

**Supplementary information**

---

**Expanded diversity of Asgard archaea and their relationships with eukaryotes**

---

In the format provided by the authors and unedited

1 **Expanded diversity of Asgard archaea and their relationships with eukaryotes**

2

3 Yang Liu<sup>1,5</sup>, Kira S. Makarova<sup>2,5</sup>, Wen-Cong Huang<sup>1,5</sup>, Yuri I. Wolf<sup>2</sup>, Anastasia Nikolskaya<sup>2</sup>, Xinxu  
4 Zhang<sup>1</sup>, Mingwei Cai<sup>1</sup>, Cui-Jing Zhang<sup>1</sup>, Wei Xu<sup>3</sup>, Zhuhua Luo<sup>3</sup>, Lei Cheng<sup>4</sup>, Eugene V. Koonin<sup>2†</sup> &  
5 Meng Li<sup>1†</sup>

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

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

10 <sup>3</sup>Key Laboratory of Marine Biogenetic Resources, Third Institute of Oceanography, Ministry of Natural  
11 Resources, Xiamen, P. R. China

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

14 <sup>5</sup>These authors contributed equally: Yang Liu, Kira S. Makarova, Wen-Cong Huang.

15 †e-mail: [koonin@ncbi.nlm.nih.gov](mailto:koonin@ncbi.nlm.nih.gov); [limeng848@szu.edu.cn](mailto:limeng848@szu.edu.cn)

16

17           **1. Description of the new taxa**

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

21    ‘**Candidatus Wukongarchaeum yapensis**’ (yap’ensis N. L. masc. adj. pertaining to Yap trench, which is  
22 the geographical position where the first type material of this species was obtained). Type material is the  
23 genome designated as As\_085 (Yap4.bin4.70) representing ‘*Candidatus Wukongarchaeum yapensis*’. The  
24 genome “As\_085” represents a MAG consisting of 2.16 Mbps in 277 contigs with an estimated  
25 completeness of 92.52%, an estimated contamination of 4.05%, a 16S and 23S rRNA gene and 14 tRNAs.  
26 The MAG recovered from a marine water metagenome (Yap trench, Western Pacific), with an estimated  
27 depth of coverage of 31.4, has a GC content of 38%.

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

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

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

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

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

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

44    **Candidatus Wukongarchaeota** (Wu.kong.ar.chae.o’ta. N.L. neut. n. Wukongarchaeum a (*Candidatus*)  
45 type genus of the class of the phylum; -ota ending to denote the phylum; N.L. neut. pl. n.  
46 Wukongarchaeota the Wukongarchaeum phylum)

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

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

51    ‘**Candidatus Hodarchaeum mangrovi**’ (man.gro’vi N.L. fem. n. of a mangrove, referring to the isolation  
52 of the type material from mangrove soil). Type material is the genome designated as As\_027

53 (FT2\_5\_011) representing ‘*Candidatus* Hodarchaeum mangrovi’. The genome “As\_027” represents a  
54 MAG consisting of 4.01 Mbps in 348 contigs with an estimated completeness of 93.61%, an estimated  
55 contamination of 0.93%, a 23S rRNA gene and 14 tRNAs. The MAG recovered from mangrove sediment  
56 metagenomes (Futian Nature Reserve, China), with an estimated depth of coverage of 17.9, has a GC  
57 content of 32.9%.

58 ***Candidatus* Hodarchaeaceae** (Hod.ar.chae.a.ce’ae. N.L. neut. n. Hodarchaeum a (Candidatus) type  
59 genus of the family; -aceae ending to denote the family; N.L. fem. pl. n. Hodarchaeaceae the  
60 Hodarchaeum family).

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

63 ***Candidatus* Hodarchaeales** (Hod.ar.chae.a’les. N.L. neut. n. Hodarchaeum a (Candidatus) type genus of  
64 the order; -ales ending to denote the order; N.L. fem. pl. n. Hodarchaeales the Hodarchaeum order).

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

67 ***Candidatus* Hodarchaeia** (Hod.ar.chae’i.a. N.L. neut. n. Hodarchaeum a (Candidatus) type genus of the  
68 order of the class; -ia ending to denote the class; N.L. fem. pl. n. Hodarchaeia the Hodarchaeum class).

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

71 ***Candidatus* Hodarchaeota** (Hod.ar.chae.o’ta. N.L. neut. n. Hodarchaeum a (Candidatus) type genus of  
72 the class of the phylum; -ota ending to denote the phylum; N.L. neut. pl. n. Hodarchaeota the  
73 Hodarchaeum phylum)

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

76 **‘*Candidatus* Kariarchaeum’** (Ka.ri.ar.chae’um. N. L. n. Kari the god of wind in Norse mythology); N.L.  
77 neut. N. *archaeum* (from Gr. adj. archaios ancient) archaeon; N. L. neut. N. Kariarchaeum.

78 **‘*Candidatus* Kariarchaeum pelagius’** (pe.la’gi.us. L. masc. adj. of or belonging to the sea, referring to  
79 the isolation of the type material from the Ocean). Type material is the genome designated as As\_030  
80 (RS678) representing ‘*Candidatus* Kariarchaeum pelagius’. The genome “As\_030” represents a MAG  
81 consisting of 1.41 Mbps in 76 contigs, an estimated completeness of 83.18%, with an estimated  
82 contamination of 1.87%, a 23S, 16S and 5S rRNA genes and 18 tRNAs. The MAG recovered from a  
83 marine metagenome (Saudi Arabia: Red Sea) has a GC content of 30.11%.

84 ***Candidatus* Kariarchaeaceae** (Ka.ri.ar.chae.a.ce’ae. N.L. neut. n. Kariarchaeum a (Candidatus) type  
85 genus of the family; -aceae ending to denote the family; N.L. fem. pl. n. Kariarchaeaceae the  
86 Kariarchaeum family).

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

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

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

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

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

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

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

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

104 '**Candidatus Borrarchaeum yapensis**' (yap'ensis N. L. masc. adj. pertaining to Yap trench, which is the  
105 geographical position where the first type material of this species was obtained). Type material is the  
106 genome designated as As\_181 (Yap2000.bin9.141) representing '*Candidatus* Borrarchaeum yapensis'.  
107 The genome "As\_181" represents a MAG consisting of 3.63 Mbps in 125 contigs, with an estimated  
108 completeness of 95.02%, an estimated contamination of 5.61% and 11 tRNAs. The MAG, recovered from  
109 a marine water metagenome (Yap trench, Western Pacific) with an estimated depth coverage of 15.04, has  
110 a GC content of 37.1%.

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

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

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

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

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

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

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

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

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

132 '***Candidatus Baldrarchaeum yapensis***' (yap'ensis N. L. masc. adj. pertaining to Yap trench, which is the  
133 geographical position where the first type material of this species was obtained). Type material is the  
134 genome designated as As\_130 (Yap30.bin9.72) representing '*Candidatus* Baldrarchaeum yapensis'. The  
135 genome "As\_130" represents a MAG consisting of 2.27 Mbps in 100 contigs, with an estimated  
136 completeness of 93.93%, an estimated contamination of 3.74%, a 23S and 16S rRNA gene and 15 tRNAs.  
137 The MAG, recovered from a marine water metagenome (Yap trench, Western Pacific) with an estimated  
138 depth coverage of 39.99, has a GC content of 45.95%.

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

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

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

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

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

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

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

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

162 **‘*Candidatus Hermodarchaeum yapensis*’** (yap’ensis N. L. masc. adj. pertaining to Yap trench, which is  
163 the geographical position where the first type material of this species was obtained). Type material is the  
164 genome designated as As\_086 (Yap4.bin9.105) representing ‘*Candidatus Hermodarchaeum yapensis*’.  
165 The genome ‘As\_086’ represent a MAG consisting of 2.71 Mbps in 77 contigs, with an estimated  
166 completeness of 92.99%, an estimated contamination of 1.87%, a 23S and 16S rRNA gene and 16 tRNAs.  
167 The MAG, recovered from a marine water metagenome (Yap trench, Western Pacific) with an estimated  
168 depth coverage of 19.24, has a GC content of 44.69%.

169 ***Candidatus Hermodarchaeaceae*** (Her.mod.ar.chae.a.ce’ae. N.L. neut. n. Hermodarchaeum a  
170 (*Candidatus*) type genus of the family; -aceae ending to denote the family; N.L. fem. pl. n.  
171 Hermodarchaeaceae the Hermodarchaeum family).

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

174 ***Candidatus Hermodarchaeales*** (Her.mod.ar.chae.a’les. N.L. neut. n. Hermodarchaeum a (*Candidatus*)  
175 type genus of the order; -ales ending to denote the order; N.L. fem. pl. n. Hermodarchaeales the  
176 Hermodarchaeum order).

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

179 ***Candidatus Hermodarchaeia*** (Her.mod.ar.chae’i.a. N.L. neut. n. Hermodarchaeum a (*Candidatus*) type  
180 genus of the order of the class; -ia ending to denote the class; N.L. fem. pl. n. Hermodarchaeia the  
181 Hermodarchaeum class).

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

184 ***Candidatus Hermodarchaeota*** (Her.mod.ar.chae.o’ta. N.L. neut. n. Hermodarchaeum a (*Candidatus*)  
185 type genus of the class of the phylum; -ota ending to denote the phylum; N.L. neut. pl. n.  
186 Hermodarchaeota the Hermodarchaeum phylum)

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

189

## 190 **2. Clusters of orthologous genes of Asgard archaea**

191 The previous analyses of Asgard genomes detected a large fraction of “dark matter” genes<sup>69</sup>. For  
192 example, in the recently published complete genome of *Candidatus P. syntrophicum*, 45% of the proteins  
193 are annotated as “hypothetical”. Here, an effort was made to improve the annotation of Asgard genomes

194 by investigating this dark matter in greater depth and developing a dedicated platform for Asgard  
195 comparative genomics. To this end, Asgard Clusters of Orthologous Genes (asCOGs) were constructed,  
196 and the most sensitive available methods of sequence analysis were employed to annotate additional  
197 Asgard proteins, attempting, in particular, to expand the catalogue of Asgard homologs of ESPs (see  
198 Methods for details).

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

205 The gene commonality plot for the asCOGs shows an abrupt drop at the right end, which reflects a  
206 surprising deficit of nearly universal genes (Extended Data Fig. 3). Such shape of the gene commonality  
207 curve appears anomalous compared to other major groups of archaea or bacteria with many sequenced  
208 genomes<sup>70</sup>. For example, in the case of the TACK superphylum of archaea, for which the number of  
209 genomes available is similar to that for Asgard, with a comparable level of diversity, the commonality  
210 plot shows no drop at the right end, but instead, presents a clear uptick, which corresponds to the core of  
211 genes represented in (almost) all genomes (Extended Data Fig. 3). Apparently, most of the Asgard  
212 genomes remain incomplete, such that conserved genes were missed randomly. Currently, there are only  
213 three gene families that are present in all Asgard MAGs, namely, a Zn-ribbon domain, a Threonyl-tRNA  
214 synthetase and an aminotransferase (Supplementary Data File 1).

215 The asCOG profiles were employed to annotate the remaining 86 Asgard MAGs, including those that  
216 were sequenced in the later stages of this work (Supplementary Table 1). On average, 89% of the proteins  
217 encoded in these genomes were covered by asCOGs (Supplementary Table 1). The protein annotation  
218 obtained using asCOGs was compared with the available annotation of '*Candidatus P. syntrophicum*'.  
219 Using asCOGs allowed at least a general functional prediction for 649 of the 1756 (37%) 'hypothetical  
220 proteins' in this organism, the only one in Asgard with a closed genome. We also identified 139 proteins,  
221 in addition to the 80 described originally, that can be considered Eukaryotic Signature Proteins, or ESPs  
222 (see below). Thus, the asCOGs database appears to be an efficient tool for annotation and comparative  
223 genomic analysis of Asgard MAGs and complete genomes.

224



225 **3. The core gene set of Asgard archaea**

226 The core set of conserved Asgard genes was arbitrarily defined as all asCOGs that are present in at least  
227 one-third of the MAGs in each of the 12 phylum-level lineages, with the mean representation across  
228 lineages >75%. Under these criteria, the Asgard core includes 378 asCOGs (Extended Data Fig. 6b ,  
229 Supplementary Table 6). As expected, most of these protein families, 293 (77%), are universal (present in  
230 bacteria, other archaea and eukaryotes), 62 (16%) are represented in other archaea and eukaryotes, but not  
231 in bacteria, 15 (4%) are found in other archaea and bacteria, but not in eukaryotes, 7 (2%) are archaea-  
232 specific, and only 1 (0.003%) is shared exclusively with eukaryotes (Extended Data Fig. 6c). Most of the  
233 core asCOGs show comparable levels of similarity to homologs from two or all three domains of life. The  
234 second largest fraction of the core asCOGs shows substantially greater sequence similarity (at least, 25%  
235 higher similarity score) to homologous proteins from archaea than to those from eukaryotes and/or  
236 bacteria (Supplementary Table 6). Compared with the 219 genes that comprise the pan-archaeal core <sup>71</sup>,  
237 the Asgard core set lacks 12 genes, each of which, however, is present in some subset of the Asgard  
238 genomes. These include three genes of diphthamide biosynthesis and 2 ribosomal proteins, L40E and  
239 L37E. The intricate evolutionary history of gene encoding translation elongation factors and enzymes of  
240 diphthamide biosynthesis in Asgard has been analyzed previously <sup>72</sup>. Also of note is the displacement of  
241 the typical archaeal glyceraldehyde-3-phosphate dehydrogenase (type II) by a bacterial one (type I) in  
242 most of the Asgard genomes (cog.001204, Supplementary Data File 1).

243 Functional distribution of the core asCOGs is shown in Extended Data Fig. 6b (also see Supplementary  
244 Data File 1). For comparison, we also derived an extended gene core for the TACK superphylum, using  
245 similar criteria (at least 50% in each of the 6 lineages and 75% of the genomes overall, Extended Data  
246 Fig. 6b). For at least half of the Asgard core genes, across most functional classes, there were no  
247 orthologs in the TACK core. The most pronounced differences were found, as expected, in the category U  
248 (intracellular trafficking, secretion, and vesicular transport). In Asgard archaea, this category includes 19  
249 core genes compared with 7 genes in TACK; 13 of these genes are specific to the Asgard archaea and  
250 include components of ESCRT I and II, 3 distinct Roadblock/longin families, 2 distinct families of small  
251 GTPases, and a few other genes implicated in related processes (Supplementary Table 6).

252

253 **4. Phylogenomic analysis of the Asgard superphylum and Asgard-eukaryote evolutionary**  
254 **relationship**

255 To assess the robustness of the Asgard phylogeny, 100 random, independent samples of 209 core asCOGs  
256 were generated, where each was sampled with the probability of  $1-e^{-1}$  (equivalent to making a bootstrap  
257 sample of the 209 asCOGs with each asCOG included once). The concatenated alignment of each  
258 subsampled set was then analyzed using IQ-tree and the automatically chosen evolutionary model. The  
259 support for the tree topology, derived from the full set of 209 core asCOGs, was estimated from the  
260 bipartitions in these 100 trees (Fig. 1a).

261 The analysis of the universal phylogeny aimed to make the species set for phylogenetic reconstruction as  
262 broadly representative as possible, while keeping its size manageable, to allow the use of powerful  
263 phylogenetic methods. The tree was constructed from alignments of 30 families of conserved orthologous  
264 proteins of 162 Asgard archaea, 286 other archaea, 98 bacteria and 72 eukaryotes (see Methods for details  
265 of the procedure including the selection of a representative species set and Supplementary Table 4). This  
266 set of universal genes, for which comparatively little HGT has been identified, was selected and  
267 employed for a classical reconstruction of the tree of life<sup>19</sup> as well as in many subsequent studies on deep  
268 phylogeny<sup>20,21</sup>. A concatenated alignment of 7411 positions was generated for these 30 protein families,  
269 after removing low information content positions (Supplementary Table 4). For the phylogenetic  
270 reconstruction, we used the IQ-tree program with several phylogenetic models (see Methods and  
271 Supplementary Table 5 for details). The resulting tree had the 3D topology, with high support values for  
272 all key bifurcations (Extended Data Fig. 7a and Supplementary Data File 2).

273 The effect of the phylogenetic marker choice on the tree topology was investigated by: 1) generating 30  
274 sets of 29 markers by removing each of the 30 markers from the original set and 2) generating 100  
275 bootstrap-like random, independent samples of markers (in the same manner as with the Asgard  
276 phylogeny markers). The concatenated alignment of each subsampled set of markers was then analyzed  
277 using IQ-tree and the automatically chosen evolutionary model. Unexpectedly, the 3D topology  
278 completely hinged on the presence of a single marker, COG0012 (Ribosome-binding ATPase YchF, an  
279 essential protein involved in translation)<sup>22</sup>. All samples that included COG0012 strongly supported the  
280 3D topology. In a sharp contrast, the set of 29 markers without COG0012 (Fig. 1c) and all 32 bootstrap-  
281 like samples that did not contain COG0012 equally strongly supported the 2D topology with varying  
282 placement of the eukaryote branch within archaea (Supplementary Table 5, Supplementary Data File 2).  
283 To further test the potential effect of branch-specific evolutionary rates on the tree topology, a 30-marker  
284 tree and a 29-marker tree without COG0012 were additionally constructed under a 4-category heterotachy  
285 model (LG+FO\*H4)<sup>73</sup>. The resulting trees showed essentially the same topologies as the trees in  
286 Extended Data Fig. 7a and Fig. 1c, respectively (Supplementary Data file 2).

287 To further assess the effect of the choice of the evolutionary model and the species selection on the tree  
288 topology, 100 trees were constructed from the same alignment (29 markers, concatenated) by randomly  
289 sampling 5 representatives of Asgard archaea, other archaea, bacteria, and eukaryotes each (using the best  
290 model selected by IQ-tree). All 100 trees showed 2D topology with at least one non-Asgard archaea clade  
291 separating eukaryotes from bacteria (Supplementary Table 5, Supplementary Data file 2).

292 To analyze the compatibility of the phylogenetic signals from the individual marker alignments with that  
293 of the 29-marker tree built from the concatenated alignment, 1000 random five-species samples were  
294 generated from each of the 29 protein families by sampling one sequence from Asgard archaea, one from  
295 the TACK archaea, one from other archaea, one from eukaryotes and one from bacteria. The topology of  
296 the five-species species trees was assessed for compatibility with the topology of the corresponding  
297 species in the 29-marker tree. The probabilities, associated with the Approximately Unbiased (AU) test <sup>74</sup>,  
298 implemented in the IQ-tree program, were averaged across the 1000 random samples. The trees for all 29  
299 markers were found to be compatible with the concatenated alignment topology, with the mean AU  
300 probabilities ranging from 0.489 to 0.513. In a sharp contrast, the COG0012 tree topology was  
301 incompatible with the 29-marker tree, with an AU probability of 0.00023 (Supplementary Table 5,  
302 Supplementary Data File 2).

303 To further probe the evolutionary relationships between eukaryotes, Asgard and other archaea, a  
304 bootstrap-like subsampling was performed on the set of the 29 markers excluding COG0012. Of the 100  
305 trees that were constructed from concatenated alignments of 14 to 23 sampled markers, 99 had the 2D  
306 topology and only one had the 3D topology (Supplementary Table 5, Supplementary Data File 2). Among  
307 these, in 62 trees, the eukaryote branch occupied different positions within the Asgard archaea, most  
308 commonly, as a sister group to the Heimdall-Gerd-Kari-Hod-Wukong clade (53 trees). In 23 trees,  
309 eukaryotes were a sister group to the Asgard-TACK clade (Supplementary Data File 2). None of the  
310 markers showed a strong association with a particular 2D tree topology akin to the dependence of the 3D  
311 topology on COG0012. The only apparent influential data point was COG0201 (Preprotein translocase  
312 subunit SecY) that was present in 40 of the 53 samples where eukaryotes grouped with the Heimdall  
313 archaea *sensu lato* and in only 6 of the 23 samples where eukaryotes were a sister branch to Asgard-  
314 TACK archaea. Not surprisingly, in the individual COG0201 tree, eukaryotes fell within the Heimdall-  
315 Gerd-Kari-Hod-Wukong clade, with a strong support (Extended Data Fig. 7c, Supplementary Data File  
316 2). Leave-one-out analysis of all possible 29 sets of 28 genes each fully supported the 2D topology and  
317 confirmed the lack of substantial effect of any of the individual genes (Supplementary Table 5). A  
318 consensus tree of the 100 bootstrap-like marker samples and 29 leave-one-out marker samples shows

319 eukaryotes branching from within Asgard archaea, as a sister to the Heimdall-Gerd-Kari-Hod-Wukong  
320 clade (Fig. 1d).

321 The sets of archaeal, bacterial and eukaryote species employed for the tree reconstruction were selected to  
322 ensure the widest possible representation of each of the domains and therefore include, among others,  
323 groups that consist (primarily) of parasites that tend to evolve fast and could hamper accurate  
324 phylogenetic reconstruction (DPANN archaea, mycoplasma, microsporidia and several others). In  
325 phylogenetic trees, the species from these groups usually form long branches with uncertain positions. To  
326 assess the potential effect of the inclusion of these species on the domain-level tree topology, we  
327 produced the tree from a reduced species sample set that excluded highly derived parasitic clades  
328 (Supplementary Table 4). In the tree constructed from the concatenated alignment of the 29 markers  
329 (excluding COG0012), eukaryotes grouped with the Heimdall-Gerd-Kari-Hod-Wukong branch within the  
330 Asgard clade (Extended Data Fig. 7d, Supplementary Table 5, Supplementary Data File 2).

331 In principle, HGT from eukaryotes to Asgard archaea or from Asgard to the ancestor of eukaryotes could  
332 produce a confounding signal resulting in apparent phylogenetic affinity between eukaryotes and Asgard  
333 and biasing the concatenated trees toward the 2D topology. To assess this possibility, a tree was  
334 constructed from a concatenated alignment of the 29 markers (excluding COG0012) excluding the Asgard  
335 sequences. In this tree, the eukaryotic branch confidently grouped with the TACK superphylum, well  
336 within the diversity of the extant Archaea (Supplementary Table 5, Extended Data Fig. 7e), indicating that  
337 for this set of markers, the 2D topology is robust and is not predicated on the specific Asgard-eukaryote  
338 affinity, whether reflecting common descent or HGT.

339

## 340 **5. Eukaryotic Signature Proteins in Asgard archaea**

341 The computational strategy for delineating an extensive yet robust ESP set is described under Methods.  
342 The set of identified ESPs contained 505 asCOGs, including 238 that were not closely similar (E-  
343 value= $10^{-10}$ , length coverage 75%) to those previously described by Zaremba-Niedzwiedzka et al. <sup>2</sup>  
344 (Supplementary Table 7). In a general agreement with previous observations, the majority of these ESPs,  
345 329 of the 505, belonged to the ‘Intracellular trafficking, secretion, and vesicular transport’ (U) functional  
346 class, followed by ‘Posttranslational modification, protein turnover, chaperones’ (O), with 101 asCOGs  
347 (Supplementary Table 7). Among the asCOGs in the U class, 130 were Roadblock/LC7 superfamily  
348 proteins, including longins, sybindin and profilins, and 94 were small GTPases of several families, such  
349 as RagA-like, Arf-like and Rab-like ones, as discussed previously <sup>25</sup>.

350 The phyletic patterns of ESP asCOGs in Asgard archaea are extremely patchy and largely lineage-specific  
351 (Extended Data Fig. 8), indicating that most of the proteins in this set are not uniformly conserved  
352 throughout Asgard evolution, but rather, are prone to frequent HGT, gene losses and duplications.  
353 Considering that the patchy distribution of the ESPs could be affected by genome incompleteness, this  
354 analysis was performed for the 76 Asgard genomes that were estimated to be at least 90% complete (see  
355 Materials and Methods for details). Capture of genes via HGT, gene loss and duplication are correlated in  
356 prokaryotes, resulting in the overall picture of dynamic evolution that is prominent in the U category  
357 COGs<sup>26</sup>. Even the most highly conserved ESP asCOG are missing in some Asgard lineages but show  
358 multiple duplications in others (Extended Data Fig. 8 and Supplementary Table 7).

359 Characteristically, many ESPs are multidomain proteins, with 37% assigned to more than one asCOG,  
360 compared to 17% among non-ESP proteins (Supplementary Table 7). Some multidomain ESPs in Asgard  
361 archaea have the same domain organizations as their homologs in eukaryotes, but these are a minority and  
362 typically contain only two domains. Examples include the fusion of two EAP30/Vps37 domains<sup>27</sup>, and  
363 Vps23 and E2 domains<sup>27</sup> in ESCRT complexes, multiple Rag family GTPases, in which longin domain is  
364 fused to the GTPase domain, and several others. However, most of the domain architectures of the  
365 multidomain ESP proteins were not detected in eukaryotes and often were found only in a narrow subset  
366 of Asgard archaea, suggesting extensive domain shuffling during Asgard evolution (Fig. 3a). For  
367 example, we identified many proteins containing a fusion of Vps28/Vps23 from ESCRT I complex<sup>27</sup>  
368 with C-terminal domains of several homologous subunits of adaptin and COPI coatomer complexes<sup>75,76</sup>,  
369 and E3 UFM1-protein ligase 1, which is involved in the UFM1 ubiquitin pathway<sup>77</sup> (Fig. 3a). Generally,  
370 a protein with such a combination of domains could be predicted to be involved in ubiquitin-dependent  
371 membrane remodeling but, because its domain architecture is unique, the precise function cannot be  
372 inferred.

373 The majority of the ESP genes of Asgard archaea do not belong to conserved genomic neighborhoods, but  
374 several such putative operons were detected. Perhaps, the most notable one is the ESCRT neighborhood  
375 which includes genes coding for subunits of ESCRT I, II and III, and often, components of the ubiquitin  
376 system<sup>2</sup>, suggesting an ancient link between the two systems that persists in eukaryotes<sup>27</sup>. Another  
377 predicted operon is conserved in most Asgard archaea and consists of genes encoding a LAMTOR1-like  
378 protein of the Roadblock superfamily, a Rab-like small GTPase, and a protein containing the DENN  
379 (differentially expressed in normal and neoplastic cells) domain that so far has been identified only in  
380 eukaryotes (Fig. 3b). Two proteins consisting of a DENN domain fused to longin are subunits of the  
381 folliculin (FLCN) complex that is conserved in eukaryotes. The FLCN complex is the sensor of amino  
382 acid starvation interacting with Rag GTPase and Ragulator lysosomal complex, and a key component of

383 the mTORC1 pathway, the central regulator of cell growth in eukaryotes<sup>78</sup>. Some Heimdallarchaea  
384 encode several proteins with the exact same domain organization as FLCN (Fig. 3b). Ragulator is a  
385 complex that consists of 5 subunits, each containing the Roadblock domain. In Asgard archaea, however,  
386 the GTPase present in the operon is from a family that is distinct from the Rag GTPases, which interact  
387 with both FLCN and Ragulator complexes in eukaryotes, despite the fact that Rag family GTPases are  
388 abundant in Asgard archaea<sup>25</sup> (Supplementary Table 7). Nevertheless, this conserved module of Asgard  
389 proteins is a strong candidate to function as a guanine nucleotide exchange factor for Rab and Rag  
390 GTPases, analogously to the eukaryotic FLCN. In eukaryotes, the DENN domain is present in many  
391 proteins with different domain architectures that interact with different partners and perform a variety of  
392 functions<sup>28,29</sup>. The Asgard archaea also encode other DENN domain proteins, and the respective genes  
393 form expanded families of paralogs in Loki, Hel and Heimdall lineages, again, with domain architectures  
394 distinct from those in eukaryotes (Fig. 3b)<sup>79</sup>.

395 Given the identification of a FLCN-like complex, a search was performed for other components of the  
396 mTORC1 regulatory pathway in Asgard archaea. The GATOR1 complex that consists of three subunits,  
397 Depdc5, Nprl2, and Nprl3, is another amino acid starvation sensor that is involved in this pathway in  
398 eukaryotes<sup>30</sup>. Nitrogen permease regulators 2 and 3 (NPRL2 and NPRL3) are homologous GATOR1  
399 subunits that contain a longin domain and a small NPRL2-specific C-terminal domain<sup>30</sup>. We identified a  
400 protein family with this domain organization in most Thor MAGs and a few Loki MAGs. Several other  
401 ESP asCOGs include proteins with high similarity to the longin domain of NPRL2. Additionally, we  
402 identified many fusions of the NPRL2-like longin domain with various domains related to prokaryotic  
403 two-component signal transduction system (Fig. 3c). Considering the absence of a homolog of  
404 phosphatidylinositol 3-kinase, the catalytic domain of the mTOR protein, it seems likely that, in Asgard  
405 archaea, the key growth regulation pathway remains centered at typical prokaryotic two-component signal  
406 transduction systems whereas at least some of the regulators and sensors in this pathway are “eukaryotic”.  
407 The abundance of NPRL2-like longin domains in Asgard archaea implies that the link between this  
408 domain and amino acid starvation regulation emerged at the onset of Asgard evolution if not earlier.

409

## 410 **6. Reconstruction of metabolic pathways in Asgard archaea**

411 Examination of the distribution of the asCOGs among the 12 Asgard archaeal phyla showed that the  
412 metabolic pathway repertoire was conserved among the MAGs of each phylum but differed between the  
413 phyla (Fig. 2a). Three distinct lifestyles were predicted by the asCOG analysis for different major  
414 branches of Asgard archaea, namely, anaerobic heterotrophy, facultative aerobic heterotrophy, and

415 chemolithotrophy (Fig. 4, Extended Data Fig. 9). For the last Asgard archaeal common ancestor  
416 (LAsCA), a mixotrophic lifestyle, including both production and consumption of H<sub>2</sub>, can be inferred from  
417 parsimony considerations (Fig. 4, Supplementary Table 8; see Methods for further details). Loki-, Thor-,  
418 Hermod-, Baldr- and Borrarchaeota encode all enzymes for the complete (archaeal) Wood-Ljungdahl  
419 pathway (WLP) and are predicted to oxidize organic substrates, likely, by using the reverse WLP, given  
420 the lack of enzymes for oxidation of inorganic compounds (e.g., hydrogen, sulfur/sulfide and  
421 nitrogen/ammonia). The genomes of these five Asgard phyla encode homologues of membrane-bound  
422 respiratory H<sub>2</sub>-evolving Group 4 [NiFe] hydrogenase and/or cytosolic cofactor-coupled bidirectional  
423 Group 3 [NiFe] hydrogenase<sup>35</sup>. Phylogenetic analysis of both group 4 and group 3 [NiFe] hydrogenases  
424 showed that Asgard archaea form distinct clades well separated from the functionally characterized  
425 hydrogenases, hampering the prediction of their specific functions in Asgard archaea (Extended Data Fig.  
426 10a and 10b, respectively). The functionally characterized group 4 [NiFe] hydrogenases in the  
427 Thermococci are involved in the fermentation of organic substrates to H<sub>2</sub>, acetate and carbon dioxide<sup>80,81</sup>.  
428 The presence of group 3 [NiFe] hydrogenases suggests that these Asgard archaea cannot use H<sub>2</sub> as an  
429 electron donor because they lack the enzyme complex coupling H<sub>2</sub> oxidation to membrane potential  
430 generation. Thus, in these organisms, bifurcate electrons from H<sub>2</sub> are likely to be used to support the  
431 fermentation of organic substrates exclusively<sup>80-82</sup>.

432 Both Wukongarchaeota genomes (As\_075 and As\_085) encode a bona fide membrane-bound Group 1k  
433 [NiFe] hydrogenase that could mediate hydrogenotrophic respiration using heterodisulfide as the terminal  
434 electron acceptor<sup>83,84</sup> (Fig. 4, Extended Data Fig. 9 and 10c). The group 1k [NiFe] hydrogenase is  
435 exclusively found in methanogens of the order Methanosarcinales (Euryarchaeota)<sup>66</sup>, and it is the first  
436 discovery of the group 1 [NiFe] hydrogenase in the Asgard archaea. Wukongarchaeota also encode all  
437 enzymes for a complete WLP and a putative ADP-dependent acetyl-CoA synthetase for acetate synthesis.  
438 Unlike all other Asgard archaea, Wukongarchaeota lack genes for citrate cycle and beta-oxidation. Thus,  
439 Wukongarchaeota appear to be obligate chemolithotrophic acetogens. The genomes of Wukongarchaeota  
440 were discovered only in seawater of the euphotic zone of the Yap trench (0 m and 125 m). Dissolved H<sub>2</sub>  
441 concentration is known to be the highest in surface seawater, where the active microbial fermentation,  
442 compared to deep sea<sup>85</sup>, could produce sufficient amounts of hydrogen for the growth of  
443 Wukongarchaeota. Hodarchaeota, Gerdarchaeota, Kariarchaeota, and Heimdallarchaeota share a common  
444 ancestor with Wukongarchaeota (Fig. 4). However, genome analysis implies different lifestyles for these  
445 organisms. Hod-, Gerd- and Kariarchaeota encode various electron transport chain components, including  
446 heme/copper-type cytochrome/quinol oxidase, nitrate reductase, and NADH dehydrogenase, most likely,  
447 allowing the use of oxygen and nitrate as electron acceptors during aerobic and anaerobic respiration,

448 respectively<sup>35</sup>. In addition, Hod-, Gerd- and Heimdallarchaeota encode phosphoadenosine phosphosulfate  
449 (PAPS) reductase and adenylylsulfate kinase for sulfate reduction, enabling the use of sulfate as electron  
450 acceptor during anaerobic respiration. Gerd-, Heimdall-, and Hodarchaeota are only found in coastal and  
451 deep-sea sedimentary environments, whereas Kariarchaeota were found also in marine water. The  
452 versatile predicted metabolic capacities of these groups suggest that Hod-, Gerd- and Kariarchaeota might  
453 occupy both anoxic and oxic niches. In contrast, Heimdallarchaeota appear to be able to thrive only in  
454 anoxic environments.

## 455 **References**

- 456 69. Makarova, K. S., Wolf, Y. I. & Koonin, E. V. Towards functional characterization of archaeal  
457 genomic dark matter. *Biochem. Soc. Trans.* **47**, 389–398 (2019).
- 458 70. Koonin, E. V. & Wolf, Y. I. Genomics of bacteria and archaea: the emerging dynamic view of the  
459 prokaryotic world. *Nucleic Acids Res.* **36**, 6688–6719 (2008).
- 460 71. Makarova, K. S., Wolf, Y. I. & Koonin, E. V. Archaeal Clusters of Orthologous Genes (arCOGs):  
461 An Update and Application for Analysis of Shared Features between Thermococcales, Methanococcales,  
462 and Methanobacteriales. *Life Basel Switz.* **5**, 818–840 (2015).
- 463 72. Narrowe, A. B. *et al.* Complex Evolutionary History of Translation Elongation Factor 2 and  
464 Diphthamide Biosynthesis in Archaea and Parabasalids. *Genome Biol. Evol.* **10**, 2380–2393 (2018).
- 465 73. Crotty, S. M. *et al.* GHOST: Recovering Historical Signal from Heterotachously Evolved  
466 Sequence Alignments. *Syst. Biol.* **69**, 249–264 (2020).
- 467 74. Shimodaira, H. An approximately unbiased test of phylogenetic tree selection. *Syst. Biol.* **51**,  
468 492–508 (2002).
- 469 75. Gomez-Navarro, N. & Miller, E. Protein sorting at the ER-Golgi interface. *J. Cell Biol.* **215**, 769–  
470 778 (2016).
- 471 76. Kibria, K. M. K. *et al.* A genome-wide analysis of coatmer protein (COP) subunits of  
472 apicomplexan parasites and their evolutionary relationships. *BMC Genomics* **20**, 98 (2019).
- 473 77. Wei, Y. & Xu, X. UFMylation: A Unique & Fashionable Modification for Life. *Genomics*  
474 *Proteomics Bioinformatics* **14**, 140–146 (2016).
- 475 78. Lawrence, R. E. *et al.* Structural mechanism of a Rag GTPase activation checkpoint by the  
476 lysosomal folliculin complex. *Science* **366**, 971–977 (2019).
- 477 79. Marat, A. L., Dokainish, H. & McPherson, P. S. DENN domain proteins: regulators of Rab  
478 GTPases. *J. Biol. Chem.* **286**, 13791–13800 (2011).
- 479 80. Yu, H. *et al.* Structure of an Ancient Respiratory System. *Cell* **173**, 1636–1649.e16 (2018).



- 480 81. Schut, G. J., Boyd, E. S., Peters, J. W. & Adams, M. W. W. The modular respiratory complexes  
481 involved in hydrogen and sulfur metabolism by heterotrophic hyperthermophilic archaea and their  
482 evolutionary implications. *FEMS Microbiol. Rev.* **37**, 182–203 (2013).
- 483 82. Bryant, F. O. & Adams, M. W. Characterization of hydrogenase from the hyperthermophilic  
484 archaeobacterium, *Pyrococcus furiosus*. *J. Biol. Chem.* **264**, 5070–5079 (1989).
- 485 83. Deppenmeier, U., Blaut, M., Schmidt, B. & Gottschalk, G. Purification and properties of a F420-  
486 nonreactive, membrane-bound hydrogenase from *Methanosarcina* strain Gö1. *Arch. Microbiol.* **157**, 505–  
487 511 (1992).
- 488 84. Thauer, R. K. *et al.* Hydrogenases from methanogenic archaea, nickel, a novel cofactor, and H<sub>2</sub>  
489 storage. *Annu. Rev. Biochem.* **79**, 507–536 (2010).
- 490 85. Conrad, R. & Seiler, W. Methane and hydrogen in seawater (Atlantic Ocean). *Deep Sea Res. Part*  
491 *Oceanogr. Res. Pap.* **35**, 1903–1917 (1988).
- 492