other archaea and eukaryotes), 62 (16%) are represented in other archaea and eukaryotes, but
207 not in bacteria, 15 (4%) are found in other archaea and bacteria, but not in eukaryotes, 7 (2%) are archaea-
208 specific, and only 1 (0.003%) is shared exclusively with eukaryotes (Supplementary Figure 5). Most of
209 the core asCOGs show comparable levels of similarity to homologs from two or all three domains of life.
210 The second largest fraction of the core asCOGs shows substantially greater sequence similarity (at least,
211 25% higher similarity score) to homologous proteins from archaea than to those from eukaryotes and/or
212 bacteria (Supplementary Table 6). Compared with the 219 genes that comprise the pan-archaeal core (27),
213 the Asgard core set lacks 12 genes, each of which, however, is present in some subset of the Asgard
214 genomes. These include three genes of diphthamide biosynthesis and 2 ribosomal proteins, L40E and
215 L37E. The intricate evolutionary history of gene encoding translation elongation factors and enzymes of
216 diphthamide biosynthesis in Asgard has been analyzed previously (28). Also of note is the displacement
217 of the typical archaeal glyceraldehyde-3-phosphate dehydrogenase (type II) by a bacterial one (type I) in
218 most of the Asgard genomes (cog.001204, additional data file 1).

219 Functional distribution of the core asCOGs is shown in Fig. 3b (also see additional data file 1). For
220 comparison, we also derived an extended gene core for the TACK superphylum, using similar criteria (at
221 least 50% in each of the 6 lineages and 75% of the genomes overall, Fig 3b). For at least half of the
222 Asgard core genes, across most functional classes, there were no orthologs in the TACK core. The most
223 pronounced differences were found, as expected, in the category U (intracellular trafficking, secretion,

224 and vesicular transport). In Asgard archaea, this category includes 19 core genes compared with 7 genes
225 in TACK; 13 of these genes are specific to the Asgard archaea and include components of ESCRT I and
226 II, 3 distinct Roadblock/longin families, 2 distinct families of small GTPases, and a few other genes
227 implicated in related processes (Supplementary Table 6).

228 We compared the protein annotation obtained using asCOGs with the available annotation of ‘*Candidatus*
229 *Prometheoarchaeum syntrophicum*’ and found that using asCOGs allowed at least a general functional
230 prediction for 649 of the 1756 (37%) ‘hypothetical proteins’ in this organism, the only one in Asgard with
231 a closed genome. We also identified 139 proteins, in addition to the 80 described originally, that can be
232 considered Eukaryotic Signature Proteins, or ESPs (see next section).

233

234 **Eukaryotic features of Asgard archaea gleaned from genome analysis**

235 The enrichment of Asgard proteomes with homologs of eukaryote signature proteins (ESPs), such as
236 ESCRTs, components of protein sorting complexes including coat proteins, complete ubiquitin
237 machinery, actins and actin-binding proteins gelsolins and profilins, might be the strongest argument in
238 support of a direct evolutionary relationship between Asgard archaea and eukaryotes (2, 29). However,
239 the definition of ESPs is fuzzy because many of these proteins, in addition to their occurrence in Asgard,
240 are either scattered among several other archaeal genomes, often, from diverse groups (30), or consist of
241 promiscuous domains that are common in archaea, bacteria and eukaryotes, such as WD40 (after the
242 conserved terminal amino acids of the repeat units, also known as beta-transducin repeats), LRR
243 (Leucine-Rich Repeats), TPR (Tetratricopeptide Repeats), HEAT (Huntingtin-EF3-protein phosphatase
244 2A-TOR1) and other, largely, repetitive domains (31, 32). Furthermore, the sequences of some ESPs have
245 diverged to the extent that they become hardly detectable with standard computational methods. Our
246 computational strategy for delineating an extensive yet robust ESP set is described under Materials and
247 Methods. The ESP set we identified contained 505 asCOGs, including 238 that were not closely similar
248 (E-value= 10^{-10} , length coverage 75%) to those previously described by Zaremba-Niedzwiedzka et al.
249 (2)(Supplementary Table 7). In a general agreement with previous observations, the majority of these
250 ESPs, 329 of the 505, belonged to the ‘Intracellular trafficking, secretion, and vesicular transport’ (U)
251 functional class, followed by ‘Posttranslational modification, protein turnover, chaperones’ (O), with 101
252 asCOGs (Supplementary Table 7). Among the asCOGs in the U class, 130 were Roadblock/LC7
253 superfamily proteins, including longins, sybindin and profilins, and 94 were small GTPases of several
254 families, such as RagA-like, Arf-like and Rab-like ones, as discussed previously (33).

255 The phyletic patterns of ESP asCOGs in Asgard archaea are extremely patchy and largely lineage-specific
256 (Fig. 4), indicating that most of the proteins in this set are not uniformly conserved throughout Asgard
257 evolution, but rather, are prone to frequent HGT, gene losses and duplications. These evolutionary
258 processes are correlated in prokaryotes, resulting in the overall picture of highly dynamic evolution (34).
259 Even the most highly conserved ESP asCOG are missing in some Asgard lineages but show multiple
260 duplications in others (Fig. 4 and Supplementary Table 7). Surprisingly, many gaps in the ESPs
261 distribution were detected in the Heimdallarchaeota that include the suspected ancestors of eukaryotes.

262 Characteristically, many ESPs are multidomain proteins, with 37% assigned to more than one asCOG,
263 compared to 17% among non-ESP proteins (Supplementary Table 7). Some multidomain ESPs in Asgard
264 archaea have the same domain organizations as their homologs in eukaryotes, but these are a minority and
265 typically contain only two domains. Examples include the fusion of two EAP30/Vps37 domains (35), and
266 Vps23 and E2 domains (35) in ESCRT complexes, multiple Rag family GTPases, in which longin domain
267 is fused to the GTPase domain, and several others. By contrast, most of the domain architectures of the
268 multidomain ESP proteins were not detected in eukaryotes and often are found only in a narrow subset of
269 Asgard archaea, suggesting extensive domain shuffling during Asgard evolution (Fig. 5a). For example,
270 we identified many proteins containing a fusion of Vps28/Vps23 from ESCRT I complex (35) with C-
271 terminal domains of several homologous subunits of adaptin and COPI coatomer complexes (36, 37), and
272 E3 UFM1-protein ligase 1, which is involved in the UFM1 ubiquitin pathway (38) (Fig. 5a). Generally, a
273 protein with such a combination of domains can be predicted to be involved in ubiquitin-dependent
274 membrane remodeling but, because its domain architecture is unique, the precise function cannot be
275 inferred.

276 The majority of the ESP genes of Asgard archaea do not belong to conserved genomic neighborhoods, but
277 several such putative operons were detected. Perhaps, the most notable one is the ESCRT neighborhood
278 which includes genes coding for subunits of ESCRT I, II and III, and often, components of the ubiquitin
279 system (2), suggesting an ancient link between the two systems that persists in eukaryotes (35). We
280 predicted another operon that is conserved in most Asgard archaea and consists of genes encoding a
281 LAMTOR1-like protein of the Roadblock superfamily, a Rab-like small GTPase, and a protein containing
282 the DENN (differentially expressed in normal and neoplastic cells) domain that so far has been identified
283 only in eukaryotes (Fig. 5b). Two proteins consisting of a DENN domain fused to longin are subunits of
284 the folliculin (FLCN) complex that is conserved in eukaryotes. The FLCN complex is the sensor of amino
285 acid starvation interacting with Rag GTPase and Ragulator lysosomal complex, and a key component of
286 the mTORC1 pathway, the central regulator of cell growth in eukaryotes (39). Some Heimdallarchaea
287 encode several proteins with the exact same domain organization as FLCN (Fig. 5b). Ragulator is a

288 complex that consists of 5 subunits, each containing the Roadblock domain. In Asgard archaea, however,
289 the GTPase present in the operon is from a family that is distinct from the Rag GTPases, which interact
290 with both FLCN and Ragulator complexes in eukaryotes, despite the fact that Rag family GTPases are
291 abundant in Asgard (33) (Supplementary Table 7). Nevertheless, this conserved module of Asgard
292 proteins is a strong candidate to function as a guanine nucleotide exchange factor for Rab and Rag
293 GTPases, analogously to the eukaryotic FLCN. In eukaryotes, the DENN domain is present in many
294 proteins with different domain architectures that interact with different partners and perform a variety of
295 functions (40, 41). The Asgard archaea also encode other DENN domain proteins, and the respective
296 genes form expanded families of paralogs in Loki, Hel and Heimdall lineages, again, with domain
297 architectures distinct from those in eukaryotes (Fig. 5b) (42).

298 Prompted by the identification of a FLCN-like complex, we searched for other components of the
299 mTORC1 regulatory pathway in Asgard archaea. The GATOR1 complex that consists of three subunits,
300 Depdc5, Nprl2, and Nprl3, is another amino acid starvation sensor that is involved in this pathway in
301 eukaryotes (43). Nitrogen permease regulators 2 and 3 (NPRL2 and NPRL3) are homologous GATOR1
302 subunits that contain a longin domain and a small NPRL2-specific C-terminal domain (43). We identified
303 a protein family with this domain organization in most Thor MAGs and a few Loki MAGs. Several other
304 ESP asCOGs include proteins with high similarity to the longin domain of NPRL2. Additionally, we
305 identified many fusions of the NPRL2-like longin domain with various domains related to prokaryotic
306 two-component signal transduction system (Fig. 5c). Considering the absence of a homolog of
307 phosphatidylinositol 3-kinase, the catalytic domain of the mTOR protein, it seems likely that, in Asgard
308 archaea, the key growth regulation pathway remains centered at typical prokaryotic two-component signal
309 transduction systems whereas at least some of the regulators and sensors in this pathway are “eukaryotic”.
310 The abundance of NPRL2-like longin domains in Asgard archaea implies that the link between this
311 domain and amino acid starvation regulation emerged at the onset of Asgard evolution if not earlier.

312

313 **Diverse metabolic repertoires, ancestral metabolism of Asgard archaea, and syntrophic evolution**

314 Examination of the distribution of the asCOGs among the 12 Asgard archaeal phyla showed that the
315 metabolic pathway repertoire was conserved among the MAGs of each phylum but differed between the
316 phyla (Fig. 3a). Three distinct lifestyles were predicted by the asCOG analysis for different major
317 branches of Asgard archaea, namely, anaerobic heterotrophy, facultative aerobic heterotrophy, and
318 chemolithotrophy (Fig. 6, Supplementary Figure 9). For the last Asgard archaeal common ancestor
319 (LAsCA), a mixotrophic life style, including both production and consumption of H₂, can be inferred

320 from parsimony considerations (Fig. 6, Supplementary Table 8; see Materials and Methods for further
321 details). Loki-, Thor-, Hermod-, Baldr- and Borrarchaeota encode all enzymes for the complete (archaeal)
322 Wood-Ljungdahl pathway (WLP) and are predicted to oxidize organic substrates, likely, by using the
323 reverse WLP, given the lack of enzymes for oxidation of inorganic compounds (e.g., hydrogen,
324 sulfur/sulfide and nitrogen/ammonia). The genomes of these five Asgard phyla encode homologues of
325 membrane-bound respiratory H₂-evolving Group 4 [NiFe] hydrogenase (Supplementary Figure 6) and/or
326 cytosolic cofactor-coupled bidirectional Group 3 [NiFe] hydrogenase (44) (Supplementary Figure 7).
327 Phylogenetic analysis of both group 3 and group 4 [NiFe] hydrogenases showed that Asgard archaea form
328 distinct clades well separated from the functionally characterized hydrogenases, hampering the prediction
329 of their specific functions in Asgard archaea. The functionally characterized group 4 [NiFe] hydrogenases
330 in the Thermococci are involved in the fermentation of organic substrates to H₂, acetate and carbon
331 dioxide (45, 46). The presence of group 3 [NiFe] hydrogenases suggests that these Asgard archaea cannot
332 use H₂ as an electron donor because they lack the enzyme complex coupling H₂ oxidation to membrane
333 potential generation. Thus, in these organisms, bifurcate electrons from H₂ are likely to be used to support
334 the fermentation of organic substrates exclusively (45–47).

335 Both Wukongarchaeota genomes (As_075 and As_085) encode a bona fide membrane-bound Group 1k
336 [NiFe] hydrogenase that could mediate hydrogenotrophic respiration using heterodisulfide as the terminal
337 electron acceptor (48, 49) (Fig. 6, Supplementary Figure 8, Supplementary Figure 10). The group 1k
338 [NiFe] hydrogenase is exclusively found in methanogens of the order Methanosarcinales (Euryarchaeota)
339 (50), and it is the first discovery of the group 1 [NiFe] hydrogenase in the Asgard archaea.

340 Wukongarchaeota also encode all enzymes for a complete WLP and a putative ADP-dependent acetyl-
341 CoA synthetase for acetate synthesis. Unlike all other Asgard archaea, Wukongarchaeota lack genes for
342 citrate cycle and beta-oxidation. Thus, Wukongarchaeota appear to be obligate chemolithotrophic
343 acetogens. The genomes of Wukongarchaeota were discovered only in seawater of the euphotic zone of
344 the Yap trench (0 m and 125 m). Dissolved H₂ concentration is known to be the highest in surface
345 seawater, where the active microbial fermentation, compared to deep sea (51), could produce sufficient
346 amounts of hydrogen for the growth of Wukongarchaeota. Hodarchaeota, Gerdarchaeota, Kariarchaeota,
347 and Heimdallarchaeota share a common ancestor with Wukongarchaeota (Fig. 6). However, genome
348 analysis implies different lifestyles for these organisms. Hod-, Gerd- and Kariarchaeota encode various
349 electron transport chain components, including heme/copper-type cytochrome/quinol oxidase, nitrate
350 reductase, and NADH dehydrogenase, most likely, allowing the use of oxygen and nitrate as electron
351 acceptors during aerobic and anaerobic respiration, respectively (44). In addition, Hod-, Gerd- and
352 Heimdallarchaeota encode phosphoadenosine phosphosulfate (PAPS) reductase and adenylylsulfate
353 kinase for sulfate reduction, enabling the use of sulfate as electron acceptor during anaerobic respiration.

354 Gerd-, Heimdall-, and Hodarchaeota are only found in coastal and deep-sea sedimentary environments,
355 whereas Kariarchaeota were found also in marine water. The versatile predicted metabolic capacities of
356 these groups suggest that Hod-, Gerd- and Kariarchaeota might occupy both anoxic and oxic niches. In
357 contrast, Heimdallarchaeota appear to be able to thrive only in anoxic environments.

358 In the widely considered syntrophy scenarios (52), eukaryogenesis has been proposed to involve
359 metabolic symbiosis (syntrophy) between an archaeon and one or two bacterial partners which, in the
360 original hydrogen/syntrophy hypothesis, were postulated to donate H₂ for methane or hydrogen sulfide
361 production by the consortium (53, 54). The syntrophic scenarios were boosted by the discovery of
362 apparent syntrophy between *Candidatus P. syntrophicum* and Deltaproteobacteria which led to the
363 proposal of the Etangle-Engulf-Endogenize (E³) model of eukaryogenesis (17, 55). Reconstruction of the
364 Lokiarchaeon metabolism has suggested that this organism was hydrogen-dependent, in accord with the
365 hydrogen-syntrophic scenarios (54). In contrast, subsequent analysis of the metabolic potentials of 4
366 Asgard phyla has led to the inference that these organisms were primarily organoheterotrophic and H₂-
367 producing, the ‘reverse flow model’ model of protoeukaryote energy metabolism that involves electron or
368 hydrogen flow from an Asgard archaeon to the alphaproteobacterial ancestor of mitochondria, in the
369 opposite direction from that in the original hydrogen-syntrophy hypotheses (44). Here, we discovered a
370 deeply branching Asgard group, Wukongarchaeota, which appears to include obligate hydrogenotrophic
371 acetogens, suggesting the possibility of the LAsCA being a hydrogen-dependent autotroph (Fig. 6). This
372 finding suggests that LAsCA both produced and consumed H₂. Thus, depending on the exact relationship
373 between Asgard archaea and eukaryotes that remains to be elucidated, our findings could be compatible
374 with different syntrophic scenarios that postulate H₂ transfer from bacteria to the archaeal symbiont or in
375 the opposite direction.

376

377 **Conclusions**

378 The Asgard archaea that were discovered only 5 years ago as a result of the painstaking assembly of
379 several Loki Castle metagenomes have grown into a highly diverse archaeal superphylum. The most
380 remarkable feature of the Asgard is their apparent evolutionary affinity with eukaryotes that has been
381 buttressed by two independent lines of evidence: phylogenetic analysis of highly conserved genes and
382 detection of multiple ESPs that are absent or far less common in other archaea. The 75 MAGs added here
383 substantially expand the phylogenetic and metabolic diversity of the Asgard superphylum. The extended
384 set of Asgard genomes provides for a phylogeny based on a far more representative species sampling than
385 available previously and a substantially expanded ESP analysis employing powerful computational

386 methods. This extended analysis reveals 6 putative additional Asgard phyla but does not immediately
387 clarify the Asgard-eukaryote relationship. Indeed, phylogenetic analysis of conserved genes in an
388 expanded set of archaea, bacteria and eukaryotes yields conflicting 3D and 2D signals, with an
389 unexpected preference for the 3D topology. Thus, the conclusion that eukaryotes emerged from within
390 Asgard archaea, in particular, from the Heimdall lineage, appears to be premature. Further phylogenomic
391 study with an even broader representation of diverse archaeal lineages as well as, possibly, even more
392 sophisticated evolutionary models are required to clarify the relationships between archaea and
393 eukaryotes.

394 Our analysis of Asgard genomes substantially expanded the set of ESPs encoded by this group of archaea
395 and revealed numerous, complex domain architectures of these proteins. These results further emphasize
396 the excess of ESPs in Asgard compared to other archaea and provide additional support to the conclusion
397 that most of the Asgard ESPs are involved in membrane remodeling and intracellular trafficking.
398 However, in parallel with the phylogenomic results, detailed analysis of the ESPs reveals a complex
399 picture. Most of the multidomain Asgard ESPs possess domain architectures distinct from typical
400 eukaryotic ones and some of these arrangements include signature prokaryotic domains, suggesting
401 substantial functional differences from the respective eukaryotic systems. Furthermore, virtually all the
402 ESPs show patchy distributions in Asgard and other archaea, indicative of a history of extensive HGT,
403 gene losses and paralogous family expansion. All these findings seem to be best compatible with the
404 model of a dispersed, dynamic archaeal ‘eukaryome’ (30) that widely spreads among archaea via HGT, so
405 far reaching the highest ESP density in the Asgard archaea.

406 The results of this work cannot rule out the possibility of the emergence of eukaryotes from within the
407 Asgard but seem to be better compatible with a different evolutionary scenario under which the conserved
408 core of eukaryote genes involved in informational processes originates from an as yet unknown ancestor
409 group that might be a deep archaeal branch or could lie outside the presently characterized archaeal
410 diversity. These hypothetical ancestral forms might have accumulated components of the mobile archaeal
411 ‘eukaryome’ to an even greater extent than the Asgard archaea, eventually, giving rise to eukaryote-like
412 cells, likely, via a form of syntrophy with one or more bacterial partners. Combined genomic and
413 (undoubtedly, far more challenging) biological study of diverse archaea is essential for further advancing
414 our understanding of eukaryogenesis.

415

416 **Author contributions**

417 ML, EVK, KSM and YL initiated the study;; YL performed metagenomic assembly, binning, metabolism
418 analysis; KSM, AN and YIW performed comparative genomic analysis; YL, KSM, YIW and WCH
419 performed phylogenetic analysis; KSM and YIW constructed asCOGs; KSM, YIW, YL, ML, and EVK
420 analyzed the data; YL, KSM, WCH, EVK, and ML wrote the manuscript that was read, edited and
421 approved by all authors.

422

423 **Acknowledgements**

424 ML and YL are supported by National Natural Science Foundation of China (Grant No. 91851105,
425 31970105 and 31700430), the Key Project of Department of Education of Guangdong Province
426 (No.2017KZDXM071), and the Shenzhen Science and Technology Program (Grant no.
427 JCYJ20170818091727570 and KQTD20180412181334790). KSM, YIW, SAS and EVK are supported
428 by the Intramural Research Program of the National Institutes of Health of the USA (National Library of
429 Medicine).

430

431 **Competing interests**

432 The authors declare no competing interests

433 **References**

- 434 1. A. Spang, J. H. Saw, S. L. Jørgensen, K. Zaremba-Niedzwiedzka, J. Martijn, A. E. Lind,
435 R. van Eijk, C. Schleper, L. Guy, T. J. G. Ettema, Complex archaea that bridge the gap between
436 prokaryotes and eukaryotes. *Nature*. **521**, 173–179 (2015).
- 437 2. K. Zaremba-Niedzwiedzka, E. F. Caceres, J. H. Saw, D. B. Ckstr m, L. Juzokaite, E.
438 Vancaester, K. W. Seitz, K. Anantharaman, P. Starnawski, K. U. Kjeldsen, M. B. Stott, T.
439 Nunoura, J. F. Banfield, A. Schramm, B. J. Baker, A. Spang, T. J. G. Ettema, Asgard archaea
440 illuminate the origin of eukaryotic cellular complexity. *Nature*. **541**, 353–358 (2017).
- 441 3. F. MacLeod, G. S. Kindler, H. L. Wong, R. Chen, B. P. Burns, Asgard archaea:
442 Diversity, function, and evolutionary implications in a range of microbiomes. *AIMS Microbiol.*
443 **5**, 48–61 (2019).
- 444 4. M. Cai, Y. Liu, X. Yin, Z. Zhou, M. W. Friedrich, T. Richter-Heitmman, R. Nimzyk, A.
445 Kulkarni, X. Wang, W. Li, J. Pan, Y. Yang, J.-D. Gu, M. Li, Diverse Asgard archaea including
446 the novel phylum Gerdarchaeota participate in organic matter degradation. *Sci. China Life Sci.*
447 **63**, 886–897 (2020).
- 448 5. T. A. Williams, C. J. Cox, P. G. Foster, G. J. Szöllösi, T. M. Embley, Phylogenomics
449 provides robust support for a two-domains tree of life. *Nat Ecol Evol* (2019),
450 doi:10.1038/s41559-019-1040-x.
- 451 6. T. A. Williams, P. G. Foster, C. J. Cox, T. M. Embley, An archaeal origin of eukaryotes
452 supports only two primary domains of life. *Nature*. **504**, 231–236 (2013).
- 453 7. C. J. Cox, P. G. Foster, R. P. Hirt, S. R. Harris, T. M. Embley, The archaeobacterial origin
454 of eukaryotes. *Proc Natl Acad Sci U S A*. **105**, 20356–20361 (2008).
- 455 8. N. Yutin, K. S. Makarova, S. L. Mekhedov, Y. I. Wolf, E. V. Koonin, The deep archaeal
456 roots of eukaryotes. *Mol. Biol. Evol.* **25**, 1619–1630 (2008).
- 457 9. V. Da Cunha, M. Gaia, D. Gabelle, A. Nasir, P. Forterre, Lokiarchaea are close relatives
458 of Euryarchaeota, not bridging the gap between prokaryotes and eukaryotes. *PLoS Genet.* **13**,
459 e1006810 (2017).
- 460 10. V. D. Cunha, M. Gaia, A. Nasir, P. Forterre, Asgard archaea do not close the debate
461 about the universal tree of life topology. *PLoS Genet.* **14**, e1007215 (2018).
- 462 11. P. Forterre, The universal tree of life: an update. *Front. Microbiol.* **6**, 717 (2015).
- 463 12. J. Lombard, P. López-García, D. Moreira, The early evolution of lipid membranes and
464 the three domains of life. *Nat. Rev. Microbiol.* **10**, 507–515 (2012).
- 465 13. J. A. Burns, A. A. Pittis, E. Kim, Gene-based predictive models of trophic modes suggest
466 Asgard archaea are not phagocytotic. *Nat. Ecol. Evol.* **6**, a016006 (2018).

- 467 14. C. Akil, R. C. Robinson, Genomes of Asgard archaea encode profilins that regulate actin.
468 *Nature*. **562**, 439–443 (2018).
- 469 15. C. Akil, L. T. Tran, M. Orhant-Prioux, Y. Baskaran, E. Manser, L. Blanchoin, R. C.
470 Robinson, Insights into the evolution of regulated actin dynamics via characterization of
471 primitive gelsolin/cofilin proteins from Asgard archaea. *Proc Natl Acad Sci U S A*. **117**, 19904–
472 19913 (2020).
- 473 16. Z. Lu, T. Fu, T. Li, Y. Liu, S. Zhang, J. Li, J. Dai, E. V. Koonin, G. Li, H. Chu, M. Li,
474 Coevolution of Eukaryote-like Vps4 and ESCRT-III Subunits in the Asgard Archaea. *mBio*. **11**,
475 e00417-20, /mbio/11/3/mBio.00417-20.atom (2020).
- 476 17. H. Imachi, M. K. Nobu, N. Nakahara, Y. Morono, M. Ogawara, Y. Takaki, Y. Takano,
477 K. Uematsu, T. Ikuta, M. Ito, Y. Matsui, M. Miyazaki, K. Murata, Y. Saito, S. Sakai, C. Song, E.
478 Tasumi, Y. Yamanaka, T. Yamaguchi, Y. Kamagata, H. Tamaki, K. Takai, Isolation of an
479 archaeon at the prokaryote-eukaryote interface. *Nature*. **577**, 519–525 (2020).
- 480 18. S. M. Karst, M. S. Dueholm, S. J. McIlroy, R. H. Kirkegaard, P. H. Nielsen, M.
481 Albertsen, *Nat. Biotechnol.*, in press, doi:10.1038/nbt.4045.
- 482 19. R.-Y. Zhang, B. Zou, Y.-W. Yan, C. O. Jeon, M. Li, M. Cai, Z.-X. Quan, Design of
483 targeted primers based on 16S rRNA sequences in meta-transcriptomic datasets and
484 identification of a novel taxonomic group in the Asgard archaea. *BMC microbiology*. **20**, 25
485 (2020).
- 486 20. K. S. Makarova, Y. I. Wolf, E. V. Koonin, Towards functional characterization of
487 archaeal genomic dark matter. *Biochem Soc Trans*. **47**, 389–398 (2019).
- 488 21. E. V. Koonin, Y. I. Wolf, Genomics of bacteria and archaea: the emerging dynamic view
489 of the prokaryotic world. *Nucleic Acids Res*. **36**, 6688–6719 (2008).
- 490 22. A. R. Nabhan, I. N. Sarker, The impact of taxon sampling on phylogenetic inference: a
491 review of two decades of controversy. *Brief Bioinform*. **13**, 122–134 (2012).
- 492 23. R. Gouy, D. Baurain, H. Philippe, Rooting the tree of life: the phylogenetic jury is still
493 out. *Philos Trans R Soc Lond B Biol Sci*. **370**, 20140329 (2015).
- 494 24. J. Felsenstein, *Inferring Phylogenies* (Oxford University Press, Oxford, New York, 2nd
495 Edition., 2003).
- 496 25. F. D. Ciccarelli, T. Doerks, C. von Mering, C. J. Creevey, B. Snel, P. Bork, Toward
497 automatic reconstruction of a highly resolved tree of life. *Science (New York, N.Y.)*. **311**, 1283–
498 1287 (2006).
- 499 26. C. Rinke, M. Chuvochina, A. J. Mussig, P.-A. Chaumeil, D. W. Waite, W. B. Whitman,
500 D. H. Parks, P. Hugenholtz, *bioRxiv*, in press, doi:10.1101/2020.03.01.972265.

- 501 27. K. S. Makarova, Y. I. Wolf, E. V. Koonin, Archaeal Clusters of Orthologous Genes
502 (arCOGs): An Update and Application for Analysis of Shared Features between
503 Thermococcales, Methanococcales, and Methanobacteriales. *Life (Basel)*. **5**, 818–840 (2015).
- 504 28. A. B. Narrowe, A. Spang, C. W. Stairs, E. F. Caceres, B. J. Baker, C. S. Miller, T. J. G.
505 Ettema, Complex Evolutionary History of Translation Elongation Factor 2 and Diphthamide
506 Biosynthesis in Archaea and Parabasalids. *Genome Biol Evol.* **10**, 2380–2393 (2018).
- 507 29. L. Eme, A. Spang, J. Lombard, C. W. Stairs, T. J. G. Ettema, Archaea and the origin of
508 eukaryotes. *Nat. Rev. Microbiol.* **15**, 711–723 (2017).
- 509 30. E. V. Koonin, N. Yutin, The dispersed archaeal eukaryome and the complex archaeal
510 ancestor of eukaryotes. *Cold Spring Harb Perspect Biol.* **6**, a016188 (2014).
- 511 31. R. L. Tatusov, N. D. Fedorova, J. D. Jackson, A. R. Jacobs, B. Kiryutin, E. V. Koonin, D.
512 M. Krylov, R. Mazumder, S. L. Mekhedov, A. N. Nikolskaya, B. S. Rao, S. Smirnov, A. V.
513 Sverdlov, S. Vasudevan, Y. I. Wolf, J. J. Yin, D. A. Natale, The COG database: an updated
514 version includes eukaryotes. *BMC Bioinformatics.* **4**, 41 (2003).
- 515 32. E. V. Koonin, N. D. Fedorova, J. D. Jackson, A. R. Jacobs, D. M. Krylov, K. S.
516 Makarova, R. Mazumder, S. L. Mekhedov, A. N. Nikolskaya, B. S. Rao, I. B. Rogozin, S.
517 Smirnov, A. V. Sorokin, A. V. Sverdlov, S. Vasudevan, Y. I. Wolf, J. J. Yin, D. A. Natale, A
518 comprehensive evolutionary classification of proteins encoded in complete eukaryotic genomes.
519 *Genome Biol.* **5**, R7 (2004).
- 520 33. C. M. Klinger, A. Spang, J. B. Dacks, T. J. G. Ettema, Tracing the Archaeal Origins of
521 Eukaryotic Membrane-Trafficking System Building Blocks. *Mol Biol Evol.* **33**, 1528–1541
522 (2016).
- 523 34. P. Puigbò, A. E. Lobkovsky, D. M. Kristensen, Y. I. Wolf, E. V. Koonin, Genomes in
524 turmoil: quantification of genome dynamics in prokaryote supergenomes. *BMC Biol.* **12**, 66
525 (2014).
- 526 35. L. Christ, C. Raiborg, E. M. Wenzel, C. Campsteijn, H. Stenmark, Cellular Functions and
527 Molecular Mechanisms of the ESCRT Membrane-Scission Machinery. *Trends Biochem Sci.* **42**,
528 42–56 (2017).
- 529 36. N. Gomez-Navarro, E. Miller, Protein sorting at the ER-Golgi interface. *J Cell Biol.* **215**,
530 769–778 (2016).
- 531 37. K. M. K. Kibria, J. Ferdous, R. Sardar, A. Panda, D. Gupta, A. Mohammed, P. Malhotra,
532 A genome-wide analysis of coatomer protein (COP) subunits of apicomplexan parasites and their
533 evolutionary relationships. *BMC Genomics.* **20**, 98 (2019).
- 534 38. Y. Wei, X. Xu, UFMylation: A Unique & Fashionable Modification for Life. *Genomics*
535 *Proteomics Bioinformatics.* **14**, 140–146 (2016).
- 536 39. R. E. Lawrence, S. A. Fromm, Y. Fu, A. L. Yokom, D. J. Kim, A. M. Thelen, L. N.
537 Young, C.-Y. Lim, A. J. Samelson, J. H. Hurley, R. Zoncu, Structural mechanism of a Rag

- 538 GTPase activation checkpoint by the lysosomal folliculin complex. *Science*. **366**, 971–977
539 (2019).
- 540 40. M.-Y. Su, S. A. Fromm, R. Zoncu, J. H. Hurley, Structure of the C9orf72 ARF GAP
541 complex that is haploinsufficient in ALS and FTD. *Nature*. **585**, 251–255 (2020).
- 542 41. N. de Martín Garrido, C. H. S. Aylett, Nutrient Signaling and Lysosome Positioning
543 Crosstalk Through a Multifunctional Protein, Folliculin. *Front Cell Dev Biol*. **8**, 108 (2020).
- 544 42. A. L. Marat, H. Dokainish, P. S. McPherson, DENN domain proteins: regulators of Rab
545 GTPases. *J Biol Chem*. **286**, 13791–13800 (2011).
- 546 43. K. Shen, R. K. Huang, E. J. Brignole, K. J. Condon, M. L. Valenstein, L. Chantranupong,
547 A. Bomaliyamu, A. Choe, C. Hong, Z. Yu, D. M. Sabatini, Architecture of the human GATOR1
548 and GATOR1-Rag GTPases complexes. *Nature*. **556**, 64–69 (2018).
- 549 44. A. Spang, C. W. Stairs, N. Dombrowski, L. Eme, J. Lombard, E. F. Caceres, C.
550 Greening, B. J. Baker, T. J. G. Ettema, Proposal of the reverse flow model for the origin of the
551 eukaryotic cell based on comparative analyses of Asgard archaeal metabolism. *Nat. Microbiol*. **4**,
552 1138–1148 (2019).
- 553 45. H. Yu, C.-H. Wu, G. J. Schut, D. K. Haja, G. Zhao, J. W. Peters, M. W. W. Adams, H.
554 Li, Structure of an Ancient Respiratory System. *Cell*. **173**, 1636-1649.e16 (2018).
- 555 46. G. J. Schut, E. S. Boyd, J. W. Peters, M. W. W. Adams, The modular respiratory
556 complexes involved in hydrogen and sulfur metabolism by heterotrophic hyperthermophilic
557 archaea and their evolutionary implications. *FEMS Microbiol Rev*. **37**, 182–203 (2013).
- 558 47. F. O. Bryant, M. W. Adams, Characterization of hydrogenase from the hyperthermophilic
559 archaeobacterium, *Pyrococcus furiosus*. *J Biol Chem*. **264**, 5070–5079 (1989).
- 560 48. U. Deppenmeier, M. Blaut, B. Schmidt, G. Gottschalk, Purification and properties of a
561 F420-nonreactive, membrane-bound hydrogenase from *Methanosarcina* strain Gö1. *Arch*
562 *Microbiol*. **157**, 505–511 (1992).
- 563 49. R. K. Thauer, A.-K. Kaster, M. Goenrich, M. Schick, T. Hiromoto, S. Shima,
564 Hydrogenases from methanogenic archaea, nickel, a novel cofactor, and H₂ storage. *Annu Rev*
565 *Biochem*. **79**, 507–536 (2010).
- 566 50. D. Søndergaard, C. N. S. Pedersen, C. Greening, HydDB: A web tool for hydrogenase
567 classification and analysis. *Scientific Reports*. **6**, 34212 (2016).
- 568 51. R. Conrad, W. Seiler, Methane and hydrogen in seawater (Atlantic Ocean). *Deep Sea*
569 *Res. Part I Oceanogr. Res. Pap*. **35**, 1903–1917 (1988).
- 570 52. P. López-García, D. Moreira, The Syntrophy hypothesis for the origin of eukaryotes
571 revisited. *Nat Microbiol*. **5**, 655–667 (2020).
- 572 53. W. Martin, M. Müller, The hydrogen hypothesis for the first eukaryote. *Nature*. **392**, 37–
573 41 (1998).

574 54. D. Moreira, P. López-García, Symbiosis Between Methanogenic Archaea and δ -
575 Proteobacteria as the Origin of Eukaryotes: The Syntrophic Hypothesis. *J Mol Evol.* **47**, 517–530
576 (1998).

577 55. P. López-García, D. Moreira, Cultured Asgard Archaea Shed Light on Eukaryogenesis.
578 *Cell.* **181**, 232–235 (2020).

579

580

581

582 .

583

584

585 **Figures**

586

587 **Figure 1. Gene commonality plots for Asgard archaea and the TACK superphylum.**

588 The gene commonality plot showing the number of asCOGs in log scale (Y-axis) that include the given
589 fraction of analyzed genomes (X-axis). The Asgard plot is compared with the TACK superphylum plot
590 based on assignment of TACK genomes to arCOGs.

591

592 **Figure 2. Phylogenetic analysis of Asgard archaea and their relationships with eukaryotes.**

593 **(a)** Maximum likelihood tree, inferred with IQ-tree and LG+F+R10 model, constructed from
594 concatenated alignments of the protein sequences from 209 core Asgard Clusters of Orthologs (asCOGs).
595 Only the 12 phylum-level clades are shown, with species within each clade collapsed. See supplementary
596 Methods and Supplementary Table 5 for details.

597 **(b)** Maximum likelihood tree, inferred with IQ-tree and SYM+R8 model, based on 16S rRNA gene
598 sequences. Red stars in **(b)** denote MAGs reconstructed in the current study.

599 **(c)** Phylogenetic tree of bacteria, archaea and eukaryotes, inferred with IQ-tree under LG+R10 model,
600 constructed from concatenated alignments of the protein sequences of 30 universally conserved genes
601 (see Material and Methods for details). The tree shows the relationships between the major clades.

602 The trees are unrooted and are shown in a pseudorooted form for visualization purposes only. The actual
603 trees and alignments are in Additional data file 2 and list of the trees are provided in the Supplementary
604 Table 4 and 5.

605

606 **Figure 3. Phyletic patterns of asCOGs and functional distribution of Asgard core genes.**

607 **(a)** Classical Multidimensional Scaling analysis of binary presence-absence phyletic patterns for 13,939
608 asCOGs that are represented in at least two genomes (see Material and Methods for details).

609 **(b)** Functional breakdown of Asgard core genes (378 asCOGs) compared with TACK superphylum core
610 genes (489 arCOGs). The values were normalized as described in Materials and Methods. Functional
611 classes of genes: J, Translation, ribosomal structure and biogenesis; K, Transcription; L, Replication,

612 recombination and repair; D, Cell cycle control, cell division, chromosome partitioning; V, Defense
613 mechanisms; T, Signal transduction mechanisms; M, Cell wall/membrane/envelope biogenesis; N, Cell
614 motility; U, Intracellular trafficking, secretion, and vesicular transport; O, Posttranslational modification,
615 protein turnover, chaperones; X, Mobilome: prophages, transposons; C, Energy production and
616 conversion; G, Carbohydrate transport and metabolism; E, Amino acid transport and metabolism; F,
617 Nucleotide transport and metabolism; H, Coenzyme transport and metabolism; I, Lipid transport and
618 metabolism; P, Inorganic ion transport and metabolism; Q, Secondary metabolites biosynthesis, transport
619 and catabolism; R, General function prediction only; S, Function unknown;

620

621 **Figure 4. Phyletic patterns of Eukaryotic Signature Proteins (ESPs) encoded in Asgard genomes.**

622 All 505 ESP asCOGs are grouped by distance between binary presence-absence phyletic patterns. The
623 most highly conserved ESP asCOGs are shown within the red rectangle. Below the plot of the number of
624 ESP domains in each genome is shown. For details, see Supplementary Table 7.

625

626 **Figure 5. Domain architectures of selected Asgard ESPs.**

627 **(a) ESPs with unique domain architectures.** The schematic of each multidomain protein is roughly
628 proportional to the respective protein length. The identified domains are shown inside the arrows
629 approximately according to their location and are briefly annotated. Homologous domains are shown by
630 the same color or pattern.

631 **(b) DENN domain proteins in Asgards.** Upper part of the figure above the dashed line shows putative
632 operon encoding DENN domain proteins. Genes are shown by block arrows with the length proportional
633 to the size of the corresponding protein. For each protein, the nucleotide contig or genome partition
634 accession number, Asgard genome ID and lineage are indicated. The part of the figure below the dashed
635 line shows domain organization of diverse proteins containing DENN domain. Homologous domains are
636 shown by the same color. The inset on the right side explains identity of the domains

637 **(c) NPRL2-like proteins in Asgards.** Designations are the same as in Figure 4a.

638

639 **Figure 6. Reconstruction and evolution of the key metabolic processes in Asgard archaea.**

640 The schematic phylogeny of Asgard archaea is from Figure 1a.

641 LAsCA, Last Asgard Common Ancestor; WLP, Wood-Ljungdahl pathway.

642 **Supplementary Materials**

643 Taxonomic description

644 Materials and Methods

645 Supplementary References

646 Tables S1-S8

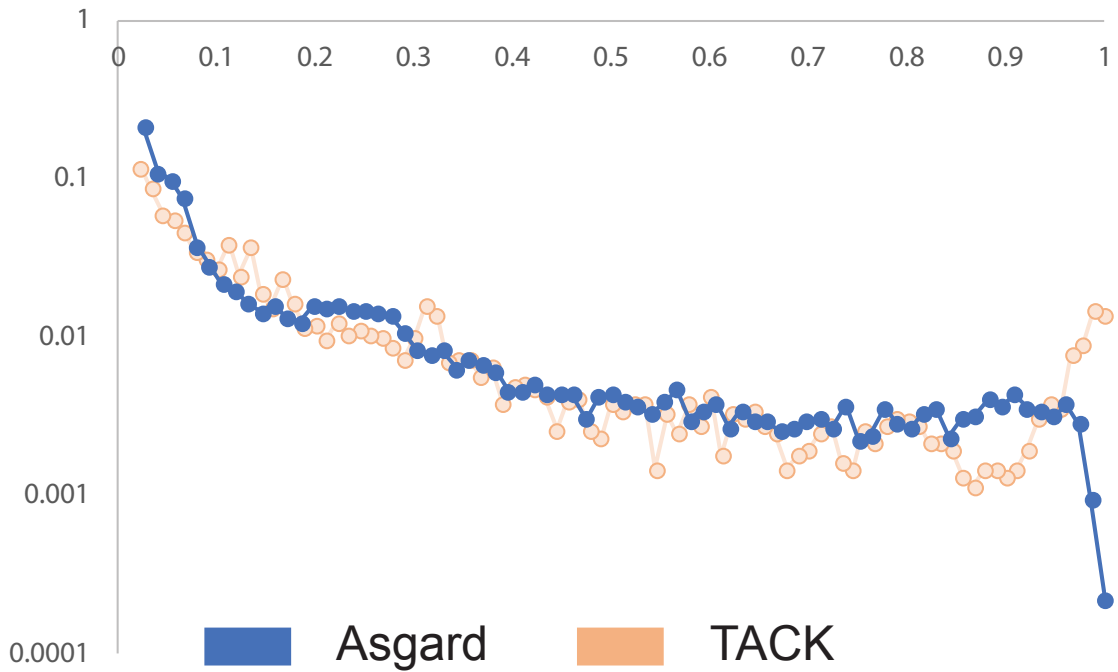
647 Figures S1-S10

648 Additional Data files are available at ftp://ftp.ncbi.nih.gov/pub/wolf/_suppl/asgard20/

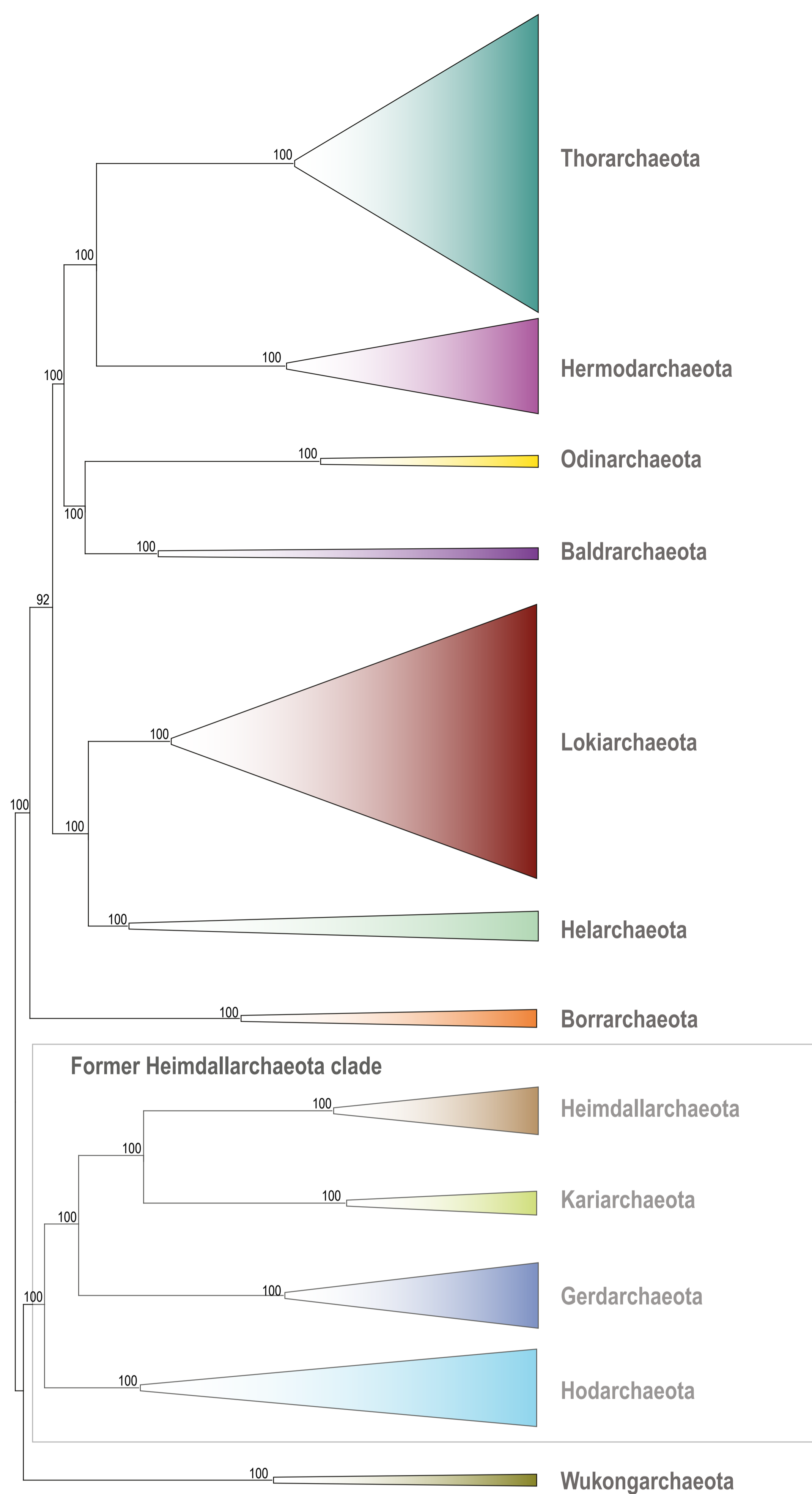
649 Additional data file 1 (Additional_data_file_1.tgz): Complete asCOG data archive

650 Additional data file 2 (Additional_data_file_2.tgz): Phylogenetic trees and alignments archive

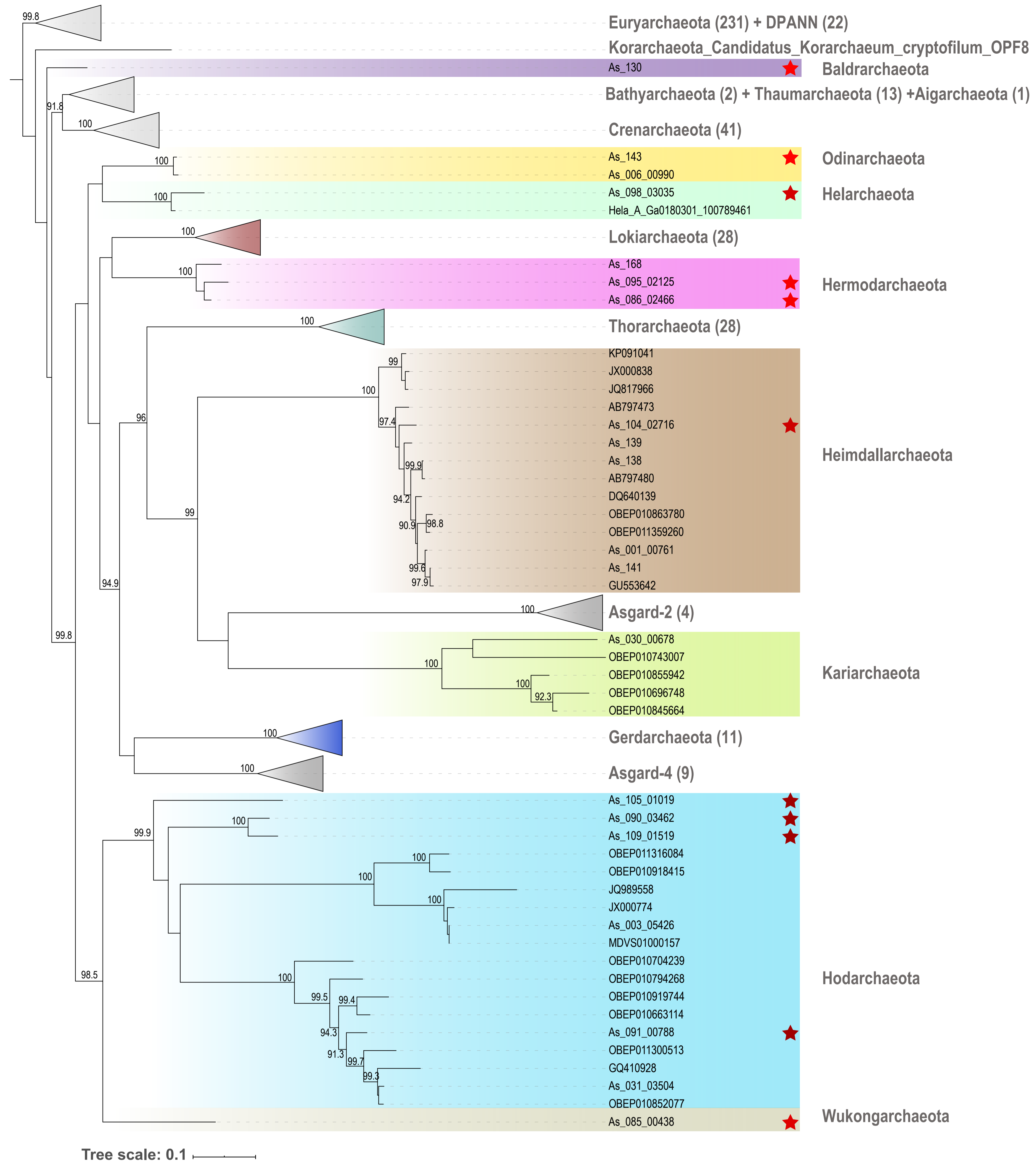
651



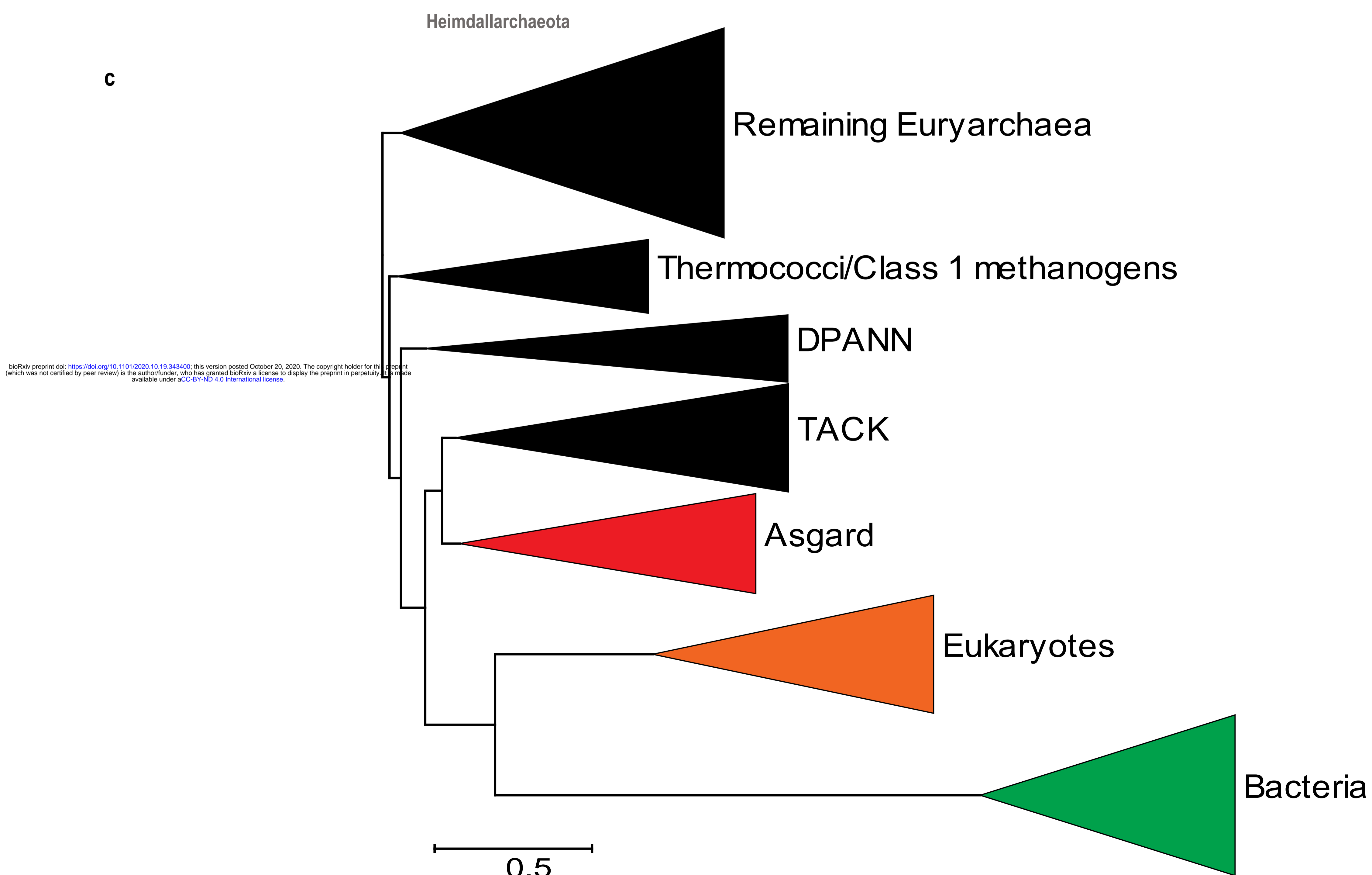
a



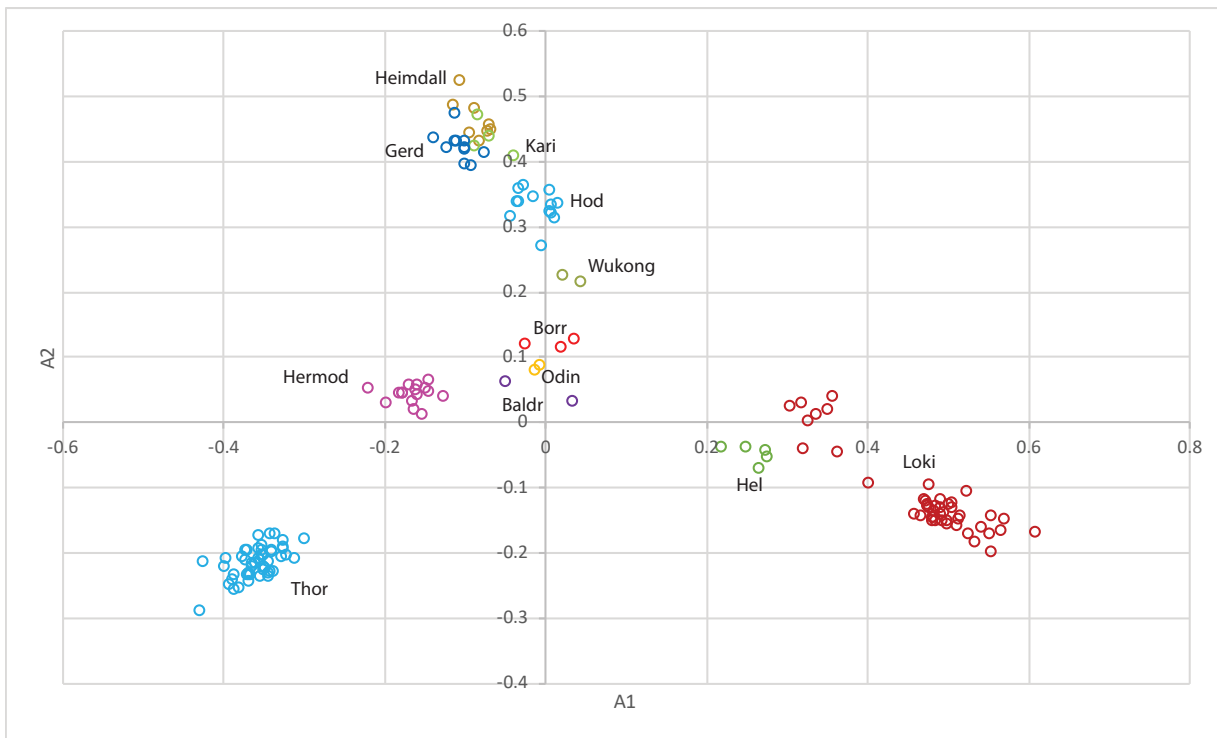
b



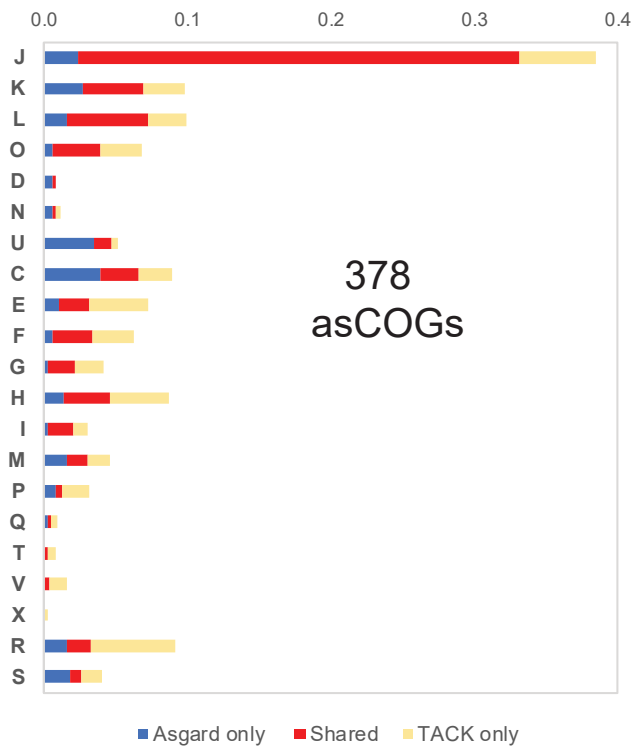
c



A



B





Thor

Hermod

Odin
Baldr

Loki

Hel

Borr

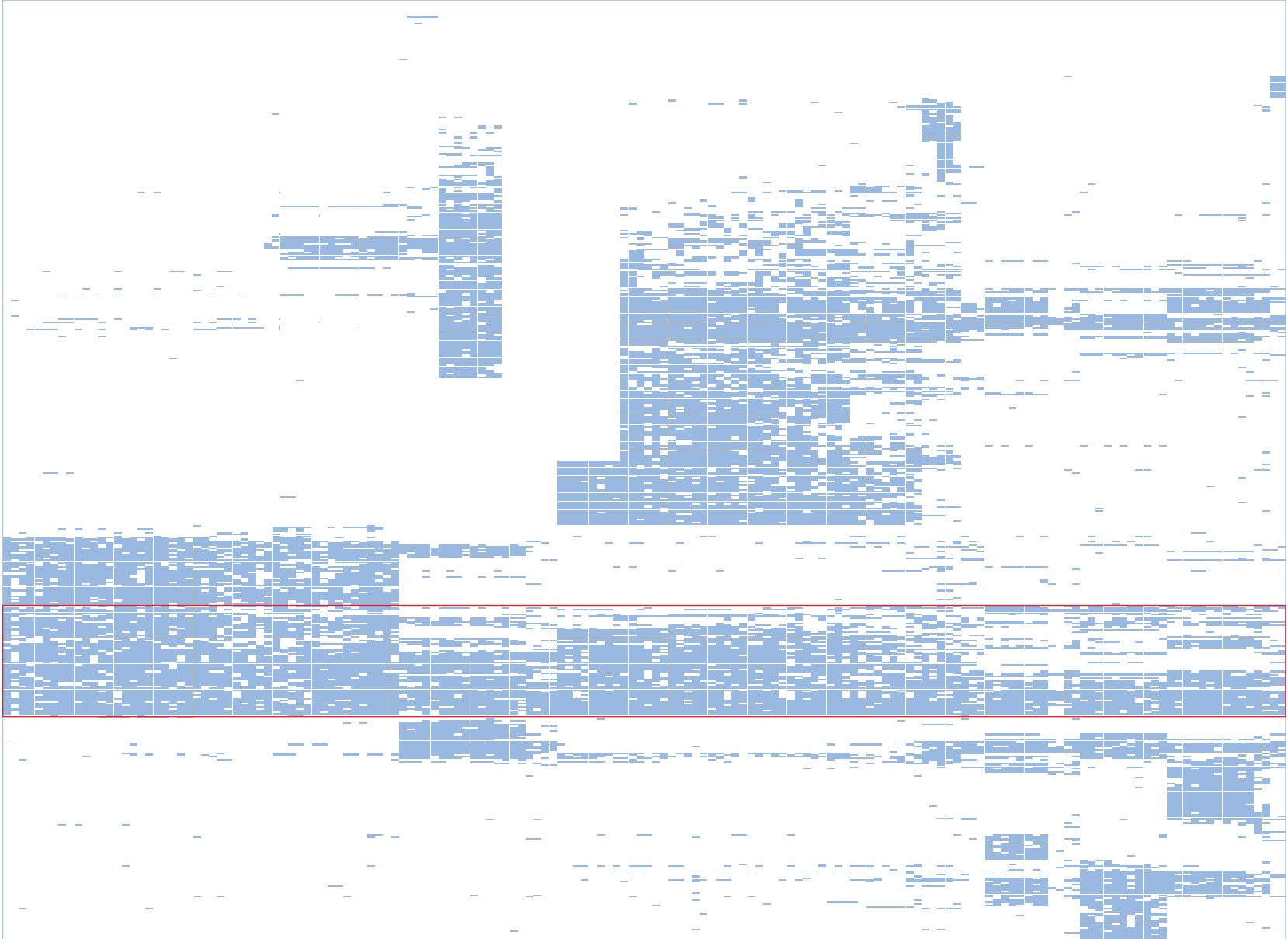
Heimdall

Kari

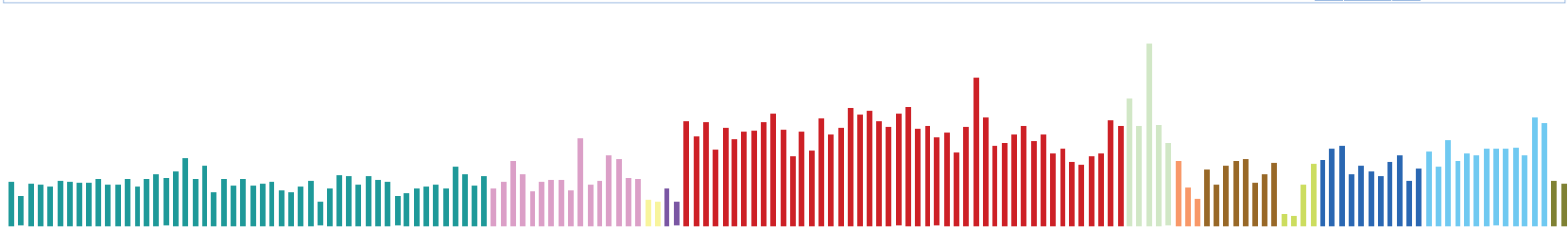
Herd

Hod

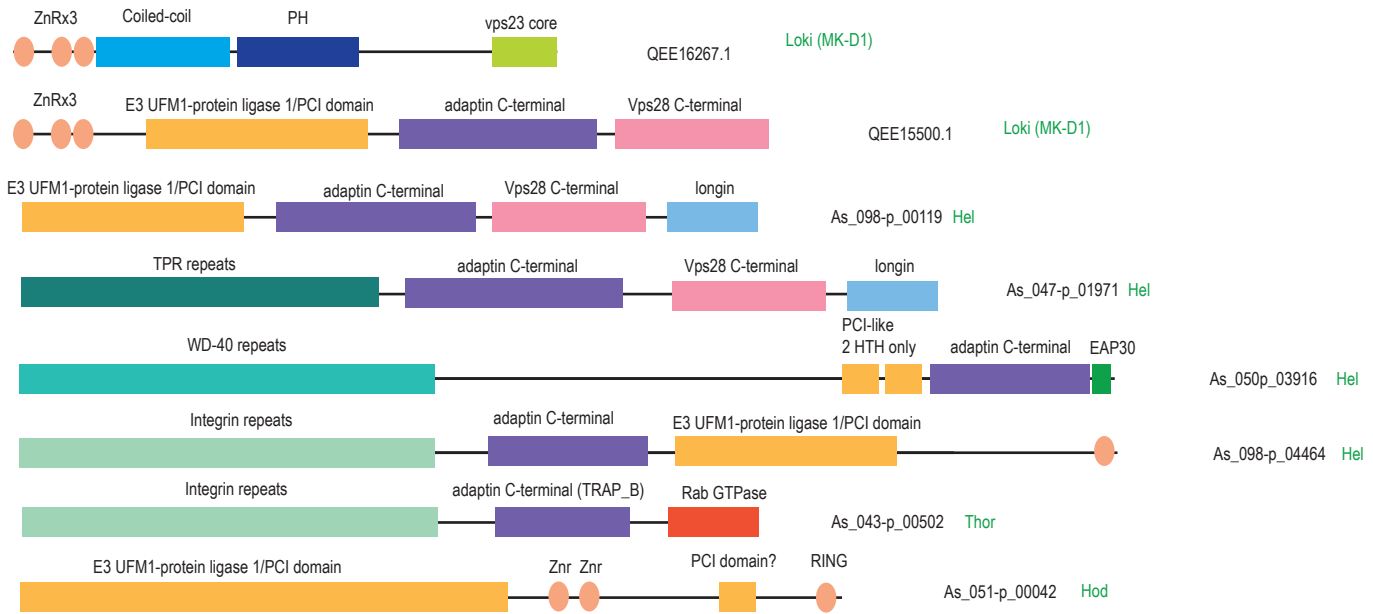
Wukong



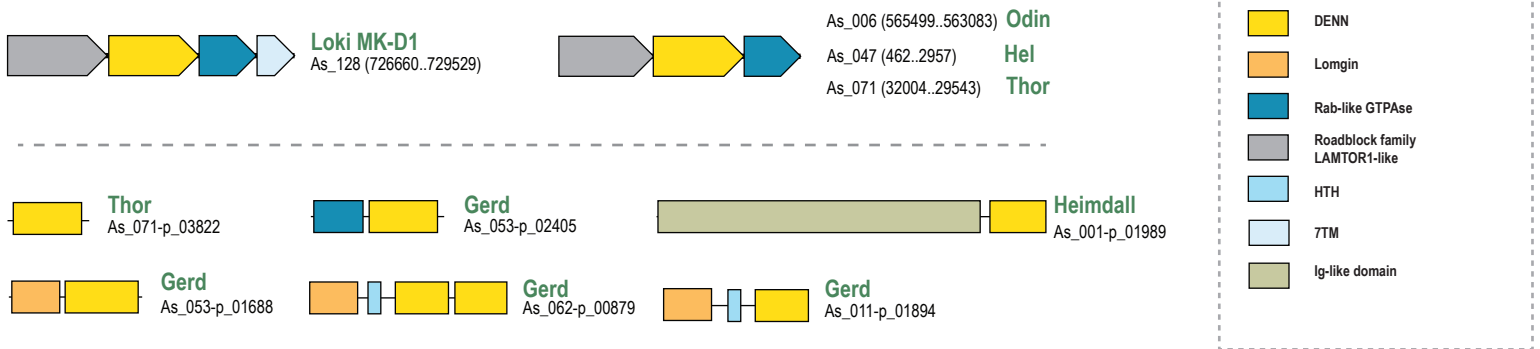
ESCRT I, II, III;
actin, gelsolins and profilins;
ubiquitin-system components E1,E2,E3;



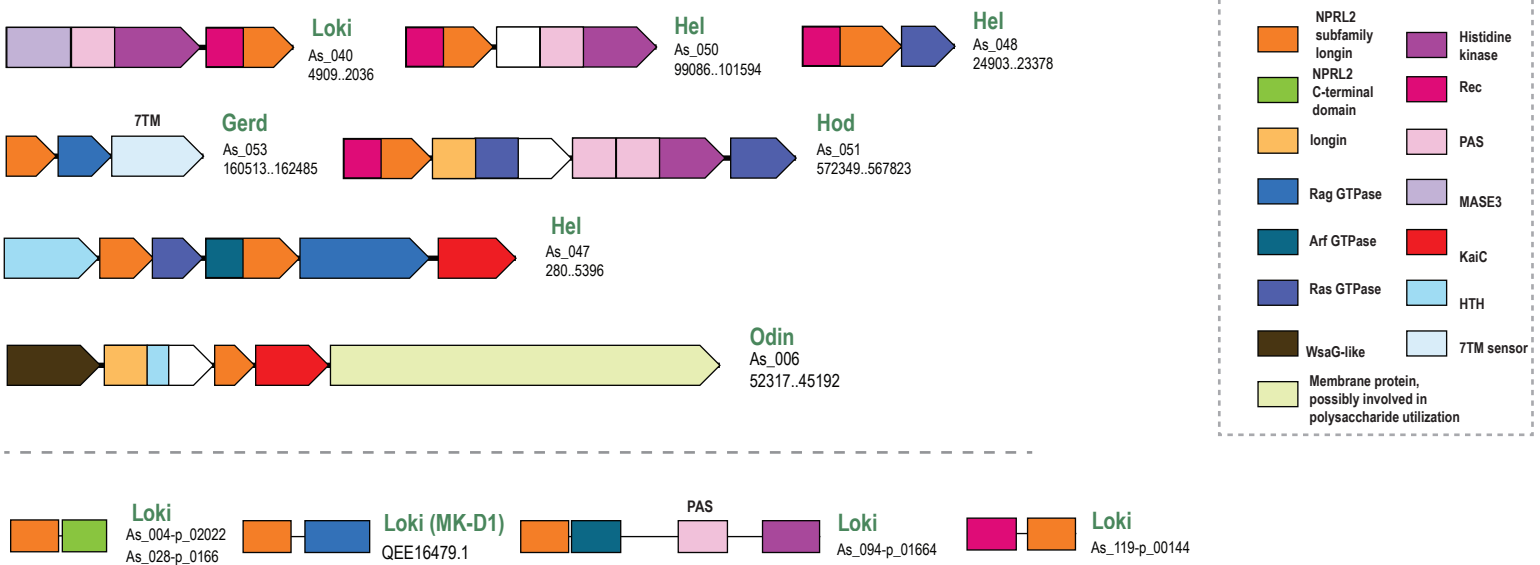
A



B



C

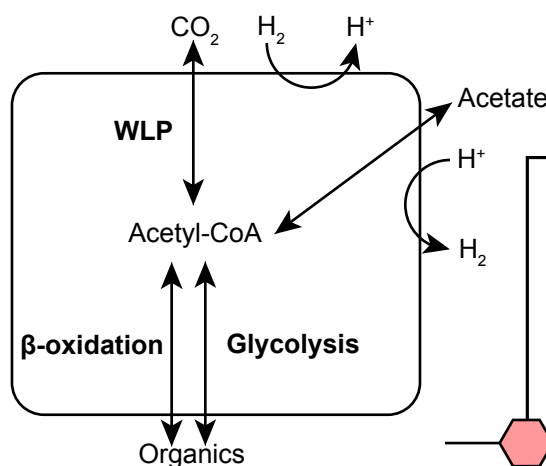


Anaerobic heterotroph
Thor-, Hermod-, Odin-, Baldr-, Loki-, Hel-, Borr-, Heimdallarchaeota

Facultative aerobic heterotroph
Kari-, Gerd-, Hodarchaeota

Chemolithotroph
Wukongarchaeota

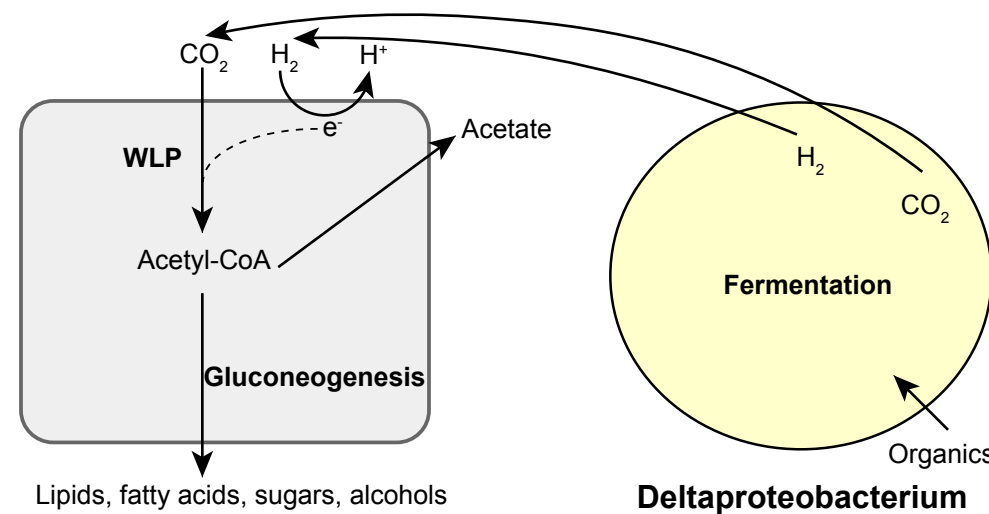
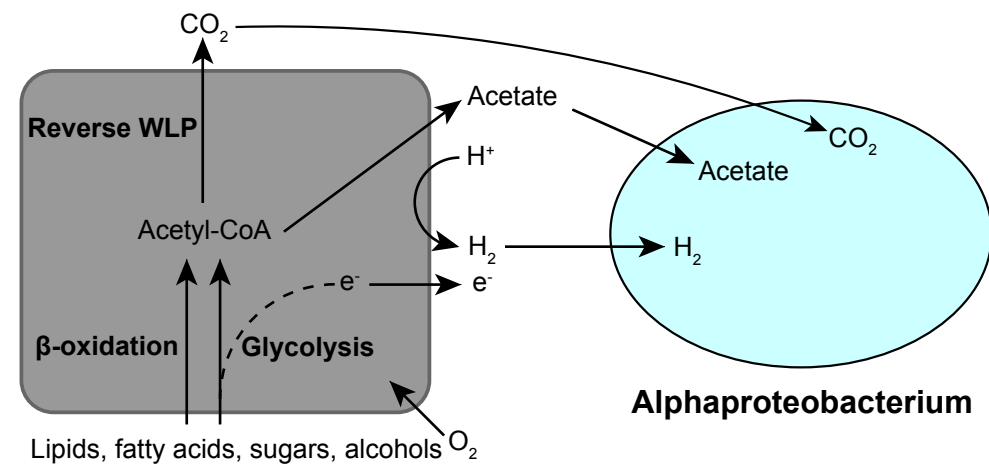
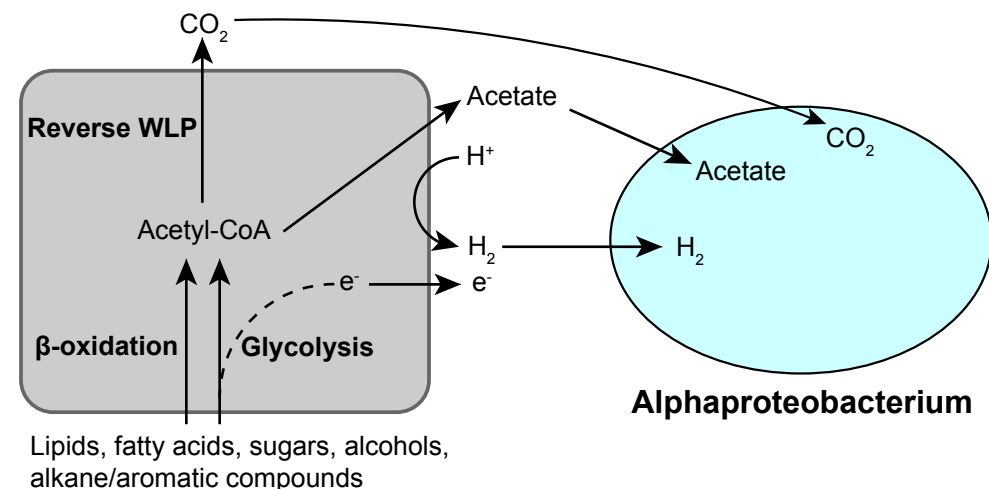
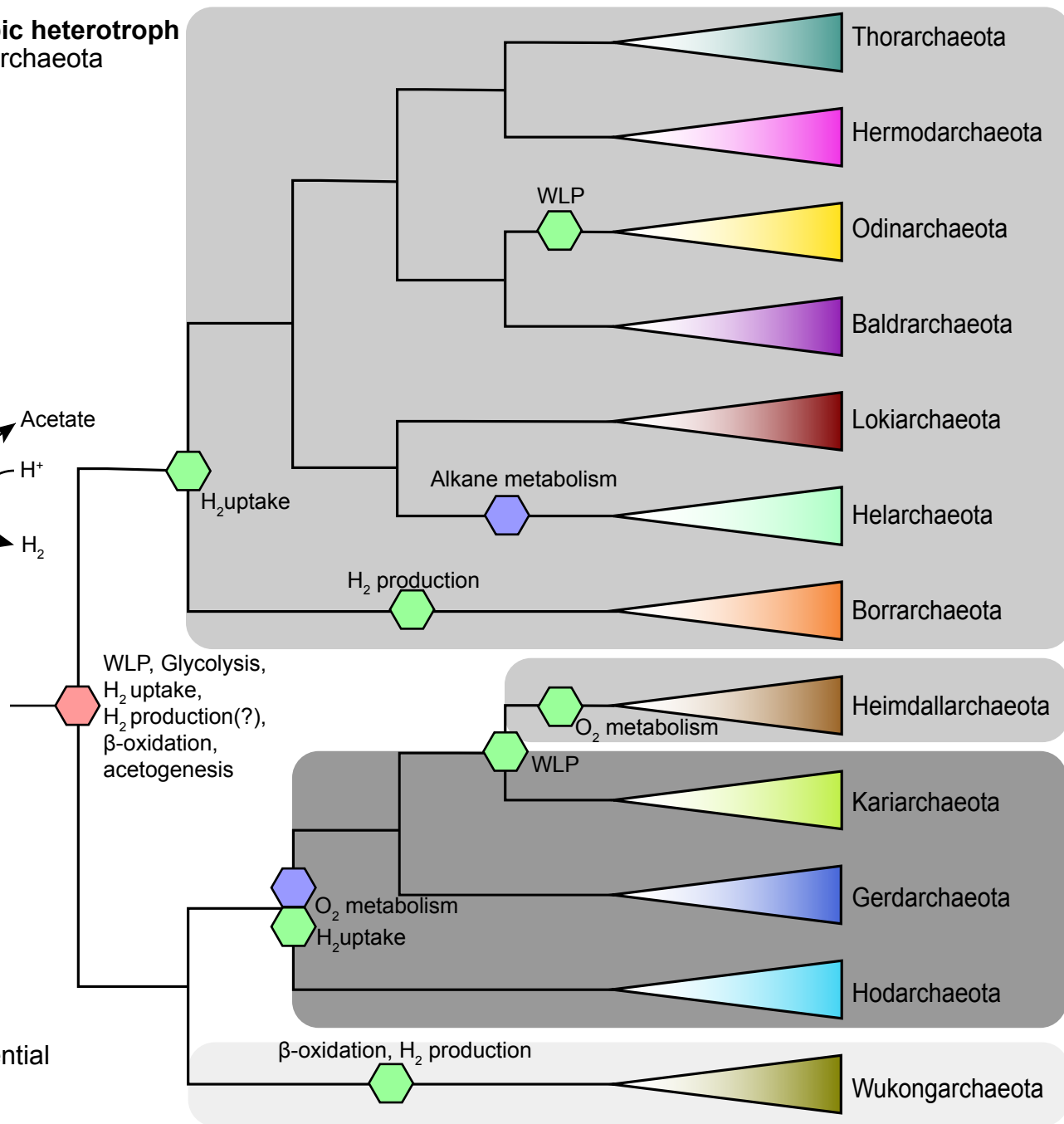
Mixotrophic LAsCA



LAsCA

Loss of metabolic potential

Gain of metabolic potential



1 **Expanding diversity of Asgard archaea and the elusive ancestry of eukaryotes**

2 Yang Liu^{1†}, Kira S. Makarova^{2†}, Wen-Cong Huang^{1†}, Yuri I. Wolf², Anastasia Nikolskaya², Xinxu
3 Zhang¹, Mingwei Cai¹, Cui-Jing Zhang¹, Wei Xu³, Zhuhua Luo³, Lei Cheng⁴, Eugene V. Koonin^{2*}, Meng
4 Li^{1*}

5 1 Shenzhen Key Laboratory of Marine Microbiome Engineering, Institute for Advanced Study, Shenzhen
6 University, Shenzhen, Guangdong, 518060, P. R. China

7 2 National Center for Biotechnology Information, National Library of Medicine, National Institutes of
8 Health, Bethesda, Maryland 20894, USA

9 3 State Key Laboratory Breeding Base of Marine Genetic Resources, Key Laboratory of Marine Genetic
10 Resources, Fujian Key Laboratory of Marine Genetic Resources, Third Institute of Oceanography, State
11 Oceanic Administration, Xiamen 361005, P. R. China

12 4 Key Laboratory of Development and Application of Rural Renewable Energy, Biogas Institute of
13 Ministry of Agriculture, Chengdu 610041, P.R. China

14 † These authors contributed equally to this work.

15 *Authors for correspondence: koonin@ncbi.nlm.nih.gov or limeng848@szu.edu.cn

16

17	Supplementary Information	
18	Taxonomic Description of new taxa.....	1
19	Materials and Methods.....	9
20	Supplementary references.....	15
21	Supplementary Figures.....	19
22	Supplementary Tables.....	21
23		
24		

25 Taxonomic Description of new taxa

26 '**Candidatus Wukongarchaeum**' (Wu.kong.ar.chae'um. N. L. n. Wukong a legendary Chinese figure,
27 also known as Monkey King, who caused havoc in the heavenly palace); N.L. neut. N. *archaeum* (from
28 Gr. adj. archaios ancient) archaeon; N. L. neut. N. Wukongarchaeum.

29 '**Candidatus Wukongarchaeum yapensis**' (yap'ensis N. L. masc. adj. pertaining to Yap trench, which is
30 the geographical position where the first type material of this species was obtained). Type material is the
31 genome designated as As_085 (Yap4.bin4.70) representing '*Candidatus Wukongarchaeum yapensis*'. The
32 genome "As_085" represents a MAG consisting of 2.16 Mbps in 277 contigs with an estimated
33 completeness of 92.52%, an estimated contamination of 4.05%, a 16S and 23S rRNA gene and 14 tRNAs.
34 The MAG recovered from a marine water metagenome (Yap trench, Western Pacific), with an estimated
35 depth of coverage of 31.4, has a GC content of 38%.

36 '**Candidatus Hodarchaeum**' (Hod.ar.chae'um. N. L. n. Hod a son of Odin in Norse mythology); N.L.
37 neut. N. *archaeum* (from Gr. adj. archaios ancient) archaeon; N. L. neut. N. Hodarchaeum.

38 '**Candidatus Hodarchaeum mangrovi**' (man.gro'vi N.L. fem. n. of a mangrove, referring to the isolation
39 of the type material from mangrove soil). Type material is the genome designated as As_027
40 (FT2_5_011) representing '*Candidatus Hodarchaeum mangrovi*'. The genome "As_027" represents a
41 MAG consisting of 4.01 Mbps in 348 contigs with an estimated completeness of 93.61%, an estimated
42 contamination of 0.93%, a 23S rRNA gene and 14 tRNAs. The MAG recovered from mangrove sediment
43 metagenomes (Futian Nature Reserve, China), with an estimated depth of coverage of 17.9, has a GC
44 content of 32.9%.

45 '**Candidatus Kariarchaeum**' (Ka.ri.ar.chae'um. N. L. n. Kari the god of wind and brother to Aegir in
46 Norse mythology); N.L. neut. N. *archaeum* (from Gr. adj. archaios ancient) archaeon; N. L. neut. N.
47 Kariarchaeum.

48 '**Candidatus Kariarchaeum pelagius**' (pe.la'gi.us. L. masc. adj. of or belonging to the sea, referring to
49 the isolation of the type material from the Ocean). Type material is the genome designated as As_030
50 (RS678) representing '*Candidatus Kariarchaeum pelagius*'. The genome "As_030" represents a MAG
51 consisting of 1.41 Mbps in 76 contigs, an estimated completeness of 83.18%, with an estimated
52 contamination of 1.87%, a 23S, 16S and 5S rRNA genes and 18 tRNAs. The MAG recovered from a
53 marine metagenome (Saudi Arabia: Red Sea) has a GC content of 30.11%.

54 '**Candidatus Borrarchaeum**' (Borr.ar.chae'um. N. L. n. Borr a creator god and father of Odin); N.L.
55 neut. N. *archaeum* (from Gr. adj. archaios ancient) archaeon; N. L. neut. N. Borrarchaeum.

56 '**Candidatus Borrarchaeum yapensis**' (yap'ensis N. L. masc. adj. pertaining to Yap trench, which is the
57 geographical position where the first type material of this species was obtained). Type material is the
58 genome designated as As_181 (Yap2000.bin9.141) representing '*Candidatus Borrarchaeum yapensis*'.
59 The genome "As_181" represents a MAG consisting of 3.63 Mbps in 125 contigs, with an estimated
60 completeness of 95.02%, an estimated contamination of 5.61% and 11 tRNAs. The MAG, recovered from
61 a marine water metagenome (Yap trench, Western Pacific) with an estimated depth coverage of 15.04, has
62 a GC content of 37.1%.

63 **‘*Candidatus Baldrarchaeum*’** (Bal.dr.ar.chae’um. N. L. n. Baldr the god of light and son of Odin and
64 borther of Thor in Norse mythology); N.L. neut. N. *archaeum* (from Gr. adj. archaios ancient) archaeon;
65 N. L. neut. N. Baldrarchaeum.

66 **‘*Candidatus Baldrarchaeum yapensis*’** (yap’ensis N. L. masc. adj. pertaining to Yap trench, which is the
67 geographical position where the first type material of this species was obtained). Type material is the
68 genome designated as As_130 (Yap30.bin9.72) representing ‘*Candidatus Baldrarchaeum yapensis*’. The
69 genome “As_130” represents a MAG consisting of 2.27 Mbps in 100 contigs, with an estimated
70 completeness of 93.93%, an estimated contamination of 3.74%, a 23S and 16S rRNA gene and 15 tRNAs.
71 The MAG, recovered from a marine water metagenome (Yap trench, Western Pacific) with an estimated
72 depth coverage of 39.99, has a GC content of 45.95%.

73 **‘*Candidatus Hermodarchaeum*’** (Her.mod.ar.chae’um. N. L. n. Hermod, messengers of the gods in the
74 Norse mythology and son of Odin and brother of Baldr in the Norse mythology); N.L. neut. N. *archaeum*
75 (from Gr. adj. archaios ancient) archaeon; N. L. neut. N. Hermodarchaeum.

76 **‘*Candidatus Hermodarchaeum yapensis*’** (yap’ensis N. L. masc. adj. pertaining to Yap trench, which is
77 the geographical position where the first type material of this species was obtained). Type material is the
78 genome designated as As_086 (Yap4.bin9.105) representing ‘*Candidatus Hermodarchaeum yapensis*’.
79 The genome ‘As_086’ represent a MAG consisting of 2.71 Mbps in 77 contigs, with an estimated
80 completeness of 92.99%, an estimated contamination of 1.87%, a 23S and 16S rRNA gene and 16 tRNAs.
81 The MAG, recovered from a marine water metagenome (Yap trench, Western Pacific) with an estimated
82 depth coverage of 19.24, has a GC content of 44.69%.

83 ***Candidatus Wukongarchaeaceae*** (Wu.kong.ar.chae.a.ce’ae. N.L. neut. n. Wukongarchaeum a
84 (*Candidatus*) type genus of the family; -aceae ending to denote the family; N.L. fem. pl. n.
85 Wukongarchaeaceae the Wukongarchaeum family).

86 The family is delineated based on 209 concatenated Asgard Cluster of Orthologs (AsCOGs) and 16S
87 rRNA gene phylogeny. The description is the same as that of its sole genus and species. Type genus is
88 *Candidatus Wukongarchaeum*.

89 ***Candidatus Wukongarchaeales*** (Wu.kong.ar.chae.a’les. N.L. neut. n. Wukongarchaeum a (*Candidatus*)
90 type genus of the order; -ales ending to denote the order; N.L. fem. pl. n. Wukongarchaeales the
91 Wukongarchaeum order).

92 The order is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
93 description is the same as that of its sole genus and species. Type genus is *Candidatus Wukongarchaeum*.

94 ***Candidatus Wukongarchaeia*** (Wu.kong.ar.chae’i.a. N.L. neut. n. Wukongarchaeum a (*Candidatus*) type
95 genus of the order of the class; -ia ending to denote the class; N.L. fem. pl. n. Wukongarchaeia the
96 Wukongarchaeum class).

97 The class is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
98 description is the same as that of its sole and type order *Candidatus Wukongarchaeales*.

- 99 ***Candidatus Wukongarchaeota*** (Wu.kong.ar.chae.o'ta. N.L. neut. n. Wukongarchaeum a (Candidatus)
100 type genus of the class of the phylum; -ota ending to denote the phylum; N.L neut. pl. n.
101 Wukongarchaeota the Wukongarchaeum phylum)
- 102 The phylum is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
103 description is the same as that of its sole and type class *Candidatus* Wukongarchaeia.
- 104 ***Candidatus Hodarchaeaceae*** (Hod.ar.chae.a.ce'ae. N.L. neut. n. Hodarchaeum a (Candidatus) type
105 genus of the family; -aceae ending to denote the family; N.L. fem. pl. n. Hodarchaeaceae the
106 Hodarchaeum family).
- 107 The family is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
108 description is the same as that of its sole genus and species. Type genus is *Candidatus* Hodarchaeum.
- 109 ***Candidatus Hodarchaeales*** (Hod.ar.chae.a'les. N.L. neut. n. Hodarchaeum a (Candidatus) type genus of
110 the order; -ales ending to denote the order; N.L fem. pl. n. Hodarchaeales the Hodarchaeum order).
- 111 The order is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
112 description is the same as that of its sole genus and species. Type genus is *Candidatus* Hodarchaeum.
- 113 ***Candidatus Hodarchaeia*** (Hod.ar.chae'i.a. N.L. neut. n. Hodarchaeum a (Candidatus) type genus of the
114 order of the class; -ia ending to denote the class; N.L fem. pl. n. Hodarchaeia the Hodarchaeum class).
- 115 The class is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
116 description is the same as that of its sole and type order *Candidatus* Hodarchaeales.
- 117 ***Candidatus Hodarchaeota*** (Hod.ar.chae.o'ta. N.L. neut. n. Hodarchaeum a (Candidatus) type genus of
118 the class of the phylum; -ota ending to denote the phylum; N.L neut. pl. n. Hodarchaeota the
119 Hodarchaeum phylum)
- 120 The phylum is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
121 description is the same as that of its sole and type class *Candidatus* Hodarchaeia.
- 122 ***Candidatus Hodarchaeaceae*** (Hod.ar.chae.a.ce'ae. N.L. neut. n. Hodarchaeum a (Candidatus) type
123 genus of the family; -aceae ending to denote the family; N.L. fem. pl. n. Hodarchaeaceae the
124 Hodarchaeum family).
- 125 The family is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
126 description is the same as that of its sole genus and species. Type genus is *Candidatus* Hodarchaeum.
- 127 ***Candidatus Hodarchaeales*** (Hod.ar.chae.a'les. N.L. neut. n. Hodarchaeum a (Candidatus) type genus of
128 the order; -ales ending to denote the order; N.L fem. pl. n. Hodarchaeales the Hodarchaeum order).
- 129 The order is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
130 description is the same as that of its sole genus and species. Type genus is *Candidatus* Hodarchaeum.
- 131 ***Candidatus Hodarchaeia*** (Hod.ar.chae'i.a. N.L. neut. n. Hodarchaeum a (Candidatus) type genus of the
132 order of the class; -ia ending to denote the class; N.L fem. pl. n. Hodarchaeia the Hodarchaeum class).

133 The class is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
134 description is the same as that of its sole and type order *Candidatus* Hodarchaeales.

135 ***Candidatus* Hodarchaeota** (Hod.ar.chae.o'ta. N.L. neut. n. Hodarchaeum a (*Candidatus*) type genus of
136 the class of the phylum; -ota ending to denote the phylum; N.L. neut. pl. n. Hodarchaeota the
137 Hodarchaeum phylum)

138 The phylum is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
139 description is the same as that of its sole and type class *Candidatus* Hodarchaeia.

140 ***Candidatus* Kariarchaeaceae** (Ka.ri.ar.chae.a.ce'ae. N.L. neut. n. Kariarchaeum a (*Candidatus*) type
141 genus of the family; -aceae ending to denote the family; N.L. fem. pl. n. Kariarchaeaceae the
142 Kariarchaeum family).

143 The family is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
144 description is the same as that of its sole genus and species. Type genus is *Candidatus* Kariarchaeum.

145 ***Candidatus* Kariarchaeales** (Ka.ri.ar.chae.a'les. N.L. neut. n. Kariarchaeum a (*Candidatus*) type genus of
146 the order; -ales ending to denote the order; N.L. fem. pl. n. Kariarchaeales the Kariarchaeum order).

147 The order is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
148 description is the same as that of its sole genus and species. Type genus is *Candidatus* Kariarchaeum.

149 ***Candidatus* Kariarchaeia** (Ka.ri.ar.chae'i.a. N.L. neut. n. Kariarchaeum a (*Candidatus*) type genus of the
150 order of the class; -ia ending to denote the class; N.L. fem. pl. n. Kariarchaeia the Kariarchaeum class).

151 The class is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
152 description is the same as that of its sole and type order *Candidatus* Kariarchaeales.

153 ***Candidatus* Kariarchaeota** (Ka.ri.ar.chae.o'ta. N.L. neut. n. Kariarchaeum a (*Candidatus*) type genus of
154 the class of the phylum; -ota ending to denote the phylum; N.L. neut. pl. n. Kariarchaeota the
155 Kariarchaeum phylum)

156 The phylum is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
157 description is the same as that of its sole and type class *Candidatus* Kariarchaeia.

158 ***Candidatus* Borrarchaeaceae** (Borr.ar.chae.a.ce'ae. N.L. neut. n. Borrarchaeum a (*Candidatus*) type
159 genus of the family; -aceae ending to denote the family; N.L. fem. pl. n. Borrarchaeaceae the
160 Borrarchaeum family).

161 The family is delineated based on 209 concatenated AsCOGs phylogeny. The description is the same as
162 that of its sole genus and species. Type genus is *Candidatus* Borrarchaeum.

163 ***Candidatus* Borrarchaeales** (Borr.ar.chae.a'les. N.L. neut. n. Borrarchaeum a (*Candidatus*) type genus of
164 the order; -ales ending to denote the order; N.L. fem. pl. n. Borrarchaeales the Borrarchaeum order).

165 The order is delineated based on 209 concatenated AsCOGs phylogeny. The description is the same as
166 that of its sole genus and species. Type genus is *Candidatus* Borrarchaeum.

167 ***Candidatus Borrarchaeia*** (Borr.ar.chae'i.a. N.L. neut. n. Borrarchaeum a (*Candidatus*) type genus of the
168 order of the class; -ia ending to denote the class; N.L fem. pl. n. Borrarchaeia the Borrarchaeum class).

169 The class is delineated based on 209 concatenated AsCOGs phylogeny. The description is the same as
170 that of its sole and type order *Candidatus* Borrarchaeales.

171 ***Candidatus Borrarchaeota*** (Borr.ar.chae.o'ta. N.L. neut. n. Borrarchaeum a (*Candidatus*) type genus of
172 the class of the phylum; -ota ending to denote the phylum; N.L neut. pl. n. Borrarchaeota the
173 Borrarchaeum phylum)

174 The phylum is delineated based on 209 concatenated AsCOGs phylogeny. The description is the same as
175 that of its sole and type class *Candidatus* Borrarchaeia.

176

177 ***Candidatus Baldrarchaeaceae*** (Bal.dr.ar.chae.a.ce'ae. N.L. neut. n. Baldrarchaeum a (*Candidatus*) type
178 genus of the family; -aceae ending to denote the family; N.L. fem. pl. n. Baldrarchaeaceae the
179 Baldrarchaeum family).

180 The family is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
181 description is the same as that of its sole genus and species. Type genus is *Candidatus* Baldrarchaeum.

182 ***Candidatus Baldrarchaeales*** (Bal.dr.ar.chae.a'les. N.L. neut. n. Baldrarchaeum a (*Candidatus*) type
183 genus of the order; -ales ending to denote the order; N.L fem. pl. n. Baldrarchaeales the Baldrarchaeum
184 order).

185 The order is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
186 description is the same as that of its sole genus and species. Type genus is *Candidatus* Baldrarchaeum.

187 ***Candidatus Baldrarchaeia*** (Bal.dr.ar.chae'i.a. N.L. neut. n. Baldrarchaeum a (*Candidatus*) type genus of
188 the order of the class; -ia ending to denote the class; N.L fem. pl. n. Baldrarchaeia the Baldrarchaeum
189 class).

190 The class is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
191 description is the same as that of its sole and type order *Candidatus* Baldrarchaeales.

192 ***Candidatus Baldrarchaeota*** (Bal.dr.ar.chae.o'ta. N.L. neut. n. Baldrarchaeum a (*Candidatus*) type genus
193 of the class of the phylum; -ota ending to denote the phylum; N.L neut. pl. n. Baldrarchaeota the
194 Baldrarchaeum phylum)

195 The phylum is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
196 description is the same as that of its sole and type class *Candidatus* Baldrarchaeia.

197

198 ***Candidatus Hermodarchaeaceae*** (Her.mod.ar.chae.a.ce'ae. N.L. neut. n. Hermodarchaeum a
199 (*Candidatus*) type genus of the family; -aceae ending to denote the family; N.L. fem. pl. n.
200 Hermodarchaeaceae the Hermodarchaeum family).

201 The family is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
202 description is the same as that of its sole genus and species. Type genus is *Candidatus* Hermodarchaeum.

203 ***Candidatus* Hermodarchaeales** (Her.mod.ar.chae.a'les. N.L. neut. n. Hermodarchaeum a (Candidatus)
204 type genus of the order; -ales ending to denote the order; N.L fem. pl. n. Hermodarchaeales the
205 Hermodarchaeum order).

206 The order is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
207 description is the same as that of its sole genus and species. Type genus is *Candidatus* Hermodarchaeum.

208 ***Candidatus* Hermodarchaeia** (Her.mod.ar.chae'i.a. N.L. neut. n. Hermodarchaeum a (Candidatus) type
209 genus of the order of the class; -ia ending to denote the class; N.L fem. pl. n. Hermodarchaeia the
210 Hermodarchaeum class).

211 The class is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
212 description is the same as that of its sole and type order *Candidatus* Hermodarchaeales.

213 ***Candidatus* Hermodarchaeota** (Her.mod.ar.chae.o'ta. N.L. neut. n. Hermodarchaeum a (Candidatus)
214 type genus of the class of the phylum; -ota ending to denote the phylum; N.L neut. pl. n.
215 Hermodarchaeota the Hermodarchaeum phylum)

216 The phylum is delineated based on 209 concatenated AsCOGs and 16S rRNA gene phylogeny. The
217 description is the same as that of its sole and type class *Candidatus* Hermodarchaeia.

218

219 **Materials and Methods**

220 **Sampling collections and DNA sequencing**

221 YT samples were obtained from the Rongcheng Nation Swan Nature Reserve (Rongcheng, China) in
222 November 2018. The sediment cores were collected using columnar samplers at depth intervals of 0–2,
223 21–26, and 36–41 cm at a seagrass meadow and a non-seagrass-covered site nearby. After collection,
224 bulk sediments were immediately sealed in plastic bags, placed in a pre-cooled icebox, and transported to
225 laboratory within 4 hours. For each sample, DNA was extracted from 10 g sediment using PowerSoil
226 DNA Isolation kit (QIAGEN, Germany), according to the manufacturer's protocol. Following extraction,
227 nucleic acids were sequenced using Illumina HiSeq2500 (Illumina, USA) PE150 by Novogene (Nanjing,
228 China).

229 MP5 samples were obtained from Mai Po Nature Reserve (Hong Kong, China) in September 2014. Three
230 subsurface sediment samples were collected from a site covering with mangrove forest at depth intervals
231 of 0-2, 10-15 and 20-25 cm. Two subsurface sediment samples were taken at an intertidal mudflat with
232 depths of 0-5 and 13-16 cm. Samples were transported back to laboratory as described for YT
233 metagenomes. DNA was extracted from 5g wet sediment per sample using the PowerSoil DNA Isolation
234 Kit (MO BIO, USA) following the manufacturer's protocol. Metagenomic sequencing data were
235 generated using Illumina HiSeq2500 (Illumina, USA) PE150 by Novogene (Tianjin, China).

236 FT samples were taken from Futian Nature Reserve (Shenzhen, Guangdong, China) in April 2017.
237 Sediment samples were collected as described for YT samples at depth intervals of 0-2, 6-8, 12-14, 20-22,
238 and 28-30 cm. DNA was extracted from 5g wet sediment per sample using DNeasy PowerSoil kit
239 (Qiagen, Germany) as per manufacturer's instructions. Nucleic acids were sequenced using Illumina
240 HiSeq2000 (Illumina, USA) PE150 by Novogene (Tianjin, China).

241 The surface sediment sample of the CJE metagenome was collected from Changjiang estuary during a
242 cruise in August 2016. The sample was grabbed from the water bed, sealed immediately in a 50 ml tubes
243 and stored in liquid nitrogen onboard. After transportation to laboratory, 10g wet sediments were used for
244 DNA extraction as per manufacturer's protocol. Nucleic acids were sequenced using Illumina HiSeq2000
245 (Illumina, USA) PE150 by Novogene (Tianjin, China).

246 An oil sand sample was collected from Shengli Oilfield (Shandong, China) into bottles and they were
247 transported to laboratory where they were stored at 4 °C. The sample was used as inoculum to perform
248 enrichment with anaerobic medium in vials as described elsewhere (1). After 253d of enrichment, the
249 genomic DNA was extracted as described elsewhere (2). Nucleic acids were sequenced Illumina
250 HiSeq2000 (Illumina, USA) PE150 by Novogene (Tianjin, China).

251 The seawater samples of Yap metagenomes were collected at Yap trench region by CTD SBE911plus
252 (Sea-Bird Electronics, USA) during the 37th Dayang cruise in 2016. 8L of seawater per sample was
253 filtered through a 0.22 µm-mesh membrane filter immediately after recovery onboard. The membrane
254 was then cut into ~0.2 cm² pieces with a flame-sterilized scissors and added to a PowerBead Tube (MO
255 BIO, USA) and the subsequent steps were implemented according to the manufacturer's protocol to
256 extract DNA. The DNA per sample was amplified in five separate reactions using REPLI-g Single Cell
257 Kits (Qiagen, Germany) following the manufacturer's protocol, given to the challenging nature of sample
258 retrieval and DNA recovery. The products were pooled together and purified using QIAamp DNA Mini

259 Kit (Qiagen, Germany) according to the manufacturer's recommendations. Parallel blank controls were
260 set for sampling, DNA extraction and ampliation with 0.22 µm-mesh membrane filtering Milli-Q water
261 (18.2 MΩ; Millipore, USA). Nucleic acids were sequenced using HiSeq X Ten (Illumina, USA) PE150.

262 **Metagenomic assembly, binning and gene calling**

263 *FT, MP5 and YT metagenomes.* These three sets of metagenomes were assembled and binned using the
264 same method. Raw shotgun metagenomic sequencing reads were trimmed with “read_qc” module from
265 metaWRAP (v.1.1) (3). All clean reads from the same set were pooled together prior to *de novo* assemble
266 to one co-assembly. Clean reads were sent out to MEGAHIT (v1.1.2) with flag “--presets meta-large” for
267 co-assembling job (4). Sequencing coverage was determined for each assembled scaffold by mapping
268 reads from each sample to the co-assembly using Bowtie2 (5). The binning analysis was carried out 8
269 times with 8 different combinations of specificity and sensitivity parameters using MetaBAT2(52) (“--
270 maxP 60 or 95” AND “--minS 60 or 95” AND “--maxEdges 200 or 500”) on the assembly with a
271 minimum length of 2000 bp (6). DAS Tool (v1.0) was used as a dereplication and aggregation strategy on
272 those eight binning results to construct accurate bins (7). Manual curation was used for reducing the
273 genome contamination based on differential coverages, GC contents, and the presence of duplicate genes.

274 The depth coverage and N50 statistics of 38 Asgard MAGs recovered from YT metagenomes range from
275 7.72 to 298.86 (median: 21.77) and from 6658 to 381755 bps (median: 26337.5 bps) respectively; for 13
276 Asgard MAGs recovered from FT metagenomes, the values range from 6.98 to 33.22 (median 16.89) and
277 from 3889 to 8957 bps (median: 5000 bps), respectively; and for 11 Asgard MAGs recovered from MP5
278 metagenomes, the values range from 7.06 to 54.42 (median: 10) and from 3898 to 19362 bps (median:
279 7581 bps) respectively. (**Supplementary Figure 2**).

280 *CJE metagenome.* Raw metagenomic shotgun sequencing reads were trimmed using Sickle
281 (<https://github.com/najoshi/sickle>) with default settings. The trimmed reads were de novo assembled
282 using IDBA-UD (v 1.1.1) with the parameters: “-mink 65, -maxk 145, -step 10” (8). Sequencing coverage
283 was determined as described above. The binning analysis were performed with MetaBAT2 12 times, with
284 12 combinations of specificity and sensitivity parameters (“--m 1500, 2000, or 2500” AND “--maxP 85 or
285 90” AND “--minS 80, 85 and 90”) for further refinement (6). All binning results were merged and refined
286 using DAS Tool (v1.0) (7).

287 2 Asgard MAGs recovered from CJE metagenome have a depth coverage of 19.16 and 20.56, and a N50
288 statistics of 8740 and 5246 bps (**Supplementary Figure 2**).

289 *J65 metagenome.* Raw metagenomic shotgun sequencing reads were trimmed with Trimmomatic (v0.38)
290 (9). The clean reads were then fed to SPAdes (v 3.12.0) for *de novo* assembly with the parameters: “--
291 meta -k 21, 33, 55, 77” (10). Sequencing coverage was determined using BBMap (v 38.24) toolkit with
292 the parameters: “bbmap.sh minid=0.99” (<https://github.com/BioInfoTools/BBMap>). MetaBAT2 (v 2.12.1)
293 was used to perform binning analysis with the parameter: “-m 2000” (6).

294 One Asgard MAG recovered from J65 metagenome had a depth coverage of 14.52 and a N50 statistics of
295 10460 bps (**Supplementary Figure 2**).

296 *Yap metagenome*. For each Yap metagenome, raw metagenomic shotgun sequencing reads were trimmed
297 with Trimmomatic (v0.38) (9). Assembly and binning analysis were performed as described for CJE
298 metagenome for each Yap metagenome.

299 The depth coverage and N50 statistics of 24 Asgard MAGs recovered from Yap metagenomes range from
300 7.78 to 82.33 (median: 13.91) and from 5097 to 889102 bps (median: 18155.5 bps) respectively.
301 **(Supplementary Figure 2)**

302 A total of 89 Asgard MAGs were reconstructed in this work. Additional 95 Asgard MAGs were
303 downloaded from public databases (e.g. NCBI FTP site). For all 184 genomes, a uniform gene calling
304 protocol was applied. Specifically, the completeness, contamination, and strain heterogeneity of the
305 genomes were estimated by using CheckM (v.1.0.12) (11) and DAS Tool under the taxonomic scope of
306 domain (i.e., Bacteria and Archaea). Protein-coding genes were predicted using Prodigal (v 2.6.3) (12)
307 embedded in Prokka (v 1.13) (13). Transfer RNAs (tRNAs) were identified with tRNAscan-SE (v1.23)
308 using the archaeal tRNA model (14). After quality screening, further analysis focused on 162 high quality
309 Asgard MAGs.

310 **Genome set**

311 The Asgard archaea genome set analyzed in this work consisted of 161 Asgard MAGs and one complete
312 Asgard genome (**Supplementary Table 1**). For comparison, selected representative genomes of archaea
313 (296), bacteria (100) and eukaryotes (76) were downloaded from Refseq and Genbank using the NCBI
314 FTP site (<ftp://ftp.ncbi.nlm.nih.gov/genomes/all/>).

315 **Average amino-acid identity**

316 The average amino-acid identity (AAI) across TACK archaeal reference genomes and the 184 Asgard
317 genomes was calculated using compareM (v0.0.23) with the “aai_wf” at default settings
318 (<https://github.com/dparks1134/CompareM>).

319 **asCOGs construction**

320 Initial clustering of 250,634 proteins encoded in 76 Asgard MAGs was performed using two approaches:
321 first, footprints of arCOG profiles were obtained by running PSI-BLAST (15), initiated with arCOG
322 alignments, against the set of predicted Asgard proteins. The footprint sequences were extracted and
323 clustered according to the arCOG best hit. The remaining protein sequences (both full-length proteins and
324 the sequence fragments outside of the footprints, if longer than 60 aa) were clustered using MMseqs2
325 (16), with similarity threshold of 0.5. Sequences within clusters were aligned using MUSCLE (17); the
326 resulting alignments were passed through several rounds of merging and splitting. The merging phase
327 involved comparing alignments to each other using HHSEARCH (18), finding full-length cluster-to-
328 cluster matches, merging the sequence sets and re-aligning the new clusters. The splitting phase consisted
329 of the construction of an approximate phylogenetic tree of the sequences using FastTree (19) (gamma-
330 distributed site rates, WAG evolutionary model) with balanced mid-point tree rooting, identification of
331 subtrees maximizing the fraction of species (MAGs) representation and minimizing the number of
332 paralogs, and pruning such subtrees as separate clusters of putative orthologs. Clusters, derived from
333 arCOGs, were prohibited from merging across distinct arCOGs to prevent distant paralogs from forming
334 mixed clusters.

335 **Phylogenetic analysis**

336 *16S rRNA gene phylogenetic analysis.* 16S rRNA gene sequences were identified in 73 genomes of
337 Asgard archaea (26 generated in this study, 47 from public database) using Barrnap (v 0.9) with the "--
338 kingdom arc" option (<https://github.com/tseemann/barrnap>). These sequences were combined with 46
339 published 16S rRNA sequences of Asgard archaea, to assess the novelty of the sequences obtained in this
340 work. The novelty of 16S rRNA gene sequences was measured in terms of their sequence identity to
341 previously identified Asgard archaeal 16S rRNA gene and phylogenetic relationships. Specifically, the
342 pairwise sequence identity of two Asgard archaeal 16S rRNA gene sequences (>1300bp) was obtained by
343 first globally aligning the sequences with "Stretcher" in EMBOSS package and then calculating the
344 percent identity excluding gaps (20). The Asgard archaeal 16S rRNA gene sequences were aligned with
345 311 reference sequences from Euryarchaeota (n=231), DPANN (n=22), Korarchaeota (n=1),
346 Crenarchaeota (n=41), Bathyarchaeota (n=2), Thaumarchaeota (n=13) and Aigarchaeota (n=1) using
347 mafft-linsi (v7.471) (21) and trimmed with BMGE (v1.12) (settings: -m DNAPAM250:4, -g 0.5) (22). The
348 alignment was used for phylogenetic inference with IQ-Tree (v2.0.6) based on SYM+R8 (selected by
349 ModelFinder) to generate a maximum-likelihood tree (23).

350 *Asgard phylogeny.* A set of asCOGs that were considered most suitable as phylogenetic markers for
351 Asgard archaea was selected using the preliminary classification of the 76 genomes in AsCOGs into
352 previously described lineages: Loki, Thor, Odin, Hel and Heimdall. The following criteria were adopted:
353 the asCOG have to be i) present in at least half of the genomes in all lineages, ii) present in at least 75%
354 among the 76 genomes, iii) the mean number of paralogs per genome not to exceed 1.25. For the 209
355 asCOGs matching these criteria, the corresponding protein sequences were obtained from the extended set
356 of 162 MAGs and aligned using MUSCLE (17); the 'index' paralog to include in the phylogeny was
357 selected for each MAG based on the similarity to the alignment consensus. Alignments were trimmed to
358 exclude columns containing more than 2/3 gap characters and with homogeneity below 0.05
359 (homogeneity is calculated from the score of the consensus amino acid against the alignment column,
360 compared to the score of the perfect match (24) and concatenated, resulting in an alignment of 50,706
361 characters from 162 sequences (one sequence per MAG). Phylogeny was reconstructed using FastTree
362 (gamma-distributed site rates, WAG evolutionary model) (19) and IQ-Tree (LG+F+R10 model, selected
363 by ModelFinder) (23), producing very similar tree topologies.

364 *Tree of life.* To elucidate the relationships between the Asgard archaea and other major clades of archaea,
365 bacteria and eukaryotes, 30 families of conserved proteins were selected that appear to have evolved
366 mostly vertically (25, 26). The set of 162 MAGs of Asgard archaea was supplemented with 66 TACK
367 archaea and 220 non-TACK archaea (the former, having been described as the closest archaeal relatives
368 of the Asgard archaea (27), were sampled more densely), 98 bacteria and 72 eukaryotes. Prokaryotic
369 genomes were sampled from the set of completely sequenced genomes to represent the maximum
370 diversity within their respective clades (briefly, all proteins, encoded by genomes within a groups were
371 clustered at 75% identity level; distances between genomes were estimated from the number of shared
372 proteins within these clusters; UPGMA trees were reconstructed from the distances, and genome sets
373 maximizing the total tree branch length were selected to represent the groups). The set of eukaryotic
374 genomes was manually selected to represent the maximum possible variety of eukaryotic taxa. Genomes,
375 in which more than four markers were missing, were excluded from the bacterial and eukaryotic sets.
376 When multiple paralogs of a marker were present in a genome, preliminary phylogenetic trees were

377 constructed from protein sequence alignments, and paralogs with the shortest branches were selected to
378 represent the corresponding genomes in the set. Sequences were aligned using MUSCLE (17), and
379 alignment columns, containing more than 2/3 gap characters or with alignment column homogeneity
380 below 0.05 were removed (24). The resulting concatenated alignments of the 30 markers consisted of
381 7411 sites. Phylogeny was reconstructed using FastTree (gamma-distributed site rates, WAG evolutionary
382 model) (19) and IQ-Tree (23) with three models: LG+R10, selected by IQ-tree ModelFinder as the best
383 fit, GTR20+F+R10 (following the suggestion of Williams et al. (28) to use GTR, we let IQ-tree to select
384 the best version of the GTR model), and LG+C20+G4+F (again, the mixture model was used following
385 the suggestion of Williams et al. (28)); we were unable to use the higher-specified C60 model due to
386 memory limitations of our hardware and used the C20 model instead).

387 **Ordination of asCOG phyletic patterns using Classical Multidimensional Scaling**

388 Binary asCOG presence-absence patterns were compared between pairs of Asgard MAGs using the
389 following procedure: first, similarity between asCOG sets $\{A\}$ and $\{B\}$ was calculated as $S_{A,B} =$
390 $|A \cap B|/\sqrt{|A||B|}$ (the number of shared AsCOGs normalized by the geometric mean of the number of
391 AsCOGs in the two MAGs); then, the distance between the patterns was calculated as $d_{A,B} = -\ln(S_{A,B})$.
392 The 162x162 distance matrix was embedded into a 2-dimensional space using Classical Multidimensional
393 Scaling analysis implemented as the *cmdscale* function in R. The projection retained 89% of the original
394 datapoint inertia.

395 **Identification and analysis of ESPs**

396 To identify eukaryotic signature proteins (ESPs), several strategies were employed. First, ESPs reported
397 by Zaremba-Niedzwiedzka et al. (29) were mapped to asCOGs using PSI-BLAST (15). These asCOGs
398 were additionally examined case by case using HHpred (30) with representative of the respective asCOG
399 or the respective asCOG alignment used as the query. Second, all asCOGs were mapped to CDD profiles
400 (31) using PSI-BLAST, hits to eukaryote-specific domains were selected. Most of the putative ESP
401 asCOGs identified in this search, and all with E-value $>1e-10$, were additionally examined using HHpred,
402 with a representative of the respective asCOG or the respective asCOG alignment as the query. Third, we
403 analyzed frequently occurring asCOGs (present in at least 50% of Asgard genomes and in at least 30% of
404 Heimdall genomes) that were not annotated automatically with the above two approaches using HHpred
405 with a representative of the respective asCOG or the respective asCOG alignment as the query. Fourth,
406 most of the putative ESP asCOGs detected with these approaches were used as queries for a PSI-BLAST
407 search that was run for three iterations (with inclusion threshold E-value=0.0001) against an Asgard only
408 protein sequence database. Additional unannotated asCOGs with similarity to the (putative) ESPs
409 identified in this search were further examined using HHpred. Fifth, the genomic neighborhoods or all
410 ESPs were examined, and proteins encoded by unannotated neighbor genes were analyzed using HHpred
411 server.

412 **Metabolic pathway reconstruction**

413 The patterns of gene presence-absence in the asCOGs were used to reconstruct the metabolic pathways of
414 Asgard archaea (Supplementary Table 8). The asCOGs were linked to the KEGG database and to the list
415 of predicted metabolic enzymes of Asgard archaea reported by Spang et al. (32) (Supplementary Table 8).
416 The classification of [NiFe] hydrogenases was performed by comparing the Asgard proteins of

417 cog.001539, cog.002254, cog.010021, cog.011939 and cog.012499 to HydDB (33). For phylogenetic
418 analysis, the reference sequences of group 1, 3 and 4 [NiFe] hydrogenases were retrieved from HydDB.
419 The sequences were filtered using cd-hit with a sequence identity cut-off of 90% prior to adding
420 orthologous genes of cog.001539, cog.002254, cog.010021, cog.011939 and cog.012499 of Asgard
421 archaea. All sequences for group 1, group 3 and group 4 [NiFe] hydrogenases were aligned using mafft-
422 LINSI (21) and trimmed with BMGE (-m BLOSUM30 -h 0.6) (22). Maximum-likelihood phylogenetic
423 analyses were performed using IQ-tree (23) with the best-fit model (group 1:LG + C60 + R + F, group 3:
424 LG + C60 + R + F and group 4 LG + C50 + R + F), respectively, according to Bayesian information
425 criterion (BIC). Support values were estimated using the SH-like approximate-likelihood ratio test and
426 ultrafast bootstraps, respectively. We adapted a relaxed common denominator approach to determine the
427 presence of certain pathway in one Asgard phyla (32), and combined with maximum parsimony principle
428 (34) to infer the metabolisms of major ancestral forms.

429

430

431 References

- 432 1. L. Cheng, T.-L. Qiu, X. Li, W.-D. Wang, Y. Deng, X.-B. Yin, H. Zhang, Isolation and
433 characterization of *Methanoculleus receptaculi* sp. nov. from Shengli oil field, China. *FEMS*
434 *Microbiology Letters*. **285**, 65–71 (2008).
- 435 2. J. Peng, Z. Lü, J. Rui, Y. Lu, Dynamics of the Methanogenic Archaeal Community during Plant
436 Residue Decomposition in an Anoxic Rice Field Soil. *Appl. Environ. Microbiol.* **74**, 2894–2901 (2008).
- 437 3. G. V. Uritskiy, J. DiRuggiero, J. Taylor, MetaWRAP—a flexible pipeline for genome-resolved
438 metagenomic data analysis. *Microbiome*. **6**, 158 (2018).
- 439 4. D. Li, C. M. Liu, R. Luo, K. Sadakane, T. W. Lam, MEGAHIT: an ultra-fast single-node solution
440 for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*. **31**, 1674–
441 1676 (2015).
- 442 5. B. Langmead, S. L. Salzberg, Fast gapped-read alignment with Bowtie 2. *Nat Methods*. **9**, 357–
443 359 (2012).
- 444 6. D. D. Kang, F. Li, E. Kirton, A. Thomas, R. Egan, H. An, Z. Wang, MetaBAT 2: an adaptive
445 binning algorithm for robust and efficient genome reconstruction from metagenome assemblies. *PeerJ*. **7**,
446 e7359 (2019).
- 447 7. C. M. K. Sieber, A. J. Probst, A. Sharrar, B. C. Thomas, M. Hess, S. G. Tringe, J. F. Banfield,
448 Recovery of genomes from metagenomes via a dereplication, aggregation and scoring strategy. *Nat*
449 *Microbiol.* **3**, 836–843 (2018).
- 450 8. Y. Peng, H. C. M. Leung, S. M. Yiu, F. Y. L. Chin, IDBA-UD: a de novo assembler for single-
451 cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics*. **28**, 1420–1428 (2012).
- 452 9. A. M. Bolger, M. Lohse, B. Usadel, Trimmomatic: a flexible trimmer for Illumina sequence data.
453 *Bioinformatics*. **30**, 2114–2120 (2014).
- 454 10. S. Nurk, D. Meleshko, A. Korobeynikov, P. A. Pevzner, metaSPAdes: a new versatile
455 metagenomic assembler. *Genome Res*. **27**, 824–834 (2017).
- 456 11. D. H. Parks, M. Imelfort, C. T. Skennerton, P. Hugenholtz, G. W. Tyson, CheckM: assessing the
457 quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Research*.
458 **25**, 1043–1055 (2015).
- 459 12. D. Hyatt, G.-L. Chen, P. F. LoCascio, M. L. Land, F. W. Larimer, L. J. Hauser, Prodigal:
460 prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics*. **11**, 119
461 (2010).
- 462 13. T. Seemann, Prokka: rapid prokaryotic genome annotation. *Bioinformatics*. **30**, 2068–2069
463 (2014).
- 464 14. P. P. Chan, T. M. Lowe, tRNAscan-SE: Searching for tRNA genes in genomic sequences.
465 *Methods Mol Biol*. **1962**, 1–14 (2019).

- 466 15. A. A. Schäffer, L. Aravind, T. L. Madden, S. Shavirin, J. L. Spouge, Y. I. Wolf, E. V. Koonin, S.
467 F. Altschul, Improving the accuracy of PSI-BLAST protein database searches with composition-based
468 statistics and other refinements. *Nucleic Acids Res.* **29**, 2994–3005 (2001).
- 469 16. M. Steinegger, J. Söding, MMseqs2 enables sensitive protein sequence searching for the analysis
470 of massive data sets. *Nat Biotechnol.* **35**, 1026–1028 (2017).
- 471 17. R. C. Edgar, MUSCLE: a multiple sequence alignment method with reduced time and space
472 complexity. *BMC Bioinformatics.* **5**, 113 (2004).
- 473 18. J. Söding, Protein homology detection by HMM-HMM comparison. *Bioinformatics.* **21**, 951–960
474 (2005).
- 475 19. M. N. Price, P. S. Dehal, A. P. Arkin, FastTree 2--approximately maximum-likelihood trees for
476 large alignments. *PLoS One.* **5**, e9490 (2010).
- 477 20. P. Rice, I. Longden, A. Bleasby, EMBOSS: the European Molecular Biology Open Software
478 Suite. *Trends Genet.* **16**, 276–277 (2000).
- 479 21. K. Katoh, D. M. Standley, MAFFT Multiple Sequence Alignment Software Version 7:
480 Improvements in Performance and Usability. *Molecular Biology and Evolution.* **30**, 772–780 (2013).
- 481 22. A. Criscuolo, S. Gribaldo, BMGE (Block Mapping and Gathering with Entropy): a new software
482 for selection of phylogenetic informative regions from multiple sequence alignments. *BMC Evolutionary*
483 *Biology.* **10**, 210 (2010).
- 484 23. L.-T. Nguyen, H. A. Schmidt, A. von Haeseler, B. Q. Minh, IQ-TREE: A Fast and Effective
485 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol Biol Evol.* **32**, 268–274
486 (2015).
- 487 24. E. Esterman, Y. I. Wolf, R. Kogay, E. V. Koonin, O. Zhaxybayeva, *bioRxiv*, in press,
488 doi:10.1101/2020.10.08.331884.
- 489 25. F. D. Ciccarelli, T. Doerks, C. von Mering, C. J. Creevey, B. Snel, P. Bork, Toward automatic
490 reconstruction of a highly resolved tree of life. *Science (New York, N.Y.).* **311**, 1283–1287 (2006).
- 491 26. P. Puigbò, Y. I. Wolf, E. V. Koonin, Search for a “Tree of Life” in the thicket of the phylogenetic
492 forest. *Journal of biology.* **8**, 59 (2009).
- 493 27. A. Spang, J. H. Saw, S. L. Jørgensen, K. Zaremba-Niedzwiedzka, J. Martijn, A. E. Lind, R. van
494 Eijk, C. Schleper, L. Guy, T. J. G. Ettema, Complex archaea that bridge the gap between prokaryotes and
495 eukaryotes. *Nature.* **521**, 173–179 (2015).
- 496 28. T. A. Williams, C. J. Cox, P. G. Foster, G. J. Szöllösi, T. M. Embley, Phylogenomics provides
497 robust support for a two-domains tree of life. *Nat Ecol Evol.* **4**, 138–147 (2020).
- 498 29. K. Zaremba-Niedzwiedzka, E. F. Caceres, J. H. Saw, D. B. Eckström, L. Juzokaite, E. Vancaester,
499 K. W. Seitz, K. Anantharaman, P. Starnawski, K. U. Kjeldsen, M. B. Stott, T. Nunoura, J. F. Banfield, A.

- 500 Schramm, B. J. Baker, A. Spang, T. J. G. Ettema, Asgard archaea illuminate the origin of eukaryotic
501 cellular complexity. *Nature*. **541**, 353–358 (2017).
- 502 30. L. Zimmermann, A. Stephens, S.-Z. Nam, D. Rau, J. Kübler, M. Lozajic, F. Gabler, J. Söding, A.
503 N. Lupas, V. Alva, A Completely Reimplemented MPI Bioinformatics Toolkit with a New HHpred
504 Server at its Core. *J Mol Biol*. **430**, 2237–2243 (2018).
- 505 31. M. Yang, M. K. Derbyshire, R. A. Yamashita, A. Marchler-Bauer, NCBI’s Conserved Domain
506 Database and Tools for Protein Domain Analysis. *Curr Protoc Bioinformatics*. **69**, e90 (2020).
- 507 32. A. Spang, C. W. Stairs, N. Dombrowski, L. Eme, J. Lombard, E. F. Caceres, C. Greening, B. J.
508 Baker, T. J. G. Ettema, Proposal of the reverse flow model for the origin of the eukaryotic cell based on
509 comparative analyses of Asgard archaeal metabolism. *Nature Microbiology*. **4**, 1138–1148 (2019).
- 510 33. D. Søndergaard, C. N. S. Pedersen, C. Greening, HydDB: A web tool for hydrogenase
511 classification and analysis. *Scientific Reports*. **6**, 34212 (2016).
- 512 34. D. L. Swofford, W. P. Maddison, Reconstructing ancestral character states under Wagner
513 parsimony. *Mathematical Biosciences*. **87**, 199–229 (1987).

514

515

516 **Data availability**

517 Asgard archaea genomes generated in this study are available in eLMSG (an eLibrary of Microbial
518 Systematics and Genomics, <https://www.biosino.org/elmsg/index>) under accession numbers from
519 LMSG_G000000521.1 to LMSG_G000000609.1 and NCBI database under the project XXXXX.

520 Additional Data files are available at ftp://ftp.ncbi.nih.gov/pub/wolf/_suppl/asgard20/

521 Additional data file 1 (Additional_data_file_1.tgz): Complete asCOG data archive

522 Additional data file 2 (Additional_data_file_2.tgz): Phylogenetic trees and alignments archive

523

524

525

526 **Supplementary Figures**

527 **Supplementary Figure 1** Global distribution of Asgard genomes analyzed in the study. The pie chart
528 shows the proportion of Asgard genomes found in a biotope. Bold letters in the map show the sampling
529 locations.

530 **Supplementary Figure 2** a. Distribution of completeness and contamination for 76 Asgard MAGs
531 assessed by CheckM (v 1.0.12). Distribution of depth coverage (b) and N50 statistics (c) for Asgard
532 MAGs reconstructed in this study. The numbers in parentheses indicate the number of Asgard genomes
533 recovered from a given sampling location. The data for this plot can be found in **Supplementary Table 1**.

534 **Supplementary Figure 3** Comparison of the mean amino-acid identity (AAI) of Asgard and TACK
535 superphyla. a. Shared AAI across Asgard and TACK lineages. Background: comparing all Asgard and
536 TACK lineages included in the analyses but excluding archaea belonging to the same lineages and the six
537 phyla proposed in the current work to investigate the distribution of AAI that defines a phylum. The AAI
538 comparison of (b) Thorarchaeota, (c) Hermodarchaeota (d) Odinarchaeota, (e) Baldrarchaeota, (f)
539 Lokiarchaeota, (g) Helarchaeota, (h) Borrarchaeota, (i) Heimdallarchaeota, (j) Kariarchaeota, (k)
540 Gerdarchaeota, (l) Hodarchaeota and (m) Wukongarchaeota to other Asgard and TACK lineages. The
541 lower and upper hinges of the boxplot correspond to the first and third quartiles. Data beyond the
542 whiskers are shown as individual data points. Number in the parenthesis indicates the number of genomes
543 in the lineages. Data for this plot is included in **Supplementary Table 2**.

544 **Supplementary Figure 4** Comparison of the 16S rRNA gene sequence (>1300 bp) identity of Asgard
545 and TACK lineages. a. Shared 16S rRNA gene sequence identity across Asgard and TACK lineages.
546 Background: comparing all Asgard and TACK lineages included in the analyses but excluding archaea
547 belonging to the same lineages and the six phyla proposed in the current study to investigate the
548 distribution of 16S rRNA gene identity that defines a phylum. 16S rRNA gene identity comparison of (b)
549 Thorarchaeota, (c) Hermodarchaeota (d) Odinarchaeota, (e) Lokiarchaeota, (f) Helarchaeota, (g)
550 Heimdallarchaeota, (h) Kariarchaeota, (i) Gerdarchaeota, (j) Hodarchaeota and (k) Wukongarchaeota to
551 other Asgard and TACK lineages. The lower and upper hinges of the boxplot correspond to the first and
552 third quartiles. Data beyond the whiskers are shown as individual data points. Number in the parenthesis
553 indicates the number of 16S rRNA gene sequences compared in the lineages. Line represents a 16S rRNA
554 gene identity of 75%. Data for this plot could be found in **Supplementary Table 3**.

555 **Supplementary Figure 5** Presence-absence of orthologs of Asgard core genes in other archaea, bacteria
556 and eukaryotes.

557 **Supplementary Figure 6** Phylogenetic analysis of group 4 [NiFe] hydrogenases in the Asgard archaea.
558 The unrooted maximum-likelihood phylogenetic tree was built from an alignment of 425 sequences
559 including 110 sequences of Asgard archaea, with 308 amino-acid positions.

560 **Supplementary Figure 7** Phylogenetic analysis of group 3 [NiFe] hydrogenases in the Asgard archaea.
561 The unrooted maximum-likelihood phylogenetic tree was built from an alignment of 813 sequences
562 including 335 sequences of Asgard archaea, 514 amino-acid positions.

563 **Supplementary Figure 8** Phylogenetic analysis of group 1 [NiFe] hydrogenases in the Asgard archaea.
564 The unrooted maximum-likelihood phylogenetic tree was built from an alignment of 541 sequences
565 including 2 sequences of Wukongarchaeota, with 745 amino-acid positions.

566 **Supplementary Figure 9** A schematic representation of the presence and absence of selected metabolic
567 features in all (putative) phyla of Asgard archaea.

568 **Supplementary Figure 10** Gene structure of the contig encoding Group 1 [Ni,Fe]-hydrogenase in
569 Wukongarchaeota. Abbreviations: Ftr, Formylmethanofuran:tetrahydromethanopterin formyltransferase;
570 FwdD Formylmethanofuran dehydrogenase subunit D; FwdB, Formylmethanofuran dehydrogenase
571 subunit B; FwdA, Formylmethanofuran dehydrogenase subunit A; FwdC, Formylmethanofuran
572 dehydrogenase subunit C; HdrC, Heterodisulfide reductase, subunit C; HyaB, Ni,Fe-hydrogenase I large
573 subunit; HyaA, Ni,Fe-hydrogenase I small subunit.

574

575

576 **Supplementary Tables**

577 **Supplementary Table 1** Genome information, proposed taxonomy and isolation data

578 **Supplementary Table 2** Mean amino-acid identity values (in %) comparing 66 TACK genomes and 184
579 Asgard genomes (162 high quality and 22 low-quality)

580 **Supplementary Table 3** The 16S rRNA gene sequence identity (in %) comparing TACK lineages and
581 Asgard lineages. The identity was calculated using sequences longer than 1300 bps

582 **Supplementary Table 4** Species and phyletic markers used for the tree of life reconstruction

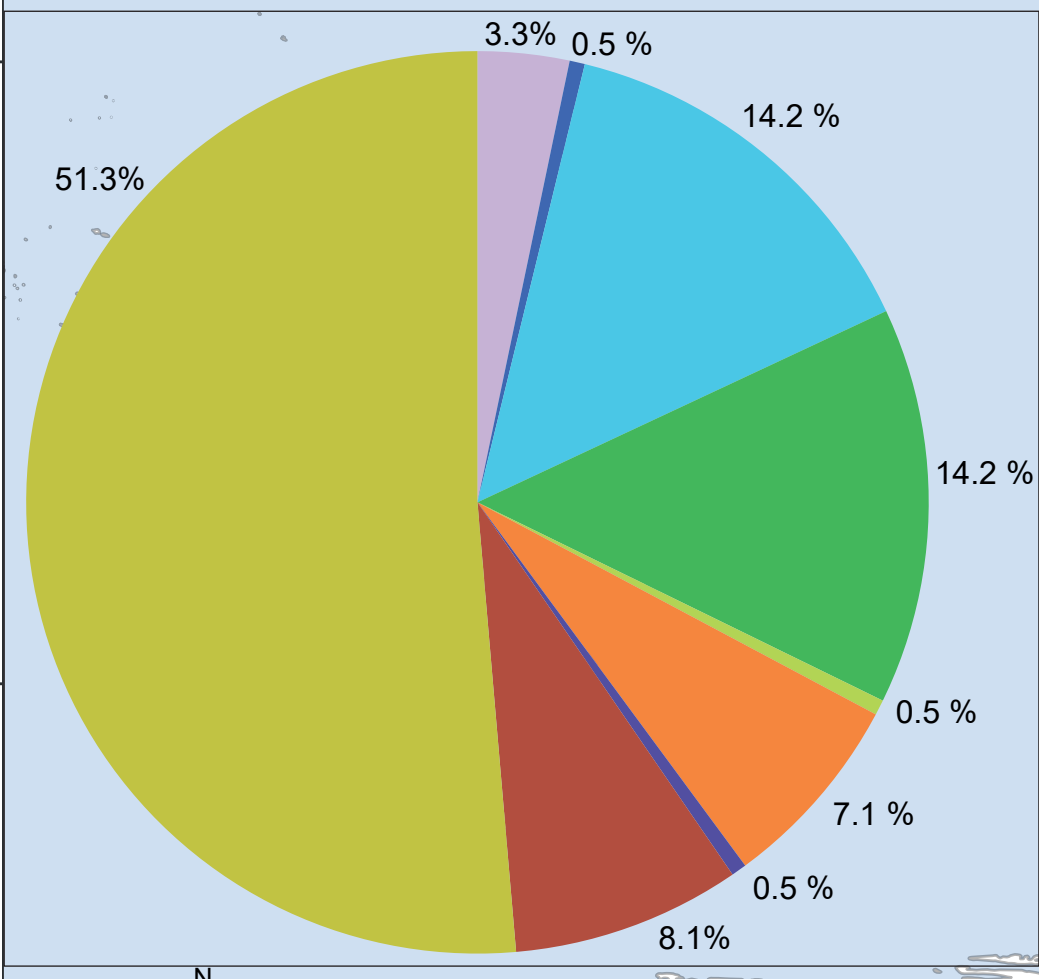
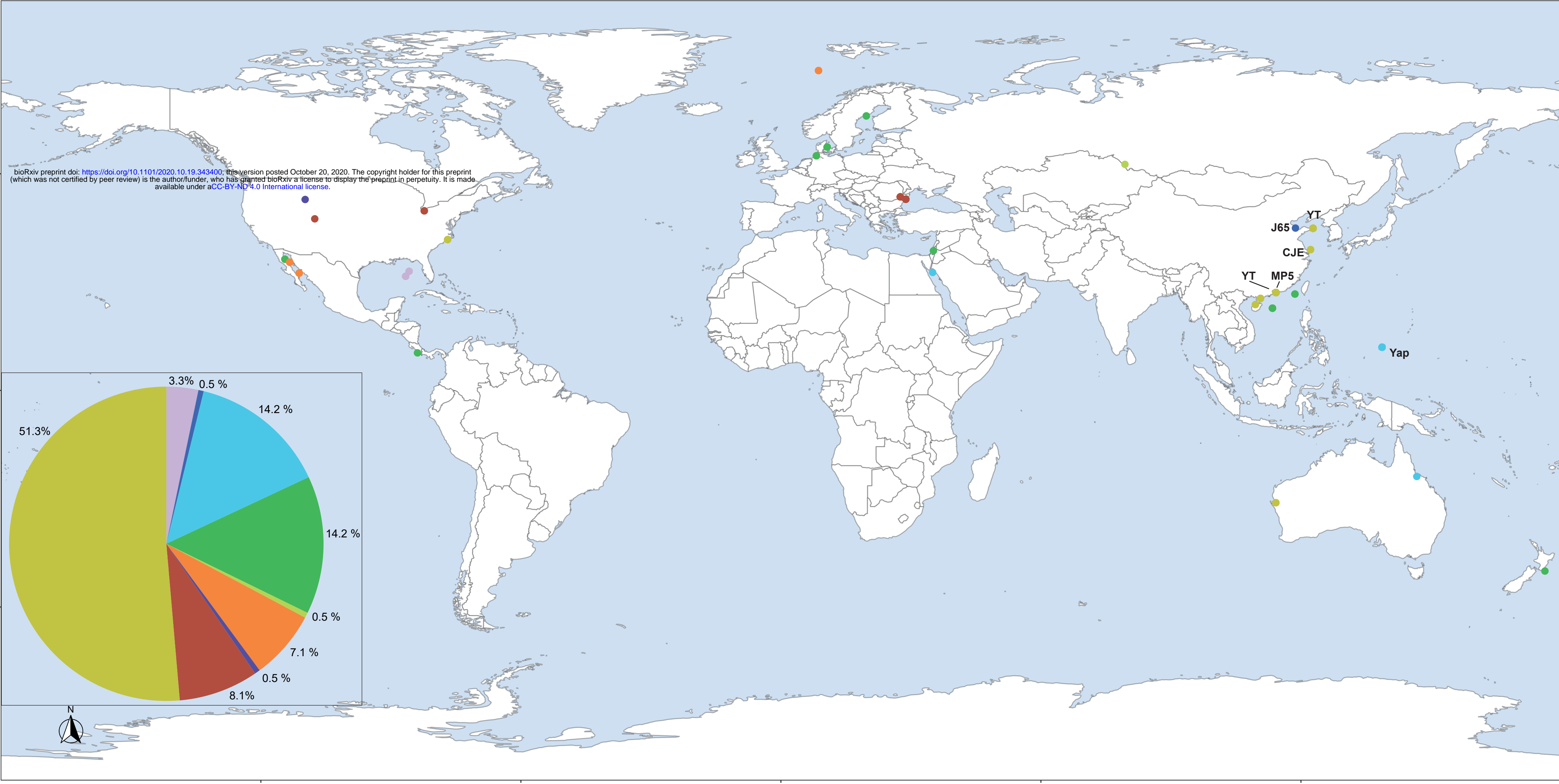
583 **Supplementary Table 5** Data for phylogenetic trees: the trees in the Newick format and the underlying
584 alignments

585 **Supplementary Table 6** The asCOGs annotation

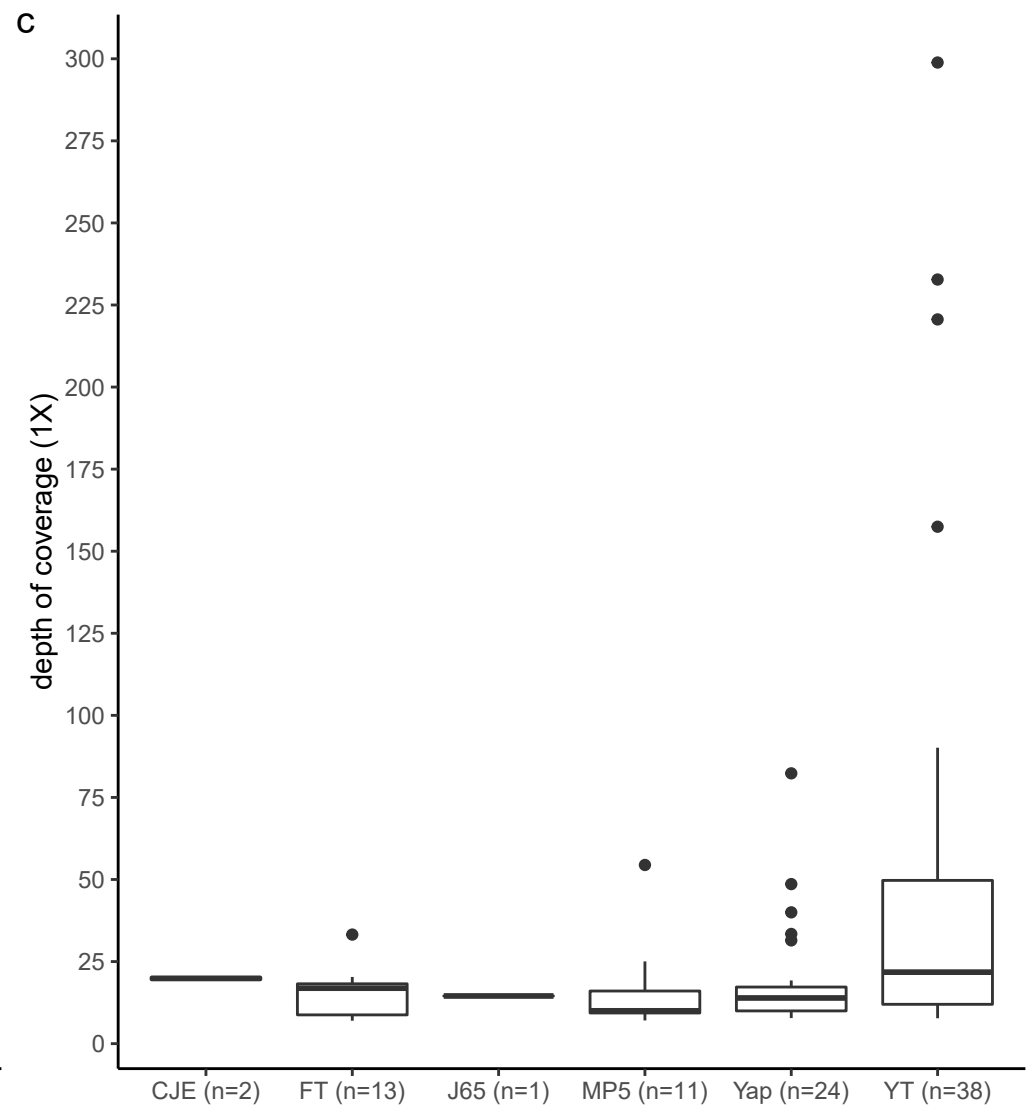
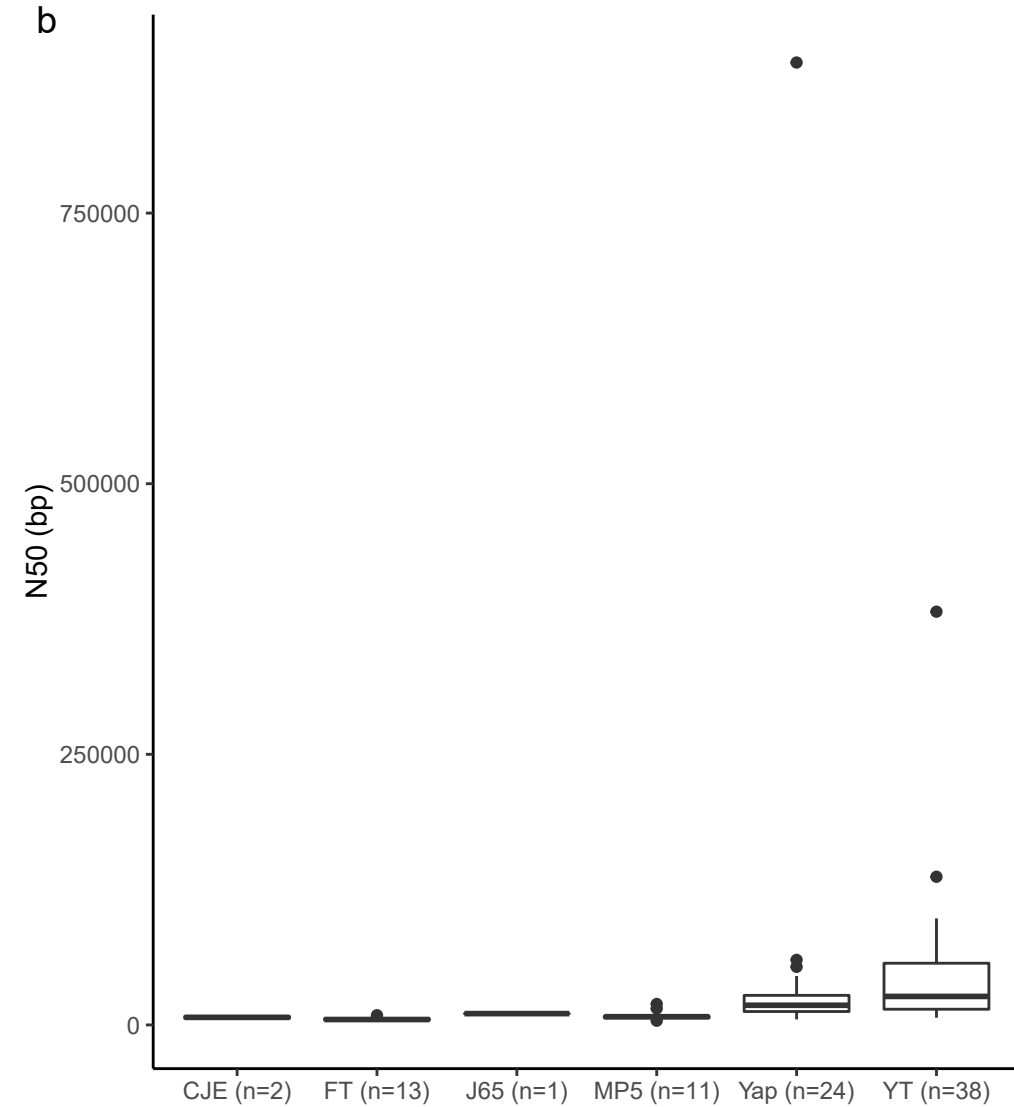
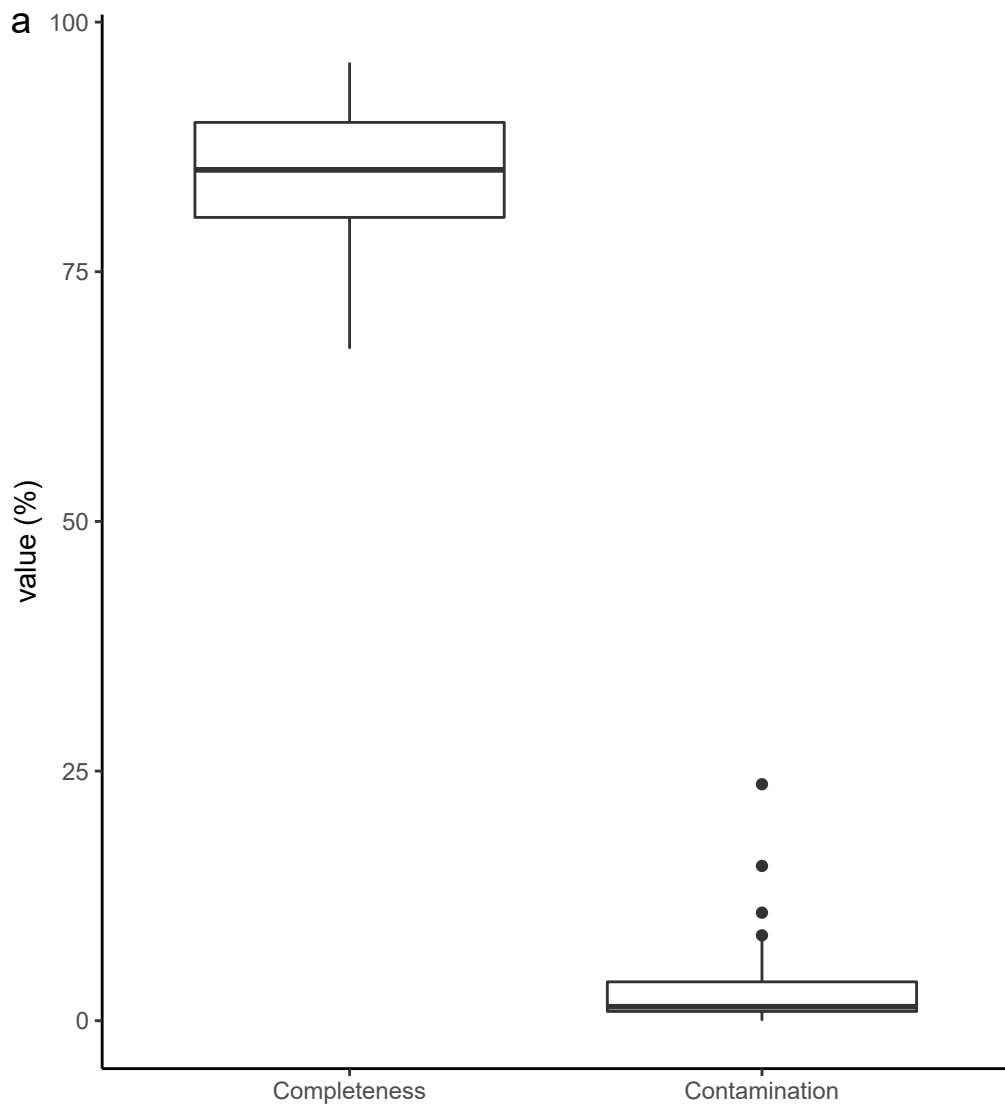
586 **Supplementary Table 7** Eukaryotic signature proteins in Asgard archaea

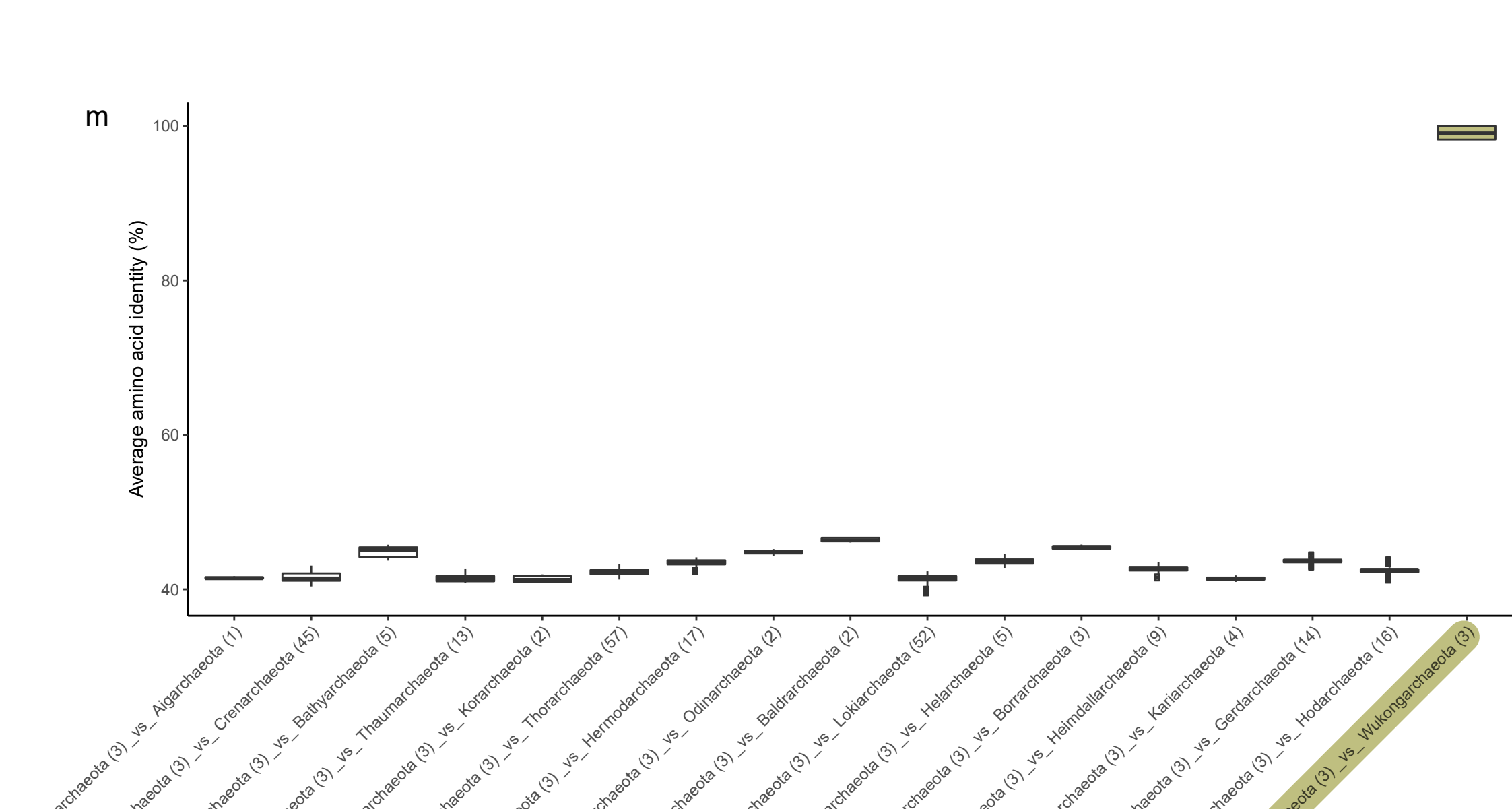
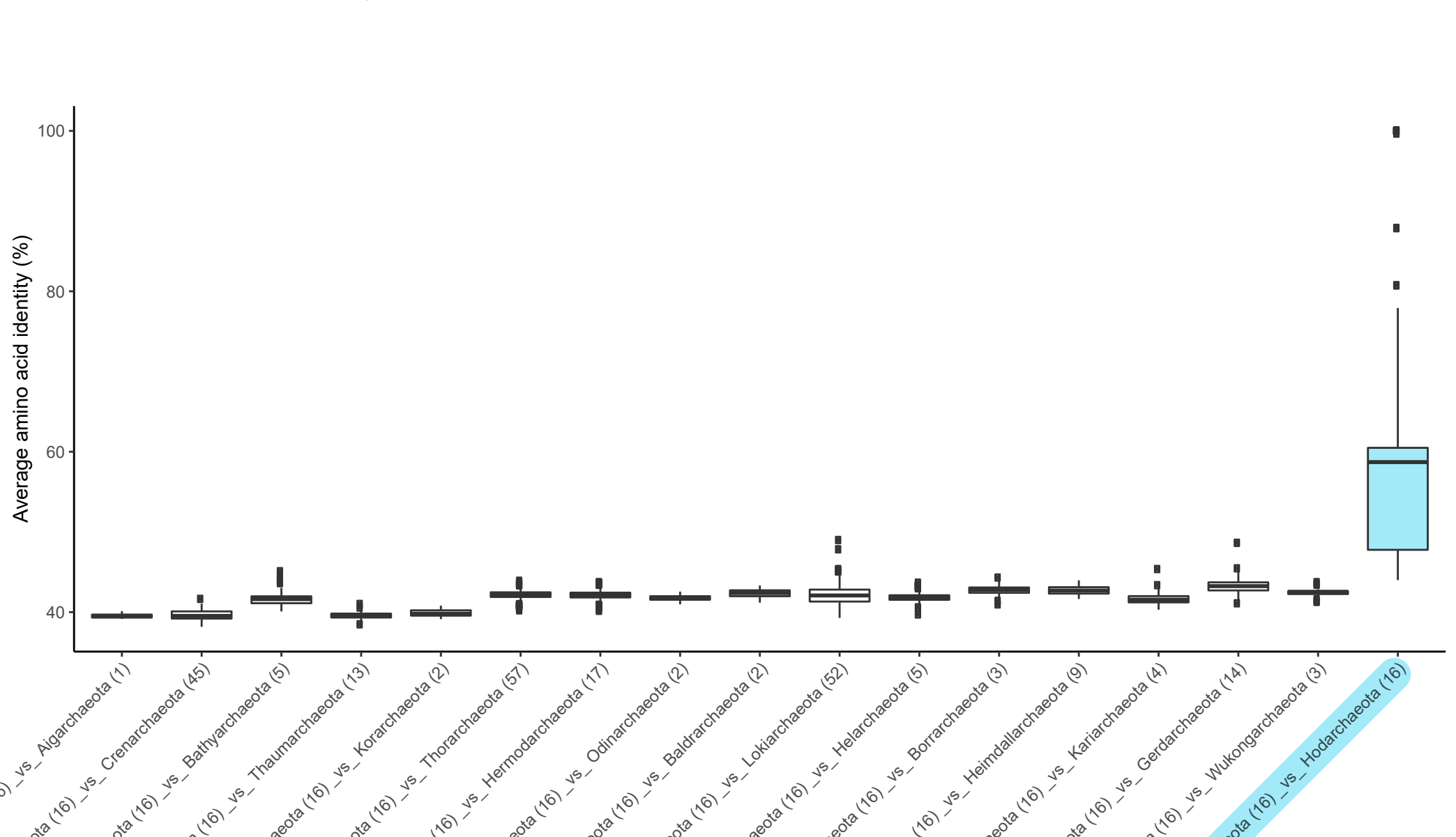
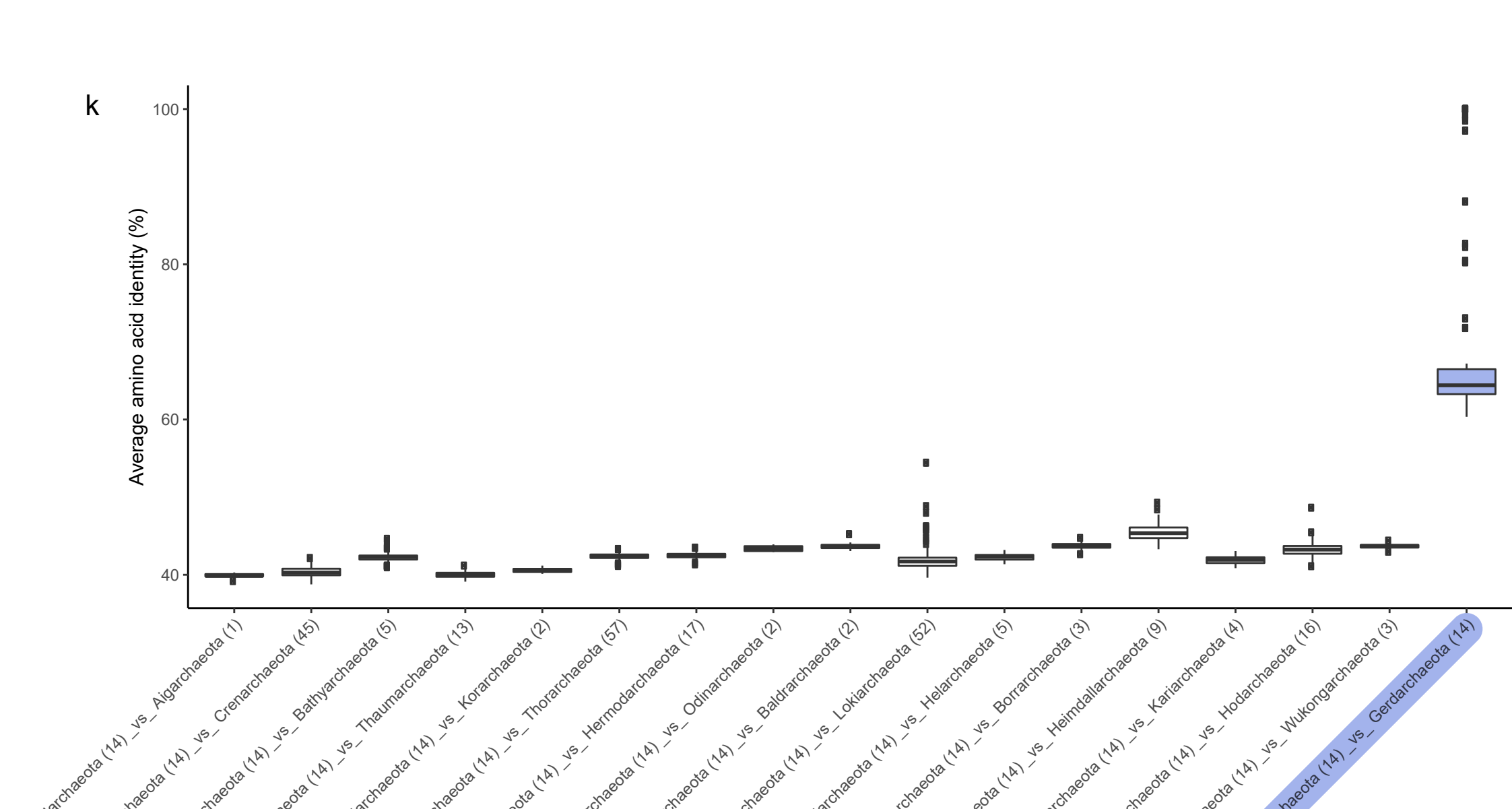
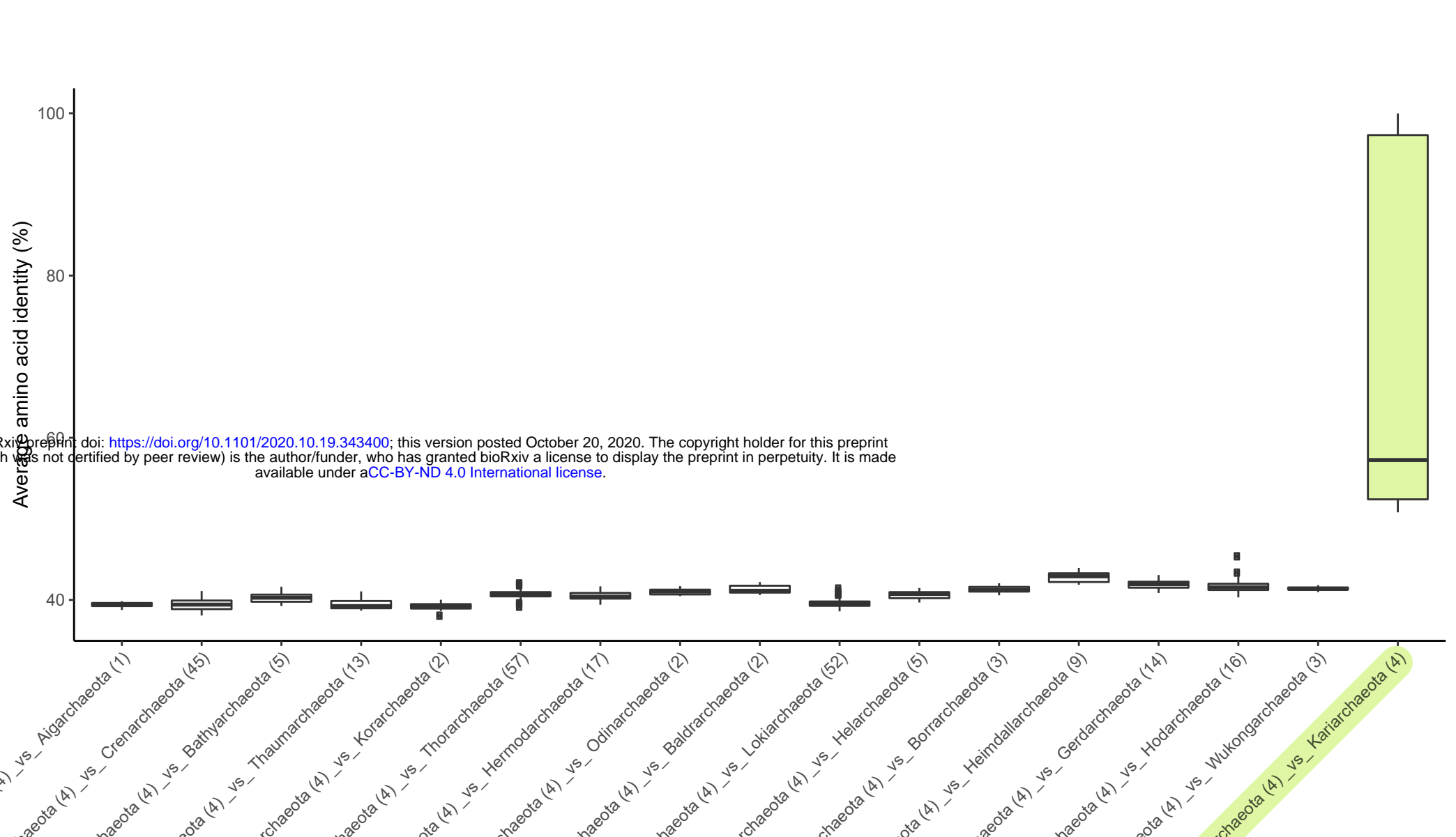
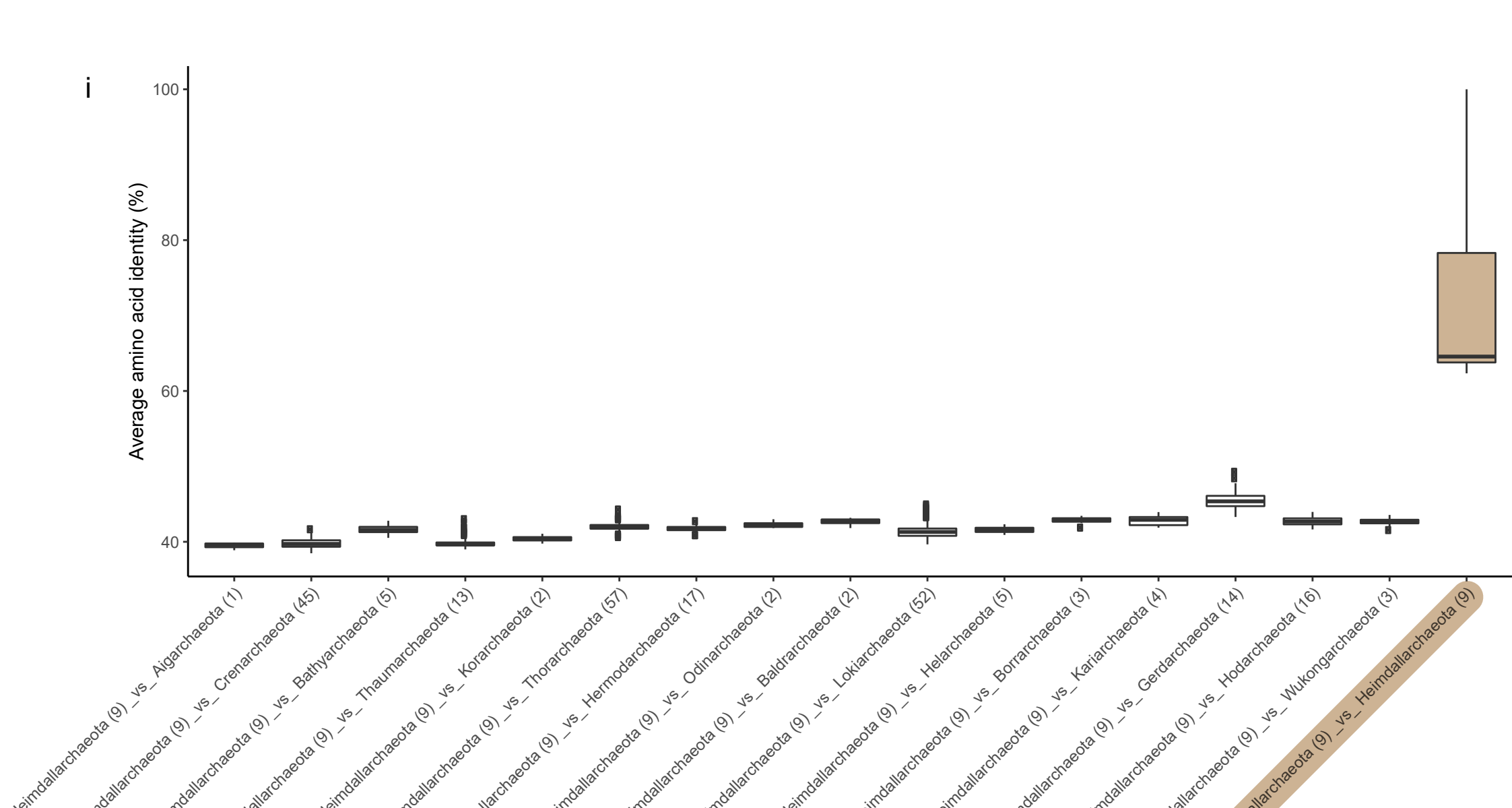
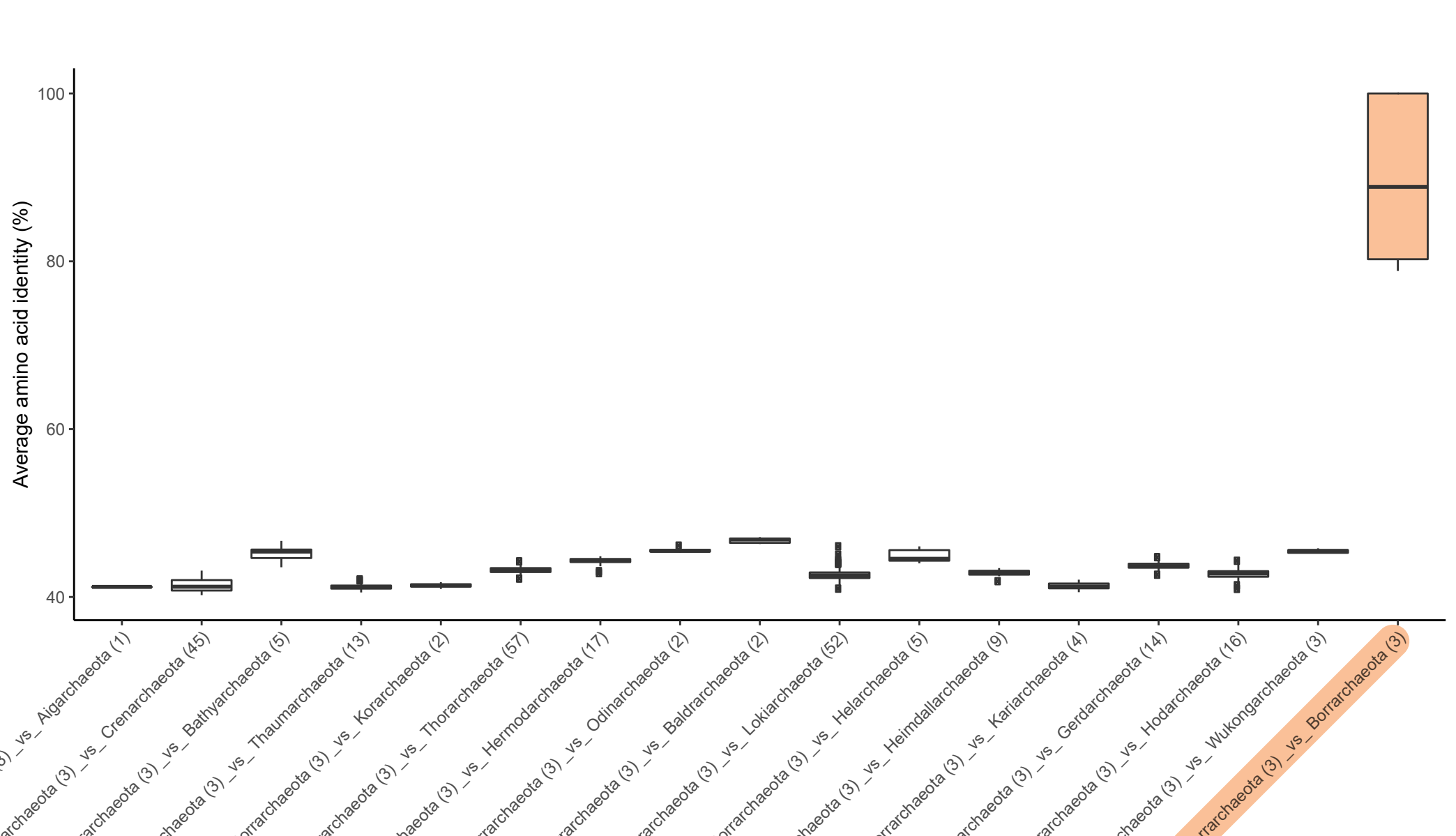
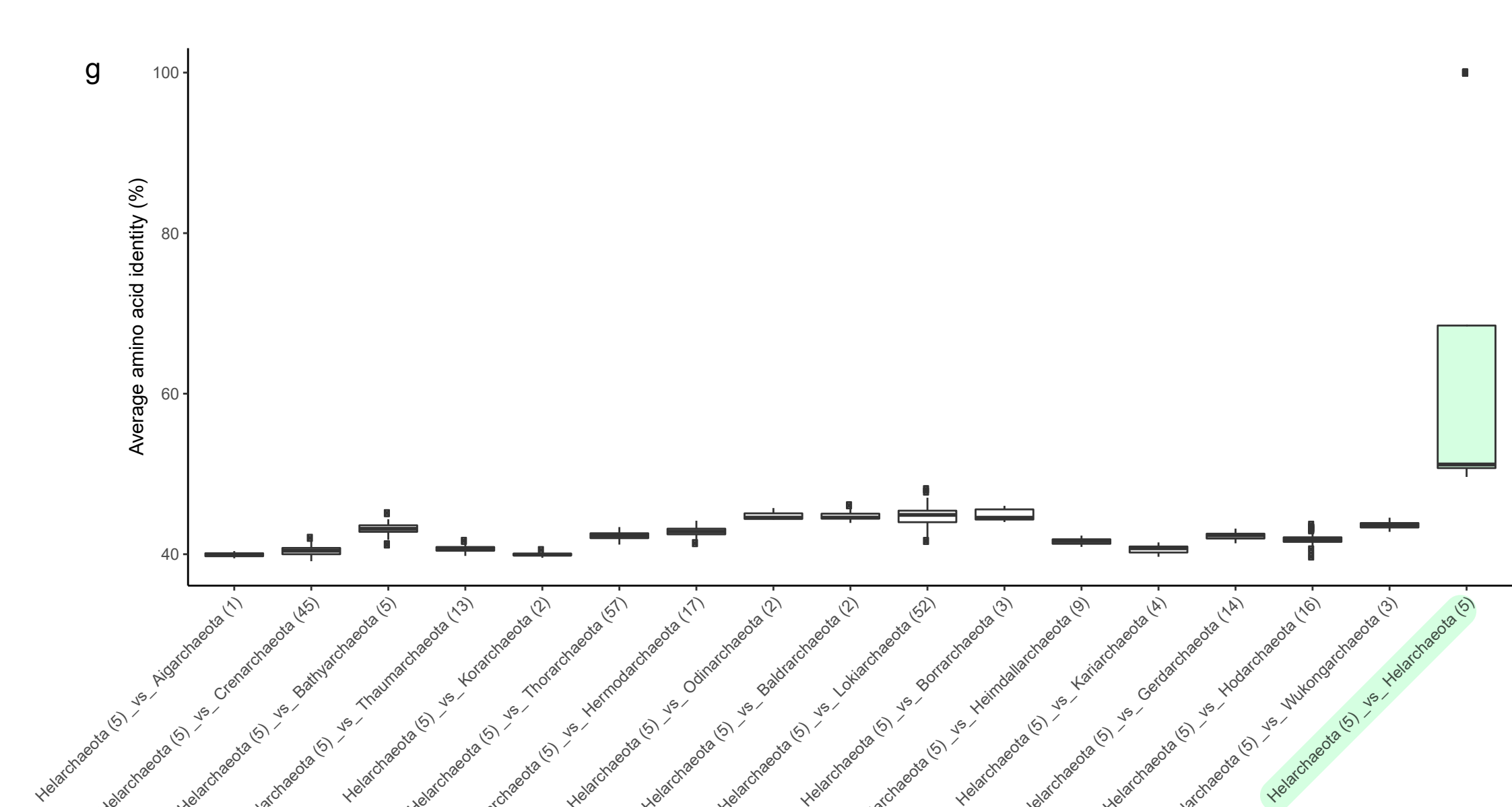
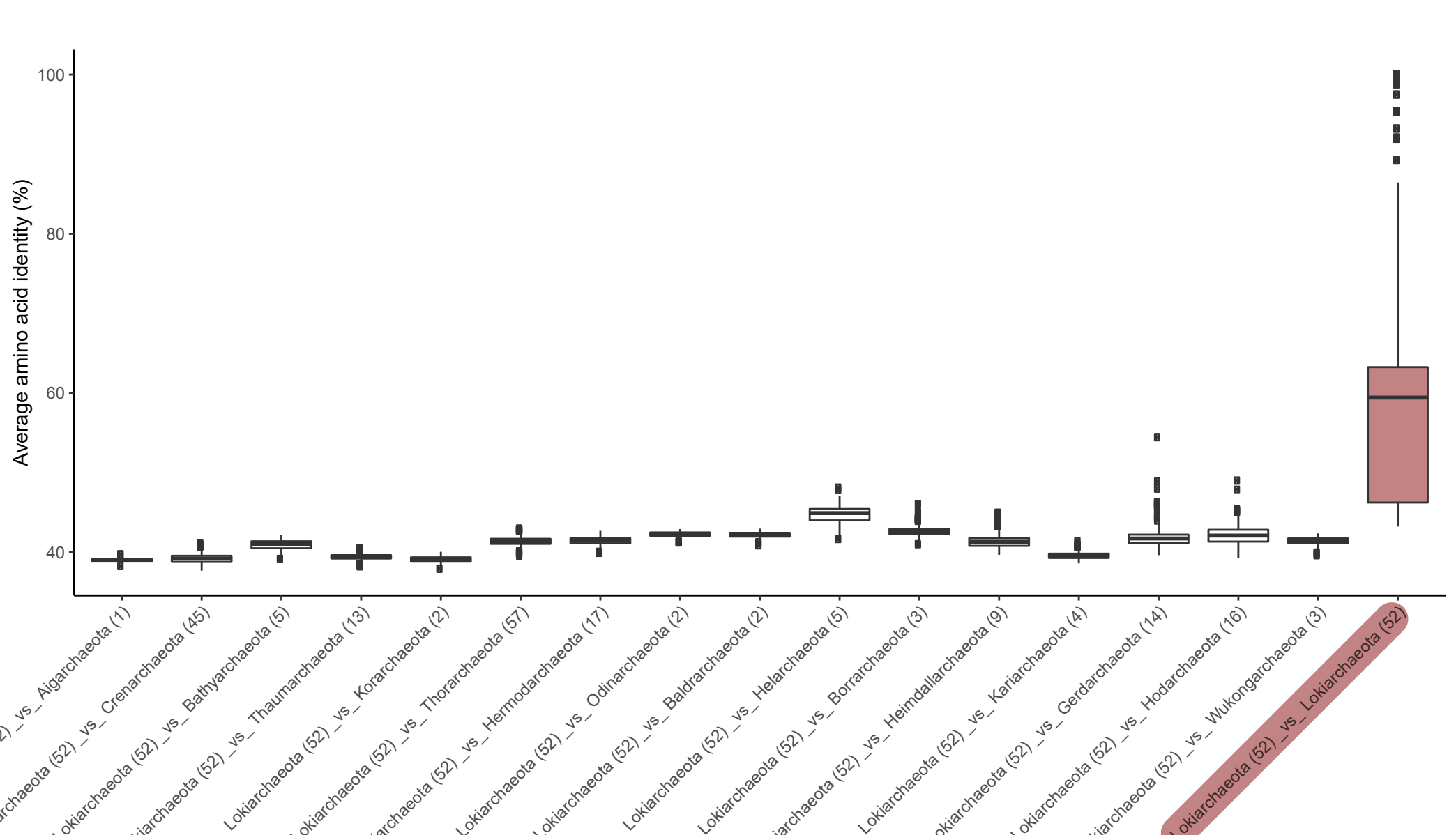
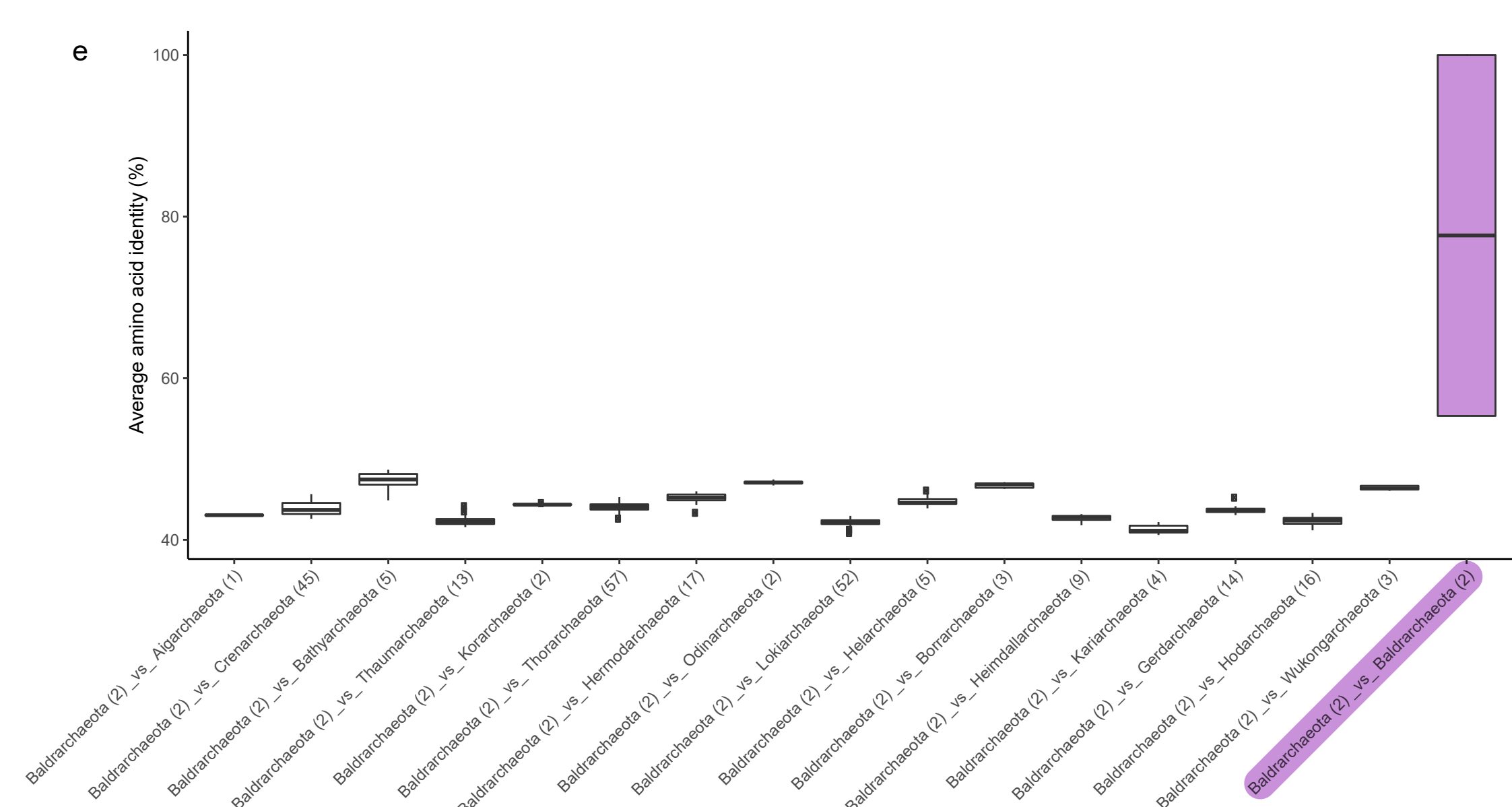
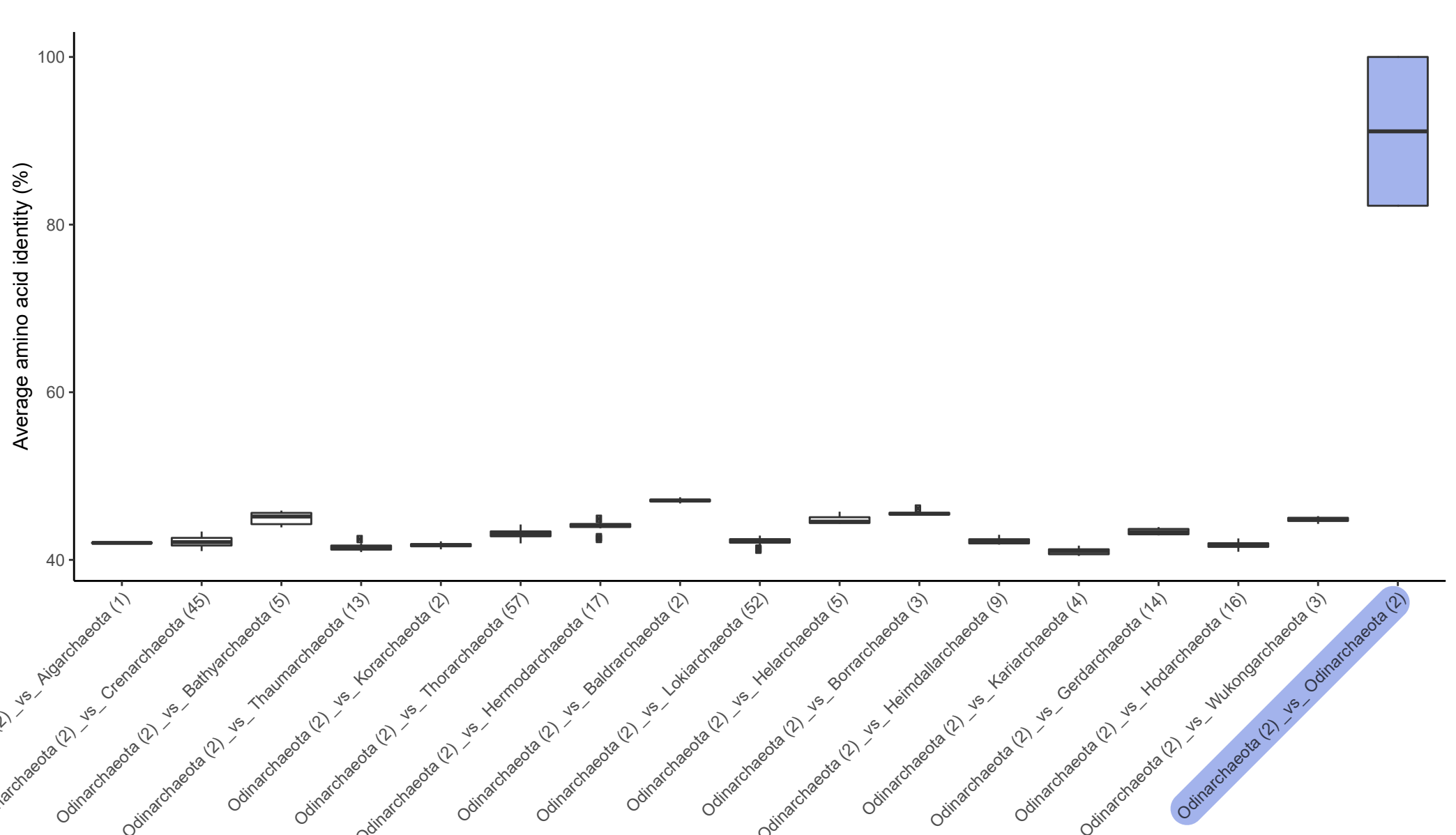
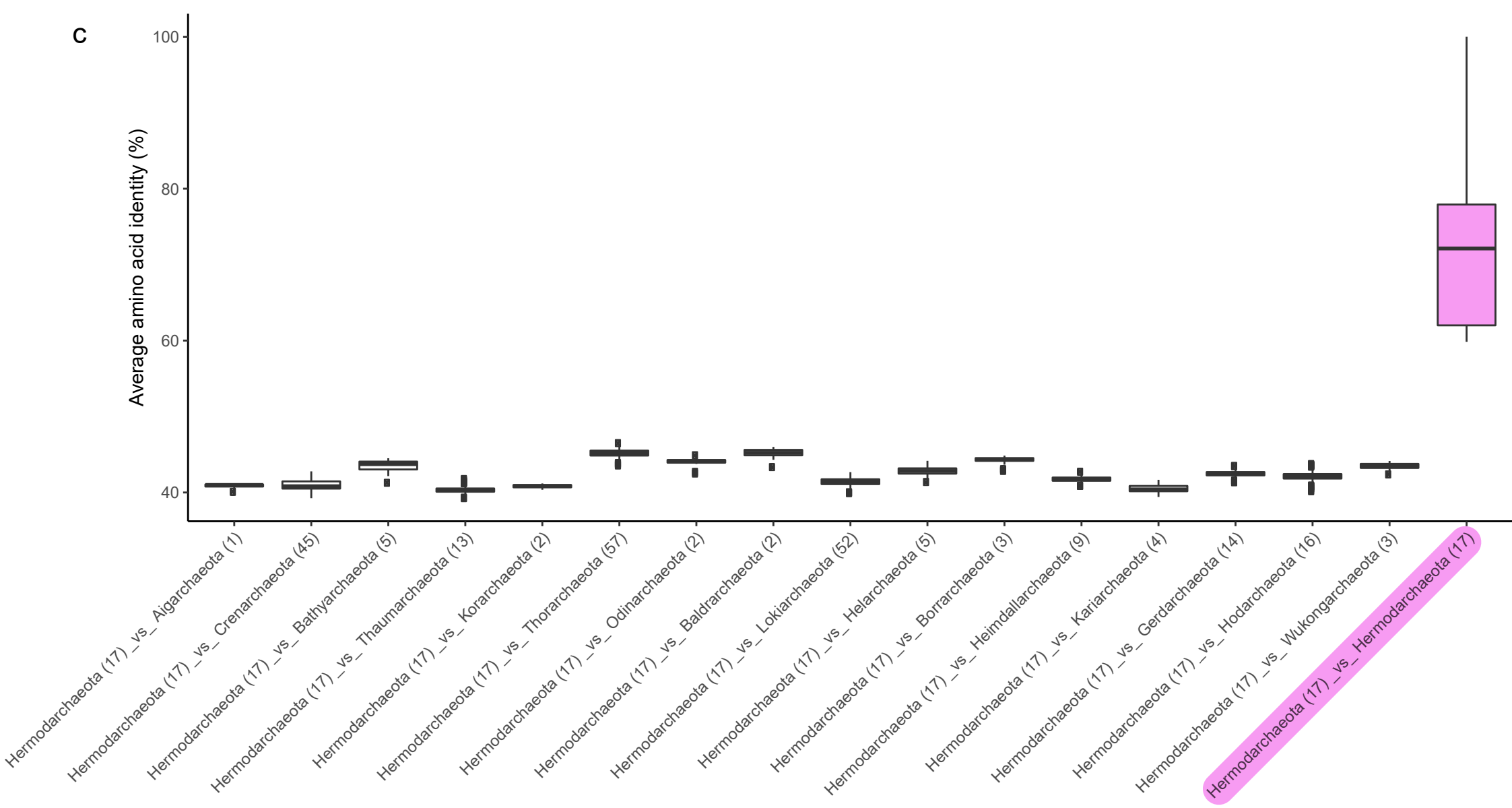
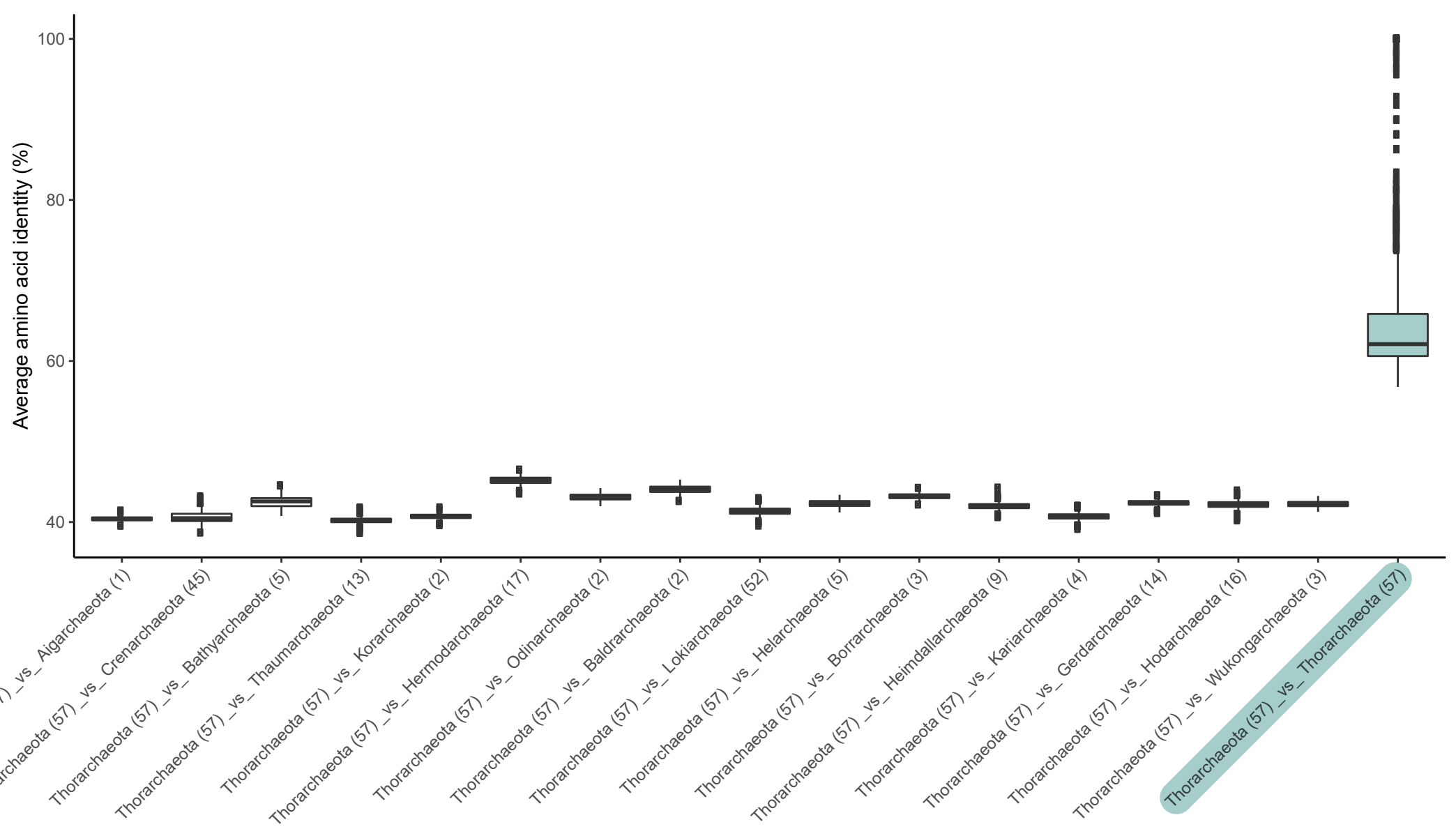
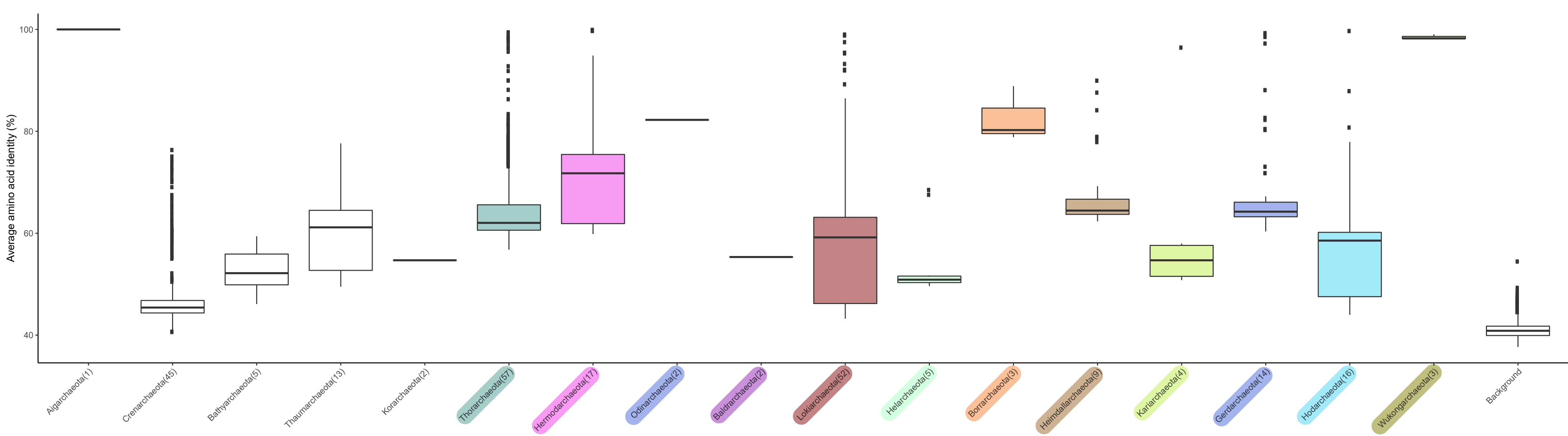
587 **Supplementary Table 8** The presence-absence of metabolic enzymes in Asgard archaea.

588

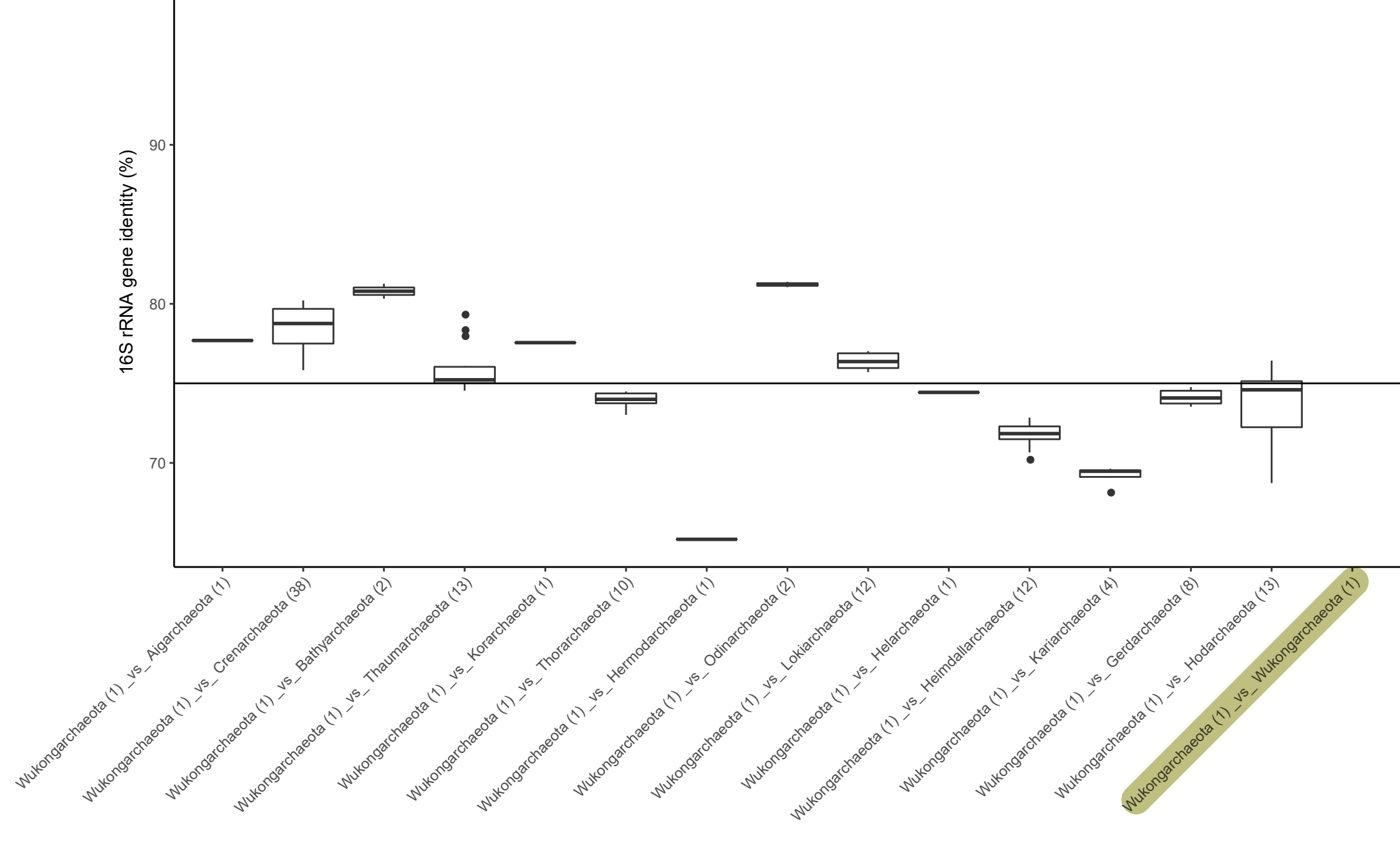
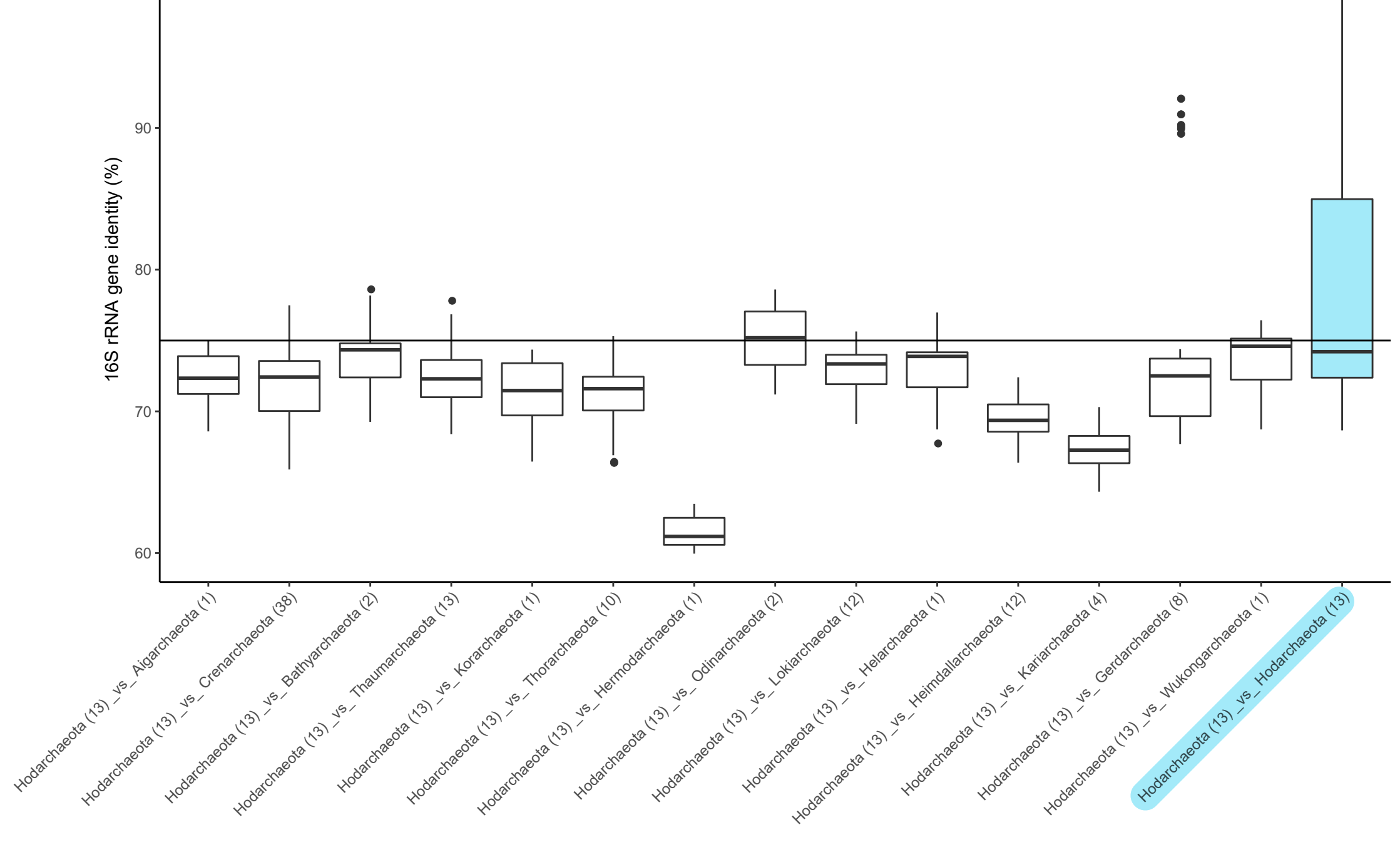
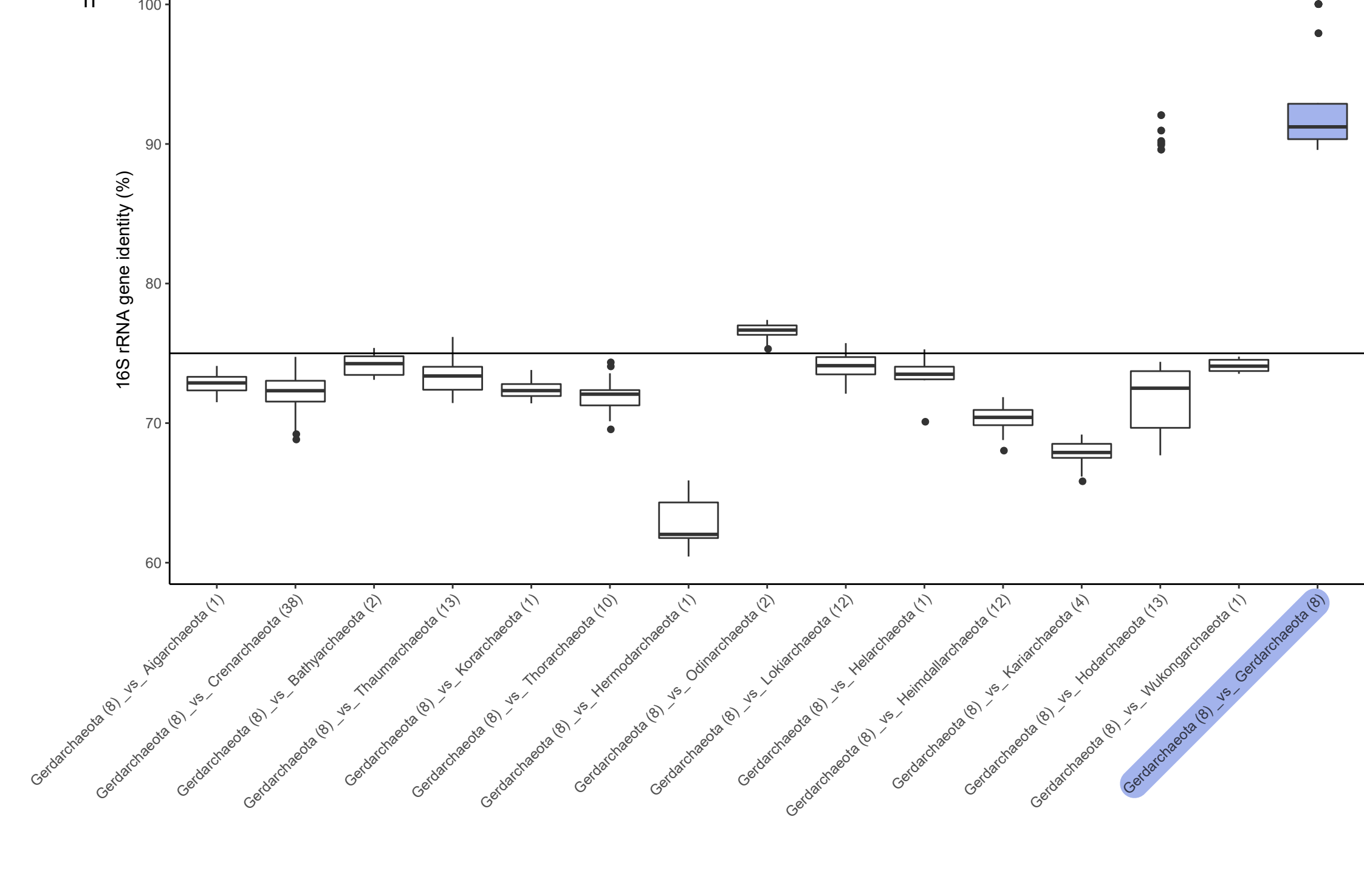
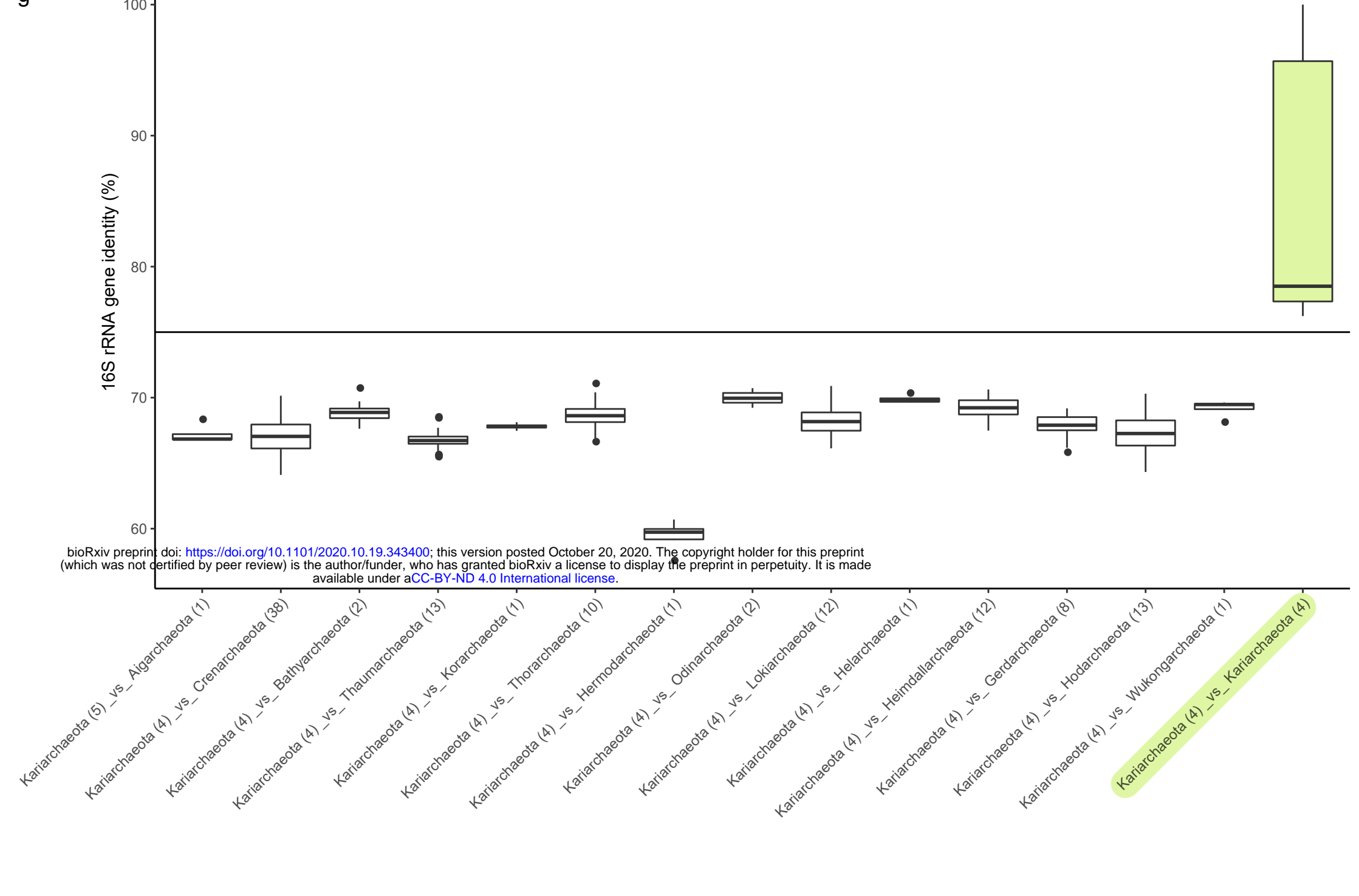
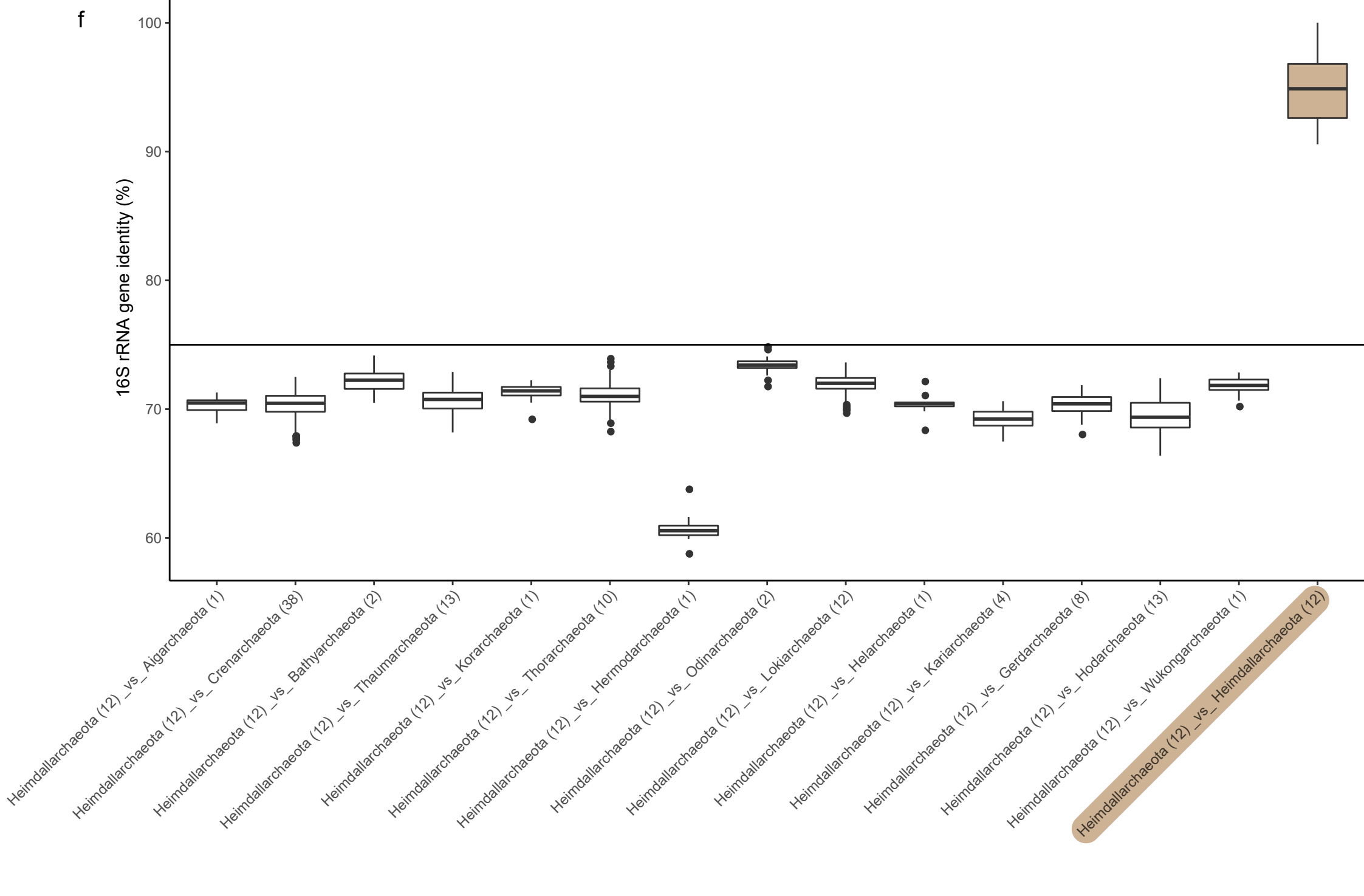
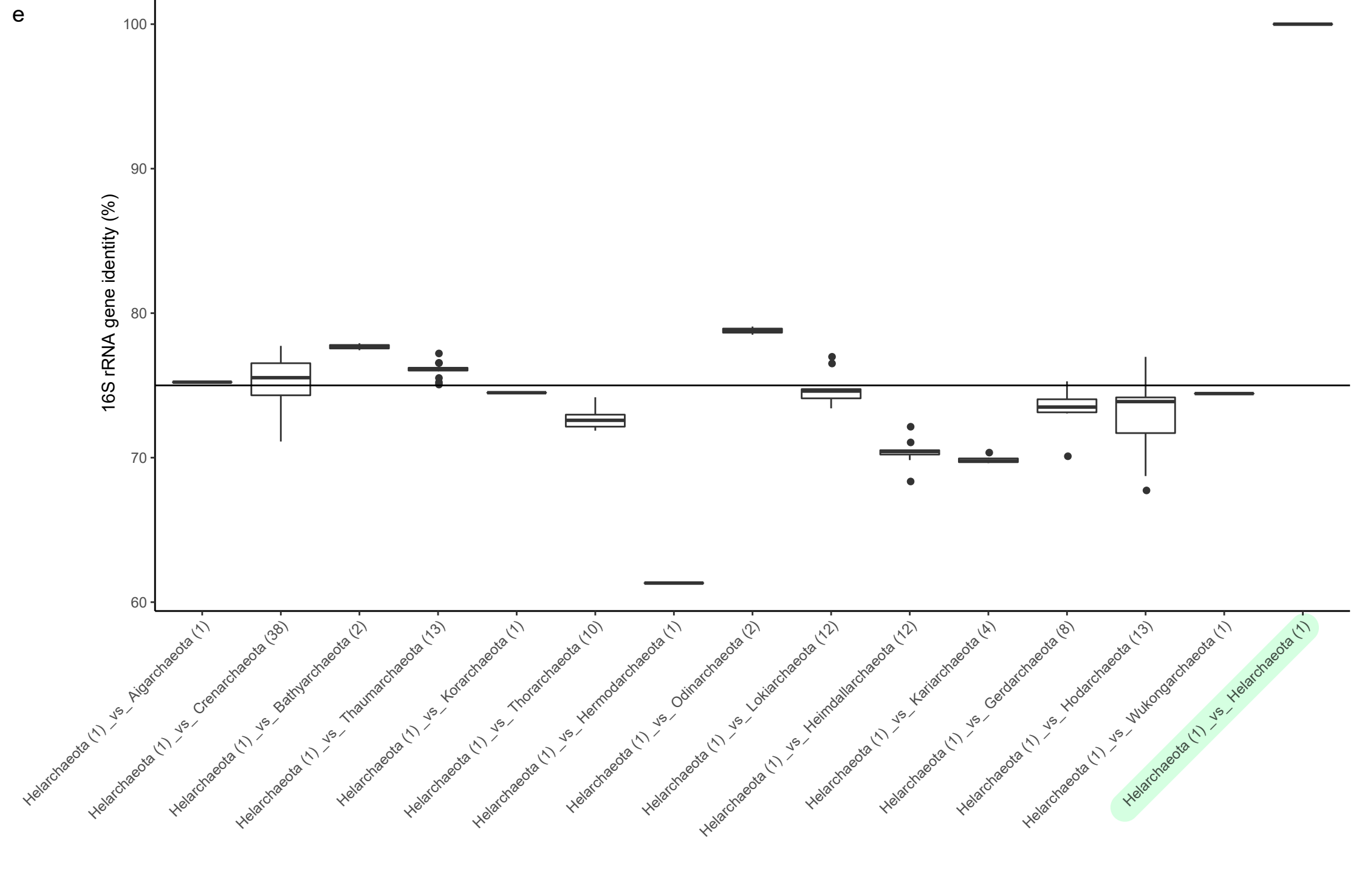
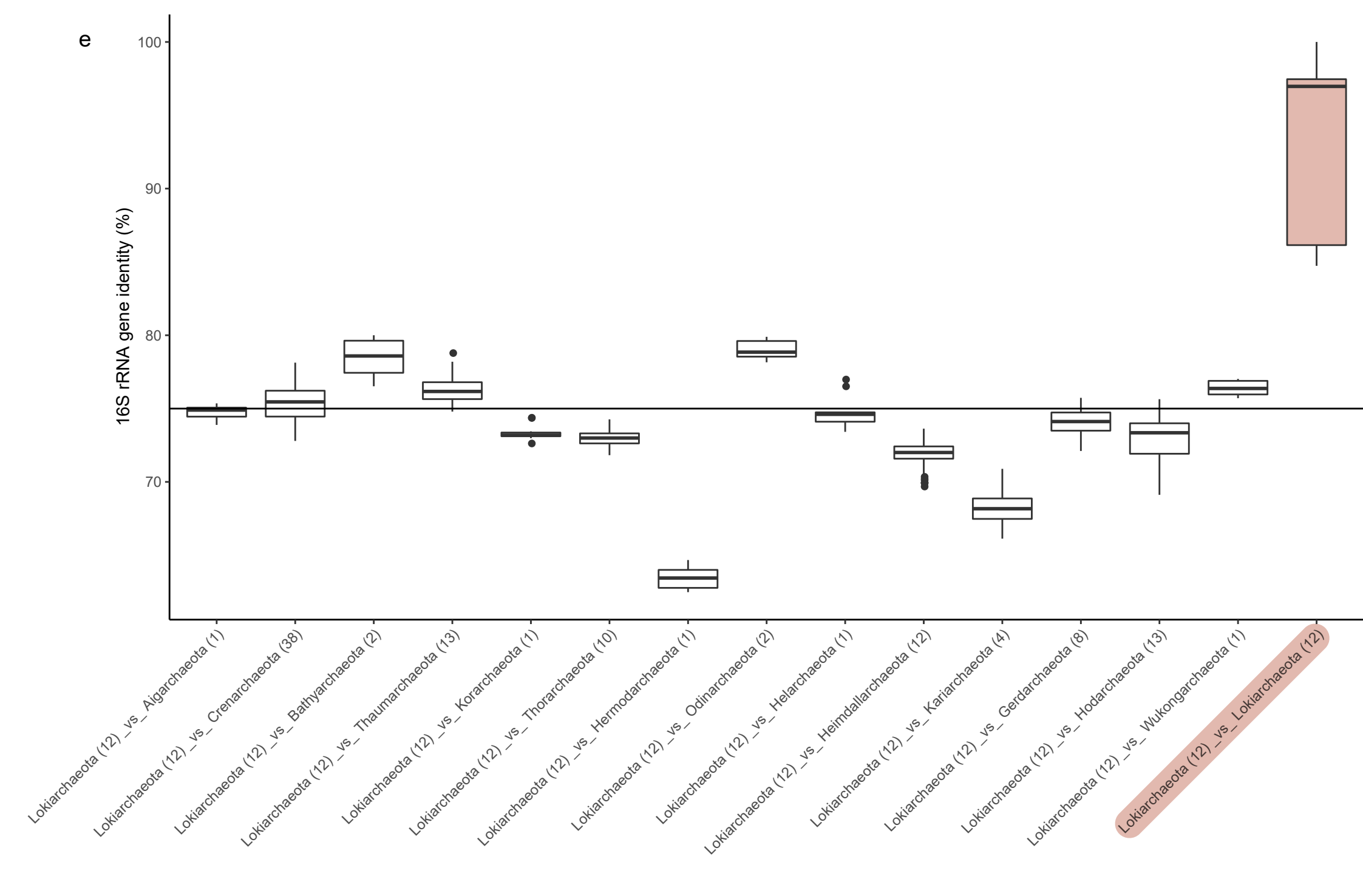
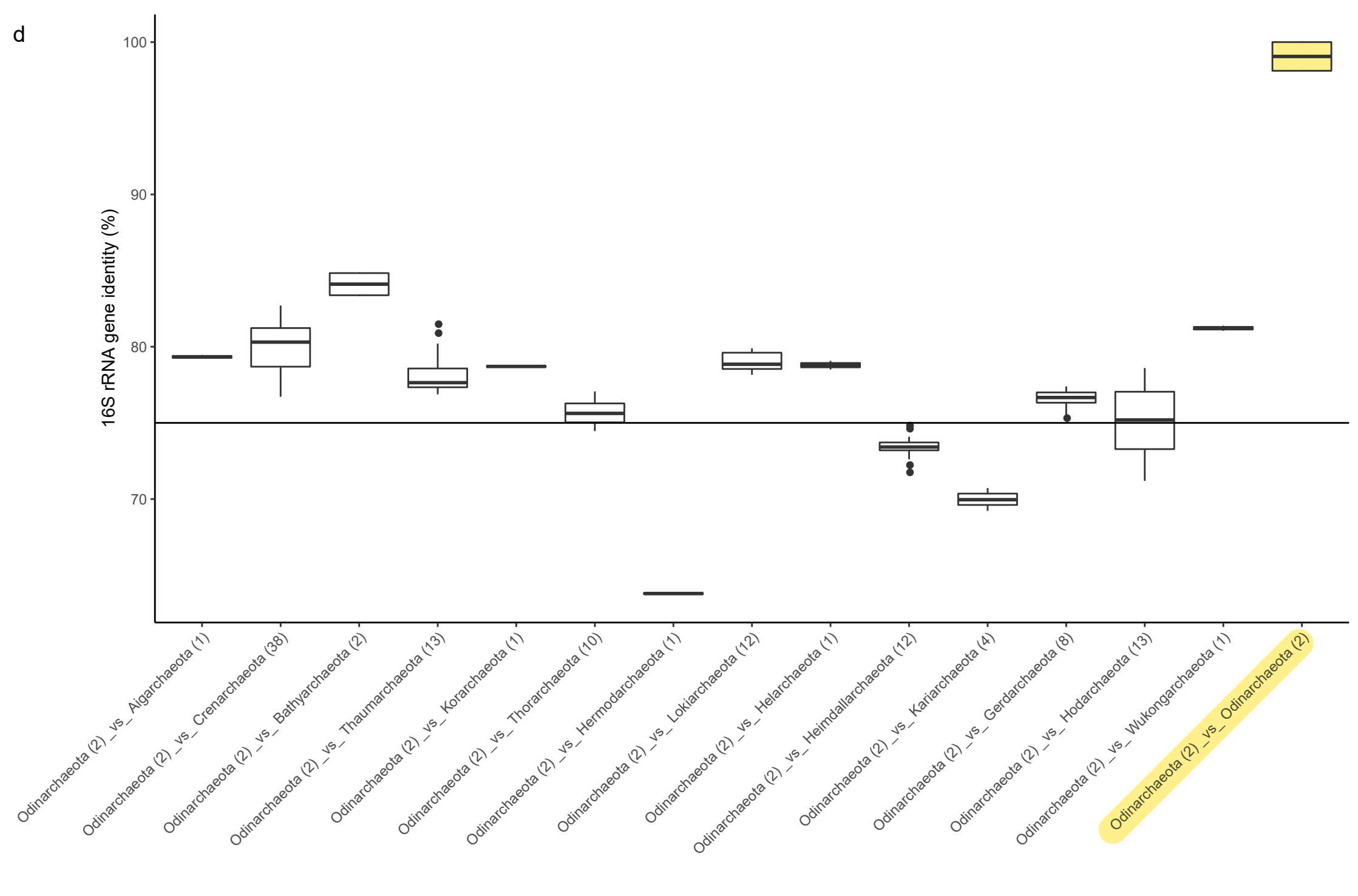
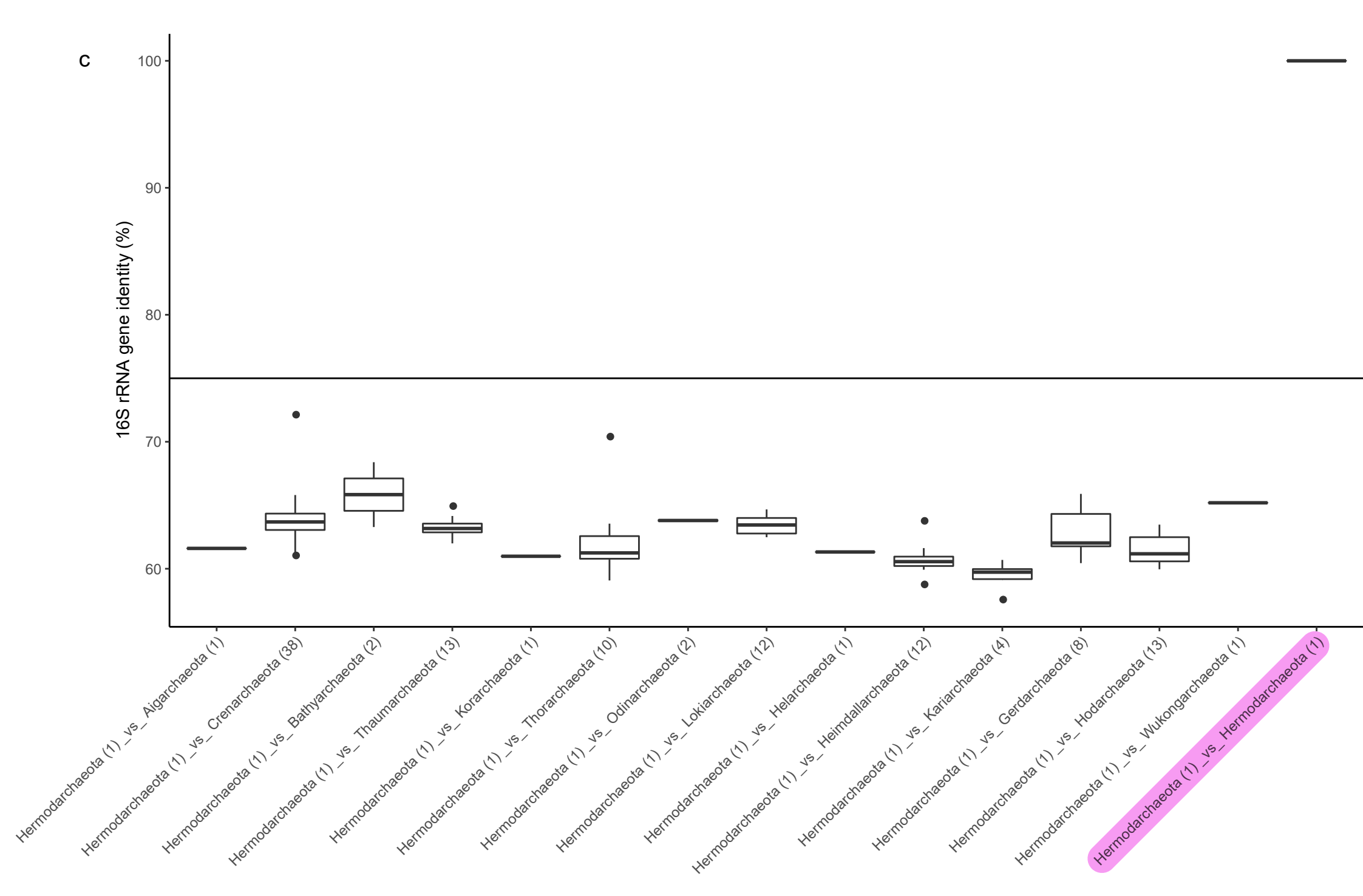
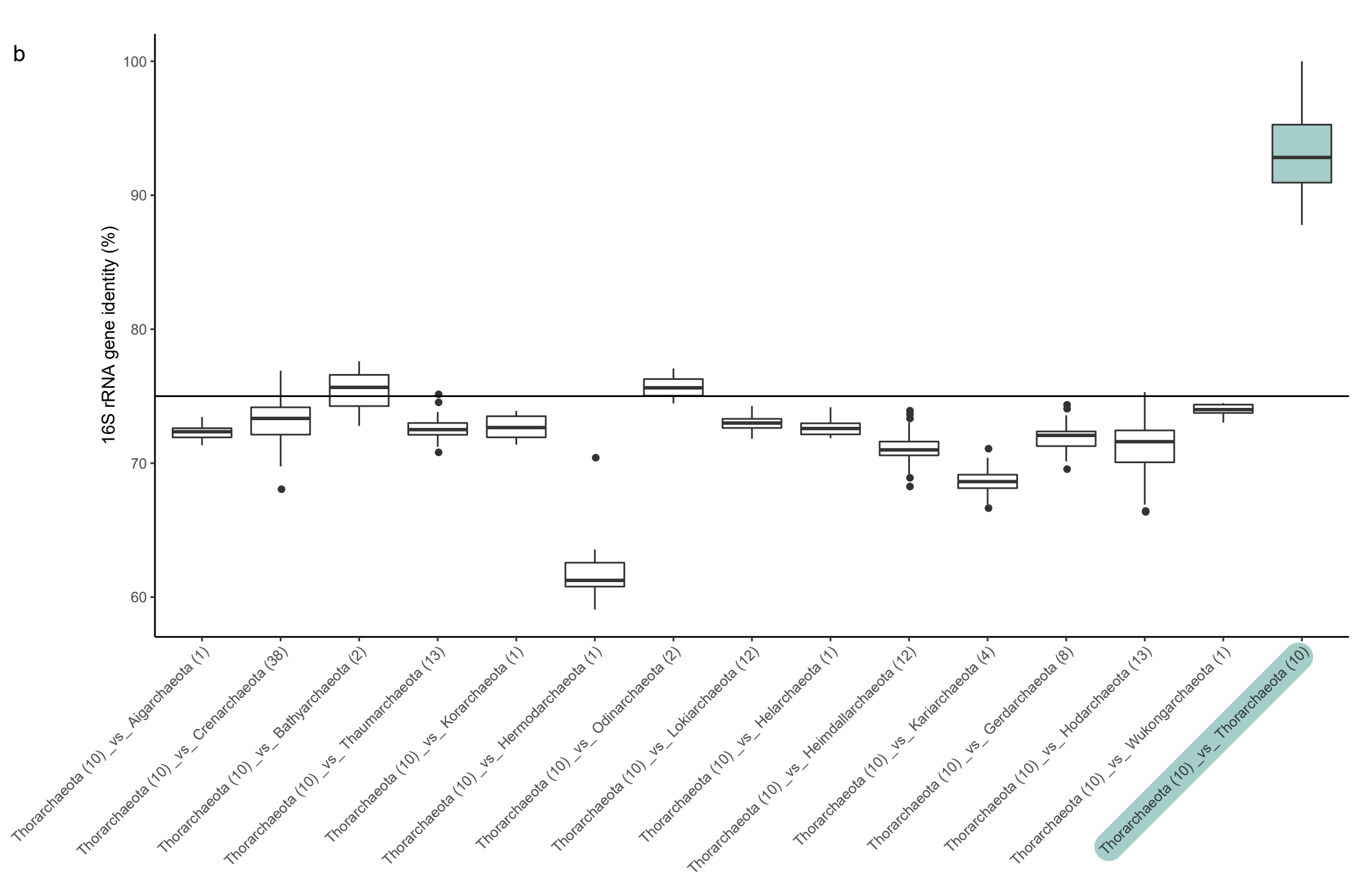
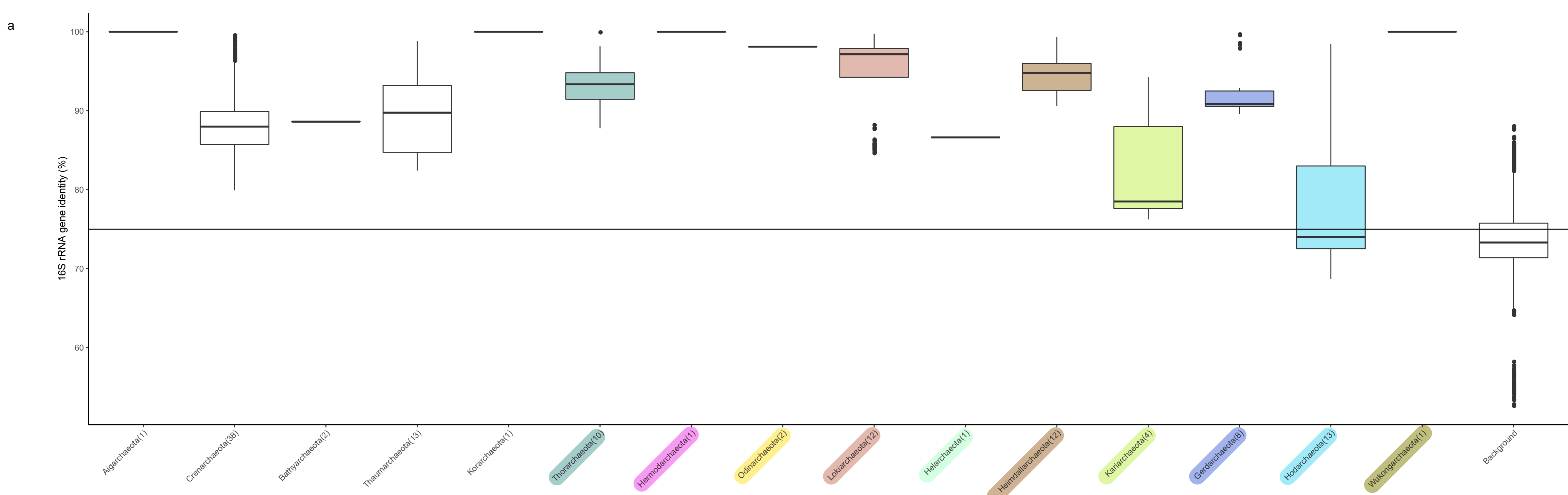


- Biotope**
- Coastal sediment
 - Hot spring
 - Hypersaline lake sediment
 - Marine water
 - Petroleum seep (Marine)
 - Freshwater sediment
 - Hydrothermal vent
 - Marine sediment
 - Petroleum field

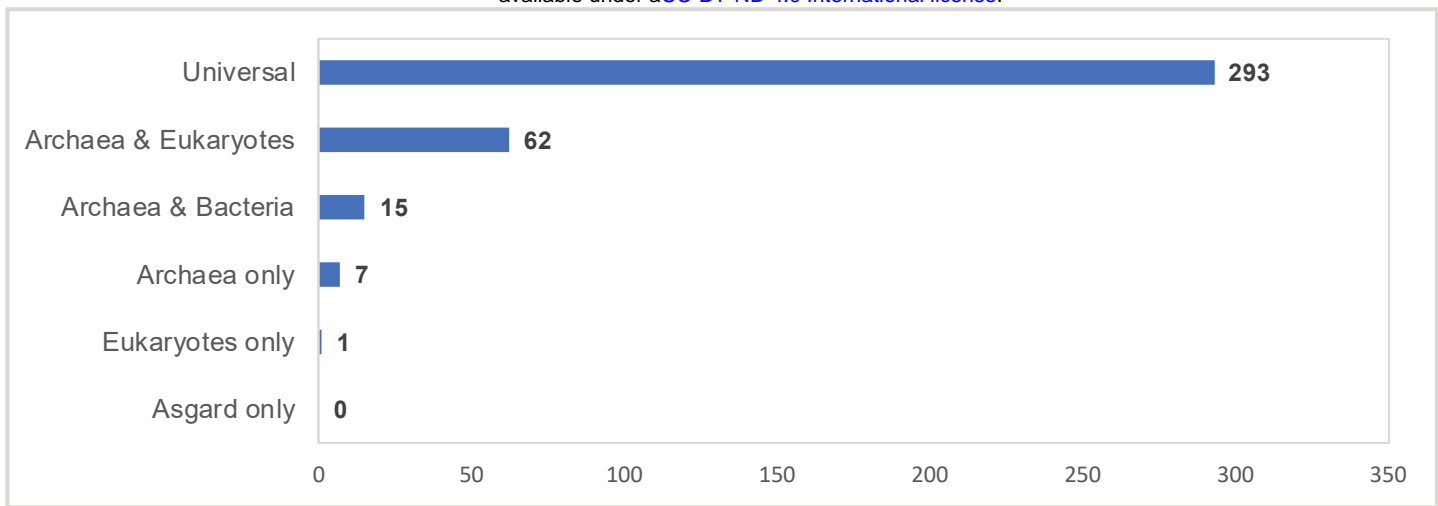




bioRxiv preprint doi: <https://doi.org/10.1101/2020.10.19.343400>; this version posted October 20, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.



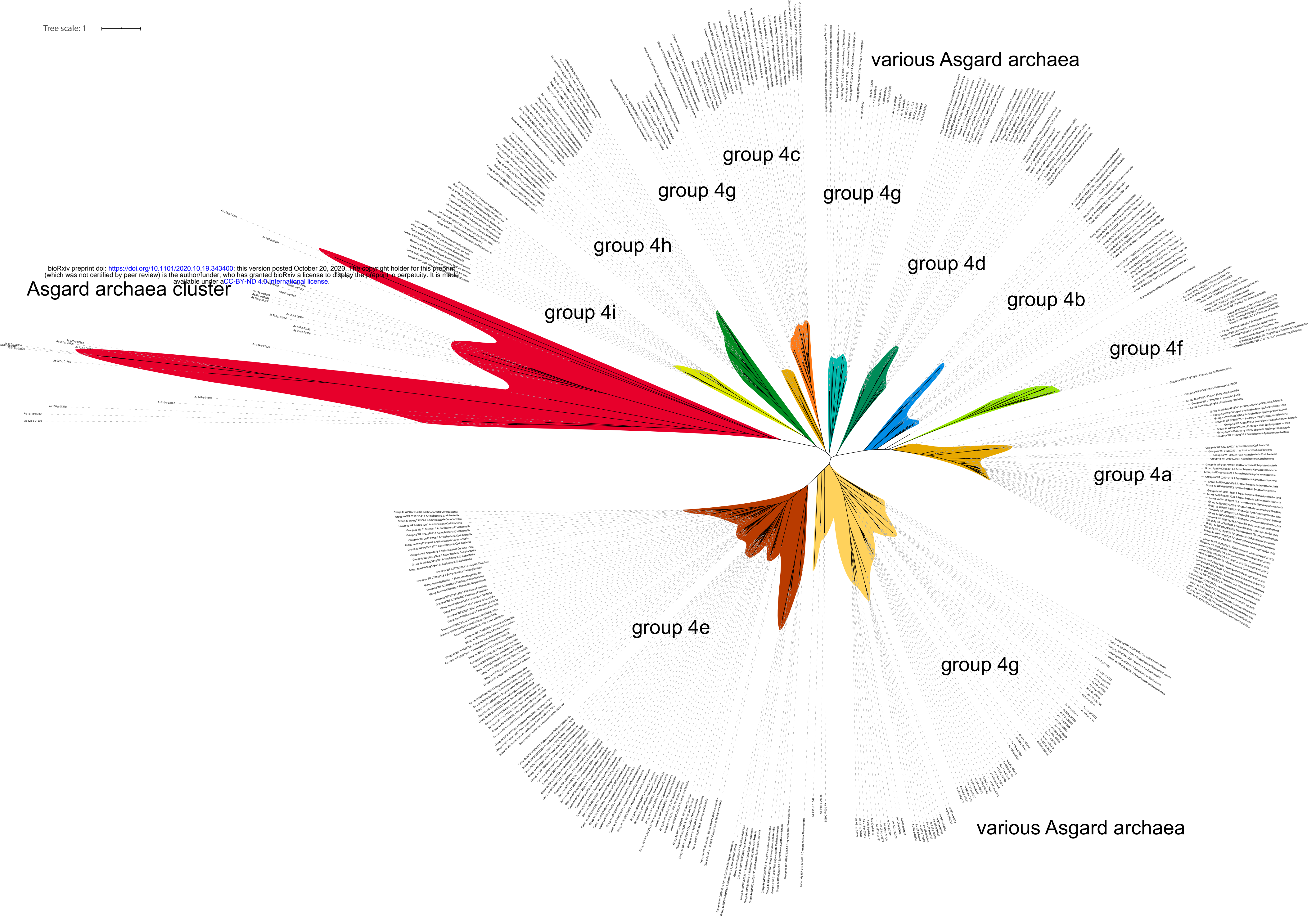
bioRxiv preprint doi: <https://doi.org/10.1101/2020.10.19.343400>; this version posted October 20, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.



Tree scale: 1

bioRxiv preprint doi: <https://doi.org/10.1101/2020.10.19.343400>; this version posted October 20, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

Asgard archaea cluster



various Asgard archaea

Lokiarchaeota

group 3c

group 3d

group 3c

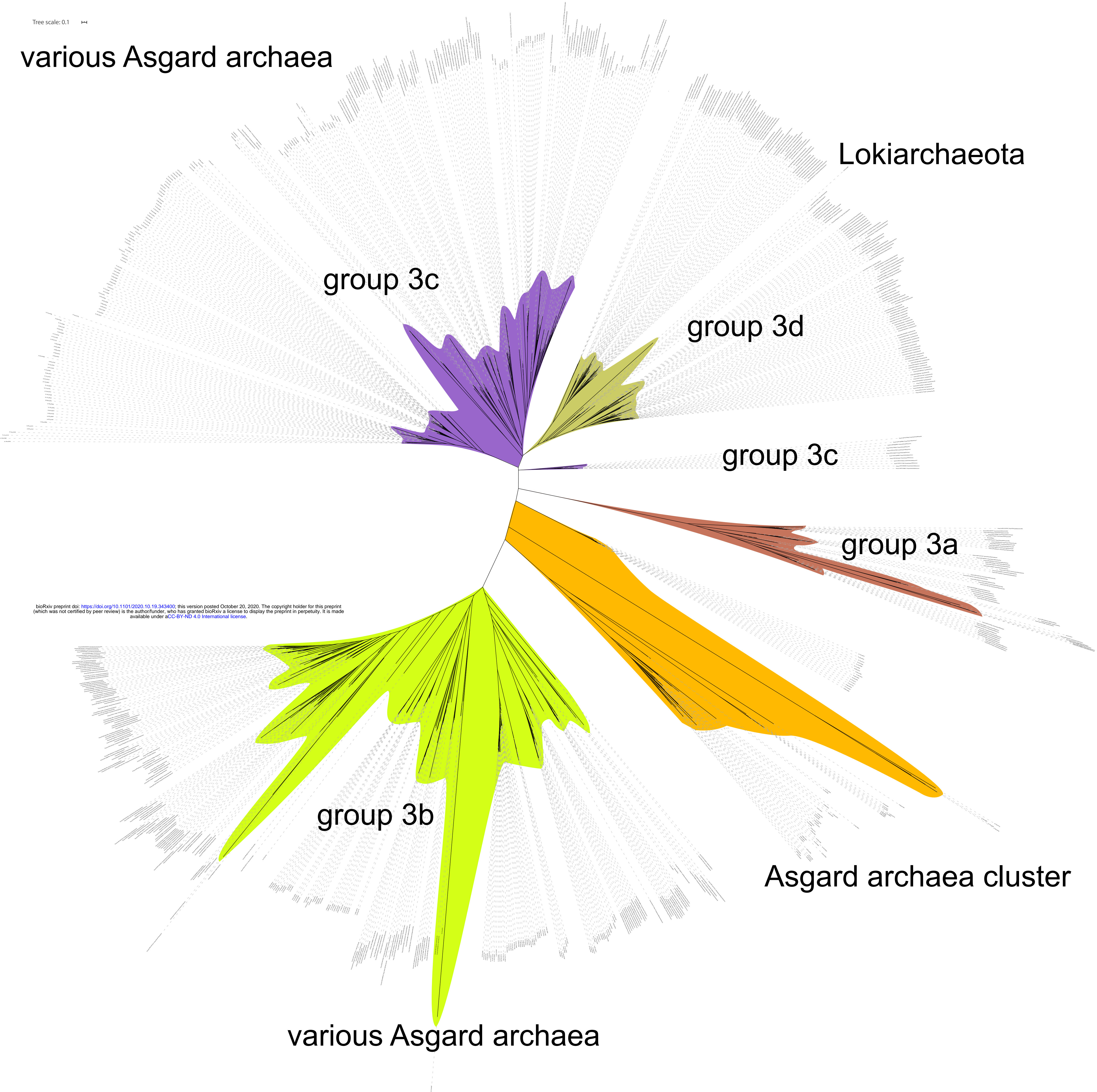
group 3a

group 3b

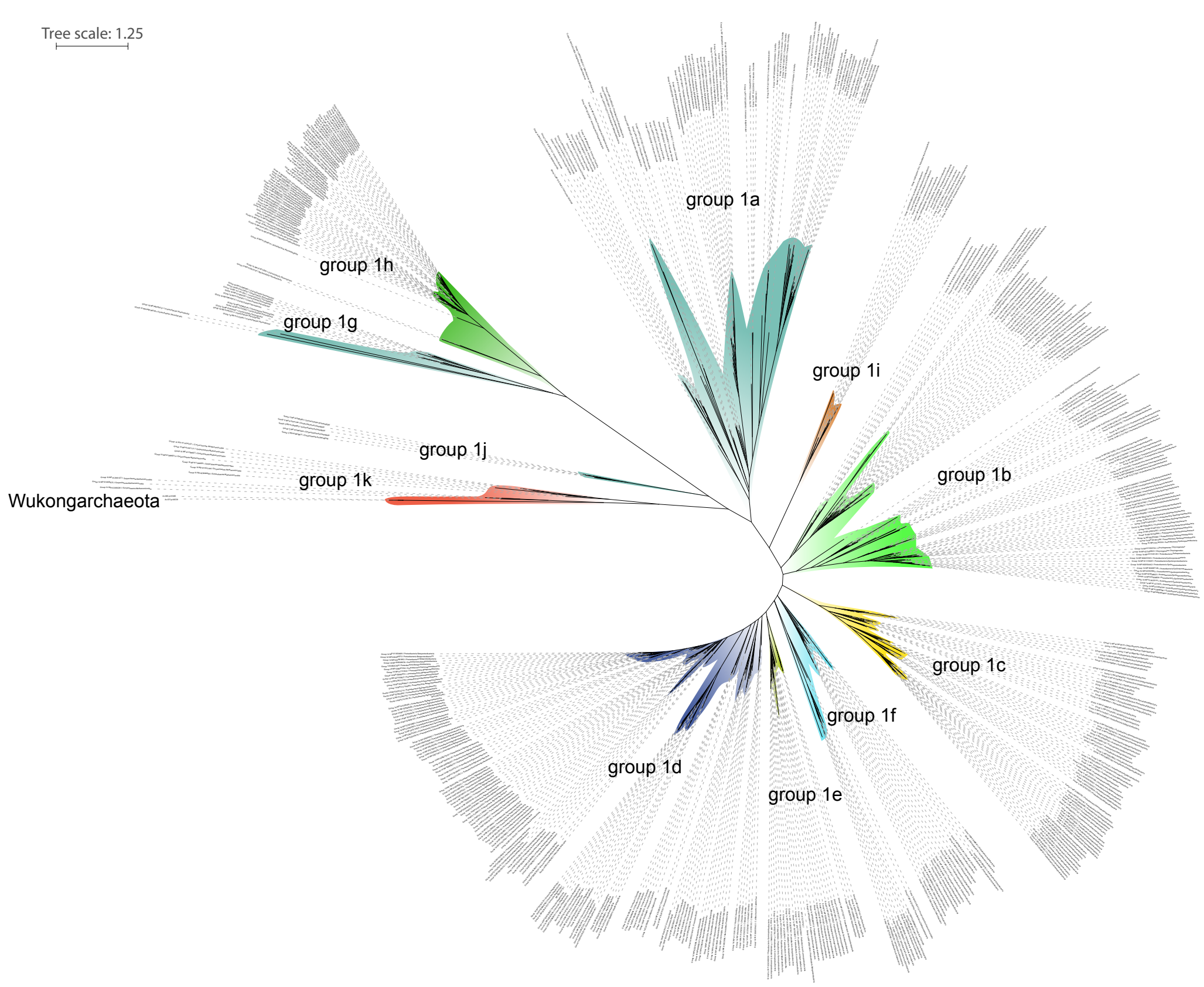
Asgard archaea cluster

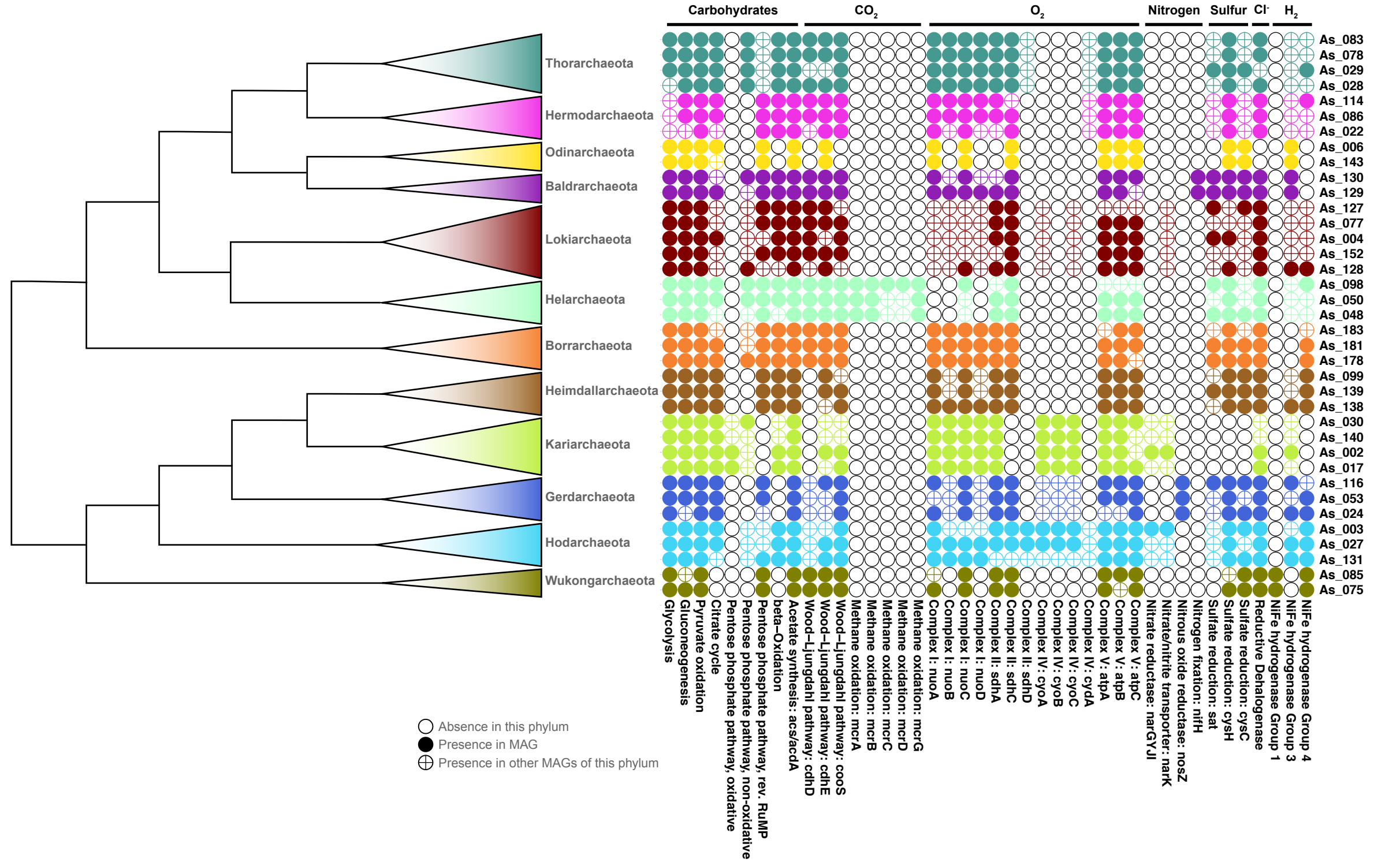
various Asgard archaea

bioRxiv preprint doi: <https://doi.org/10.1101/2020.10.19.343400>; this version posted October 20, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.



Tree scale: 1.25





Wukongarchaeota

As_075 contig_145_6336



As_085 contig_145_13955

