

 Open access • Posted Content • DOI:10.1101/2020.07.16.207365

## Early acquisition of conserved, lineage-specific proteins currently lacking functional predictions were central to the rise and diversification of archaea — [Source link](#)

[Raphaël Méheust](#), [Raphaël Méheust](#), [Cindy J. Castelle](#), [Alexander L. Jaffe](#) ...+1 more authors

**Institutions:** [University of California, Berkeley](#), [Planetary Science Institute](#)

**Published on:** 17 Jul 2020 - [bioRxiv](#) (Cold Spring Harbor Laboratory)

**Topics:** [Archaea](#) and [Phylogenetics](#)

Related papers:

- [Archaeal “Dark Matter” and the Origin of Eukaryotes](#)
- [Undinarchaeota illuminate the evolution of DPANN archaea](#)
- [A congruent phylogenomic signal places eukaryotes within the Archaea](#)
- [Rooting the Domain Archaea by Phylogenomic Analysis Supports the Foundation of the New Kingdom Proteoarchaeota](#)
- [The origin and evolution of Archaea: a state of the art](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/early-acquisition-of-conserved-lineage-specific-proteins-ejn3t0sqzb>

1  
2 **Early acquisition of conserved, lineage-specific proteins currently lacking functional**  
3 **predictions were central to the rise and diversification of archaea**  
4

5 Raphaël Méheust<sup>+1,4</sup>, Cindy J. Castelle<sup>1,2</sup>, Alexander L. Jaffe<sup>5</sup> and Jillian F. Banfield<sup>+1,2,3,4,5</sup>  
6

7 <sup>+</sup>Corresponding authors: [jbanfield@berkeley.edu](mailto:jbanfield@berkeley.edu), [raphael.meheust@berkeley.edu](mailto:raphael.meheust@berkeley.edu)  
8

9 <sup>1</sup>Department of Earth and Planetary Science, University of California, Berkeley, CA

10 <sup>2</sup>Chan Zuckerberg Biohub, San Francisco, CA

11 <sup>3</sup>Department of Environmental Science, Policy, and Management, University of California, Berkeley, CA

12 <sup>4</sup>Innovative Genomics Institute, University of California, Berkeley, CA

13 <sup>5</sup>Department of Plant and Microbial Biology, University of California, Berkeley, CA  
14  
15

16 **Abstract**

17 Recent genomic analyses of Archaea have profoundly reshaped our understanding of their  
18 distribution, functionalities and roles in eukaryotic evolution. Within the domain, major  
19 supergroups are Euryarchaeota, which includes many methanogens, the TACK, which includes  
20 Thaumarchaeota that impact ammonia oxidation in soils and the ocean, the Asgard, which  
21 includes lineages inferred to be ancestral to eukaryotes, and the DPANN, a group of mostly  
22 symbiotic small-celled archaea. Here, we investigated the extent to which clustering based on  
23 protein family content recapitulates archaeal phylogeny and identified the proteins that distinguish  
24 the major subdivisions. We also defined 10,866 archaeal protein families that will serve as a  
25 community resource. Clustering based on these families broadly recovers the archaeal  
26 phylogenetic tree. Interestingly, all major groups are distinguished primarily by the presence of  
27 families of conserved hypothetical proteins that are either novel or so highly diverged that their  
28 functions are obscured. Given that these hypothetical proteins are near ubiquitous within phyla,  
29 we conclude that they were important in the origin of most of the major archaeal lineages.  
30

## 31 **Introduction**

32           Until recently, the archaeal domain comprised only two phyla, the Euryarchaeota and the  
33 Crenarchaeota, most of which were described from extreme environments (Woese, Kandler, and  
34 Wheelis 1990; Woese and Fox 1977). The recovery of genomes from metagenomes without the  
35 prerequisite of laboratory cultivation has altered our view of diversity and function across the  
36 Archaea domain (Spang, Caceres, and Ettema 2017; Adam et al. 2017; Baker et al. 2020).  
37 Hundreds of genomes from little studied and newly discovered archaeal clades have provided  
38 new insights into archaeal metabolism and evolution. Now, Archaea include at least four major  
39 large groups, the Euryarchaeota (Cluster I and Cluster II; (Spang, Caceres, and Ettema 2017;  
40 Adam et al. 2017; Baker et al. 2020)), the TACK (Proteoarchaeota) (Petitjean et al. 2014), the  
41 Asgard (Spang et al. 2015; Zaremba-Niedzwiedzka et al. 2017), and the DPANN (Castelle et al.  
42 2015; Rinke et al. 2013), all of which comprise several distinct phylum-level lineages. These  
43 archaea are not restricted to extreme habitats, but are widely distributed in diverse ecosystems  
44 (Spang, Caceres, and Ettema 2017; Adam et al. 2017; Baker et al. 2020).

45           Most studies have focused on the metabolic potential of archaea based on analysis of  
46 proteins with known functions and revealed roles in the carbon, nitrogen, hydrogen and sulfur  
47 biogeochemical cycles. For example, Euryarchaeota includes many methanogens and non-  
48 methanogens, including heterotrophs and sulfur oxidizers (Offre, Spang, and Schleper 2013). The  
49 TACK includes Thaumarchaeota, most but not all of which oxidize ammonia (Pester, Schleper,  
50 and Wagner 2011; Brochier-Armanet et al. 2008), Aigarchaeota that tend to be chemolithotrophs  
51 that oxidize reduced sulfur compounds (Hua et al. 2018), Crenarchaeota that include thermophilic  
52 sulfur oxidizers (Woese et al. 1984), and Korarchaeota, a highly undersampled group represented  
53 by amino acid degraders that anaerobically oxidize methane and also metabolize sulfur  
54 compounds (McKay et al. 2019). The Asgard have variable metabolisms and their genomes  
55 encode pathways involved in structural components that are normally considered to be eukaryotic  
56 signatures (Spang et al. 2015; Zaremba-Niedzwiedzka et al. 2017). The DPANN are an intriguing  
57 group that typically have very small genomes and symbiotic lifestyles (Castelle et al. 2018;  
58 Dombrowski et al. 2019). Their geochemical roles are difficult to predict, given the predominance  
59 of hypothetical proteins.

60           Previously, the distribution of protein families over bacterial genomes was used to provide  
61 a function rather than phylogeny-based clustering of lineages (Méheust et al. 2019). Protein  
62 clustering allows the comparison of the gene content between genomes by converting amino acid  
63 sequences into units of a common language. The method is agnostic and unbiased by  
64 preconceptions about the importance or functions of genes.

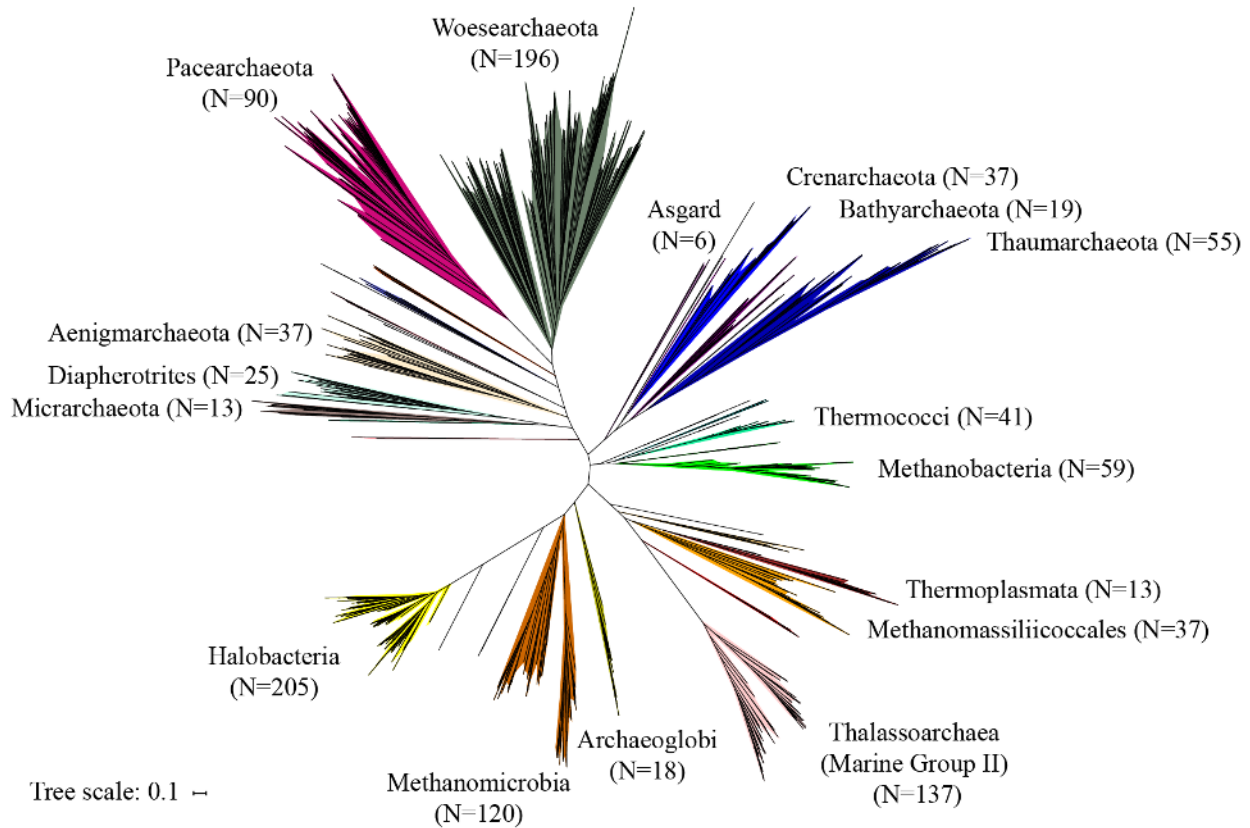
65 Here, we adapted this approach to evaluate the protein family-based coherence of the  
66 archaea and to test the extent to which a subdivision of archaea could be resolved based on  
67 shared protein family content. The analysis drew upon the large genome dataset that is now  
68 available for cultivated as well as uncultivated archaea (3,197 genomes). The observation that  
69 hypothetical proteins dominate the sets of co-occurring protein families that distinguish major  
70 archaeal groups indicates the importance of these protein sets in the rise of the major archaeal  
71 lineages.

72

## 73 **Results**

### 74 **Genome reconstruction and collection improved the resolution of the DPANN lineages**

75 We collected 2,618 genomes spanning all the recognized phyla and superphyla of the  
76 Archaea domain from the NCBI genome database (**Supplemental Dataset - Table S1**). To  
77 enable our analyses, we augmented the relatively limited sampling of the DPANN by adding 569  
78 newly available DPANN draft genomes (Castelle et al. in prep.) from low oxygen marine  
79 ecosystems, an aquifer adjacent to the Colorado River, Rifle, Colorado, and from groundwater  
80 collected at the Genasci dairy farm, Modesto CA (He et al., n.d.). The 3,197 genomes were  
81 clustered at  $\geq 95\%$  average nucleotide identity (ANI) to generate 1749 clusters. We removed  
82 genomes with  $<70\%$  completeness or  $>10\%$  contamination or if there was  $< 50\%$  of the expected  
83 columns in the alignment of 14 concatenated ribosomal proteins (**see Materials and Methods**).  
84 To avoid contamination due to mis-binning, we required that these proteins were co-encoded on  
85 a single scaffold. The average completeness of the final set of 1,179 representative genomes is  
86 95% and 928 were  $>90\%$  complete (**Supplementary Dataset - Table S1**). The 1,179  
87 representative genomes comprise 39 phylum-level lineages including 16 phyla that have more  
88 than 10 genomes (**Supplementary Dataset - Table S1 and Figure 1**).



89

90 **Figure 1. Phylogenetic tree of the 1,179 representative genomes. The maximum-likelihood**  
91 **tree was calculated based on the concatenation of 14 ribosomal proteins (L2, L3, L4, L5,**  
92 **L6, L14, L15, L18, L22, L24, S3, S8, S17, and S19) using the LG plus gamma model of**  
93 **evolution. Scale bar indicates the average substitutions per site. The complete ribosomal**  
94 **protein tree is available in rectangular format in Supplementary Figure 1.**

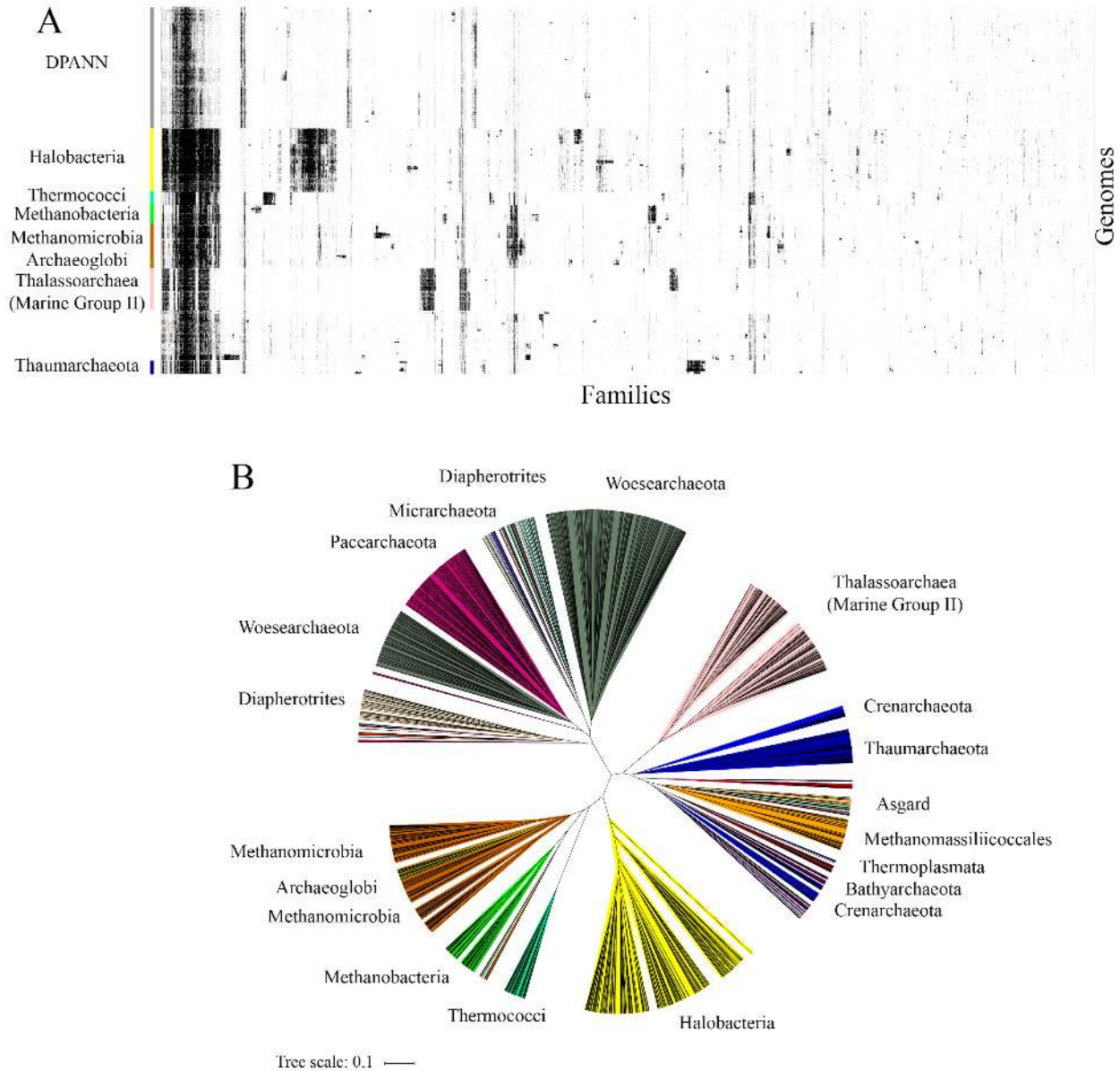
95

### 96 **Genomic content of representative genomes correlates with the phylogeny of archaea**

97 We clustered the 2,336,157 protein sequences from the representative genomes in a two-  
98 step procedure to generate groups of homologous proteins (**Supplementary Figure 2**). This  
99 resulted in 10,866 clusters (representing 2,075,863 sequences) that were present in at least five  
100 distinct genomes. These clusters are henceforth referred to as protein families.

101 We assessed the quality of the protein clustering. The rationale was that we expected  
102 protein sequences with the same function to cluster into the same protein family. We annotated  
103 our protein dataset using the KEGG annotations (Kanehisa et al. 2016) and systematically verified  
104 that the protein family groupings approximate functional annotations. The KEGG annotations in  
105 our dataset encompass 6,482 unique annotations of various biological processes, including fast-  
106 evolving defense mechanisms. For each of these 6,482 annotations, we reported the family that

107 contains the highest percentage of protein members annotated with that KEGG annotation. Most  
108 clusters were of good quality. For 87% of the KEGG annotations (5,627 out of 6,482), one family  
109 always contained >80% of the proteins (**Supplementary Figure 3A**). The contamination of each  
110 protein family was assessed by computing the percentage of the proteins with KEGG annotations  
111 that differ from the dominant annotation (percentage annotation admixture). Most of the families  
112 contain only proteins with the same annotation, and 2,654 out of 3,746 families (71%) have <20%  
113 annotation admixture (**Supplementary Figure 3B**). Although this metric is useful, we note that it  
114 is imperfect because two homologous proteins can have different KEGG annotations and thus  
115 cluster into the same protein family, increasing the apparent percentage of annotation admixture.  
116 Although we used sensitive HMM-based sequence-comparison methods and assessed the  
117 quality of the protein clustering, we cannot completely rule out the possibility that our pipeline  
118 failed to retrieve distant homology for highly divergent proteins. Small proteins and fast-evolving  
119 proteins are more likely to be affected. This lack of sensitivity would result in the separation of  
120 homologous proteins into distinct families and would impact the results. To reduce the incidence  
121 of proteins without functional predictions for which annotations should have been achieved we  
122 augmented PFAM and KEGG-based annotations by comparing sequences to PDB database and  
123 by performing HMM-HMM comparison against the eggNog database (**see Materials and**  
124 **Methods**).  
125



126

127 **Figure 2. The distribution of the 10,866 families across the 1,179 representative genomes.**

128 **A. The distribution of 10,866 widely distributed protein families (columns) in 1,179**  
129 **representative genomes (rows) from Archaea. Data are clustered based on the presence**  
130 **(black) and absence (white) profiles (Jaccard distance, complete linkage). B. Tree resulting**  
131 **from the hierarchical clustering of the genomes based on the distributions of proteins**  
132 **families in panel A.**

133

134 We visualized the distribution of the families over the genomes by constructing an array  
135 of the 1,179 representative genomes (rows) vs. 10,866 protein families (columns) and



136 hierarchically clustered the genomes based on profiles of protein family presence/absence  
137 **(Figure 2A)**. The families were also hierarchically clustered based on profiles of genome  
138 presence/absence. As previously reported for bacteria (Snel, Bork, and Huynen 1999; Méheust  
139 et al. 2019), the hierarchical clustering tree of the genomes resulting from the protein clustering  
140 **(Figure 2B)** correlated with the maximum-likelihood phylogenetic tree based on the concatenation  
141 of the 14 ribosomal proteins **(Figure 1)** (the cophenetic correlation based on a complete-linkage  
142 method is 0.83, based on average-linkage 0.84, and based on single-linkage, 0.84)  
143 **(Supplementary Figure 4)**. Although the tree resulting from the protein families correlates with  
144 the phylogenetic tree, it does not achieve the resolution of the phylogenetic tree, especially for  
145 placement of the deep branches. Interestingly, several phyla, such as the Crenarchaeota or the  
146 Woesarchaeota, are resolved into multiple groups **(Figure 2B)**. The first clade of Woesearchaeota  
147 corresponds to the Woesarchaeota-like I whereas the second clade groups together the  
148 Woesarchaeota and Woesarchaeota-like II groups. We could not evaluate the placement of  
149 Altiarchaeota relative to the DPANN because no genomes passed our quality control thresholds.

150 We defined modules as blocks of co-occurring protein families containing at least 20  
151 families (see **Materials and Methods**) (Méheust et al. 2019). Each module was assigned a  
152 taxonomic distribution based on the taxonomy of the genomes with the highest number of families  
153 (see **Materials and Methods and Supplementary Dataset - Table S2**). A block of 587 protein  
154 families that was broadly conserved across the 1,179 genomes (left side in **Figure 2A**) was  
155 designed as the module of 'core families' (Module 1) **(Supplementary Figure 5)**. Given their  
156 widespread distribution, it is unsurprising that most of the families are involved in well-known  
157 functions, including replication, transcription and translation, basic metabolism (oxidative  
158 phosphorylation chain, nucleotides, amino acids, ribosomal proteins, cofactors and vitamins,  
159 transporters, peptidases, DNA repair and chaperones). As expected, many of these easily  
160 recognized core families, primarily those involved in energy metabolism and cofactor synthesis,  
161 are absent in DPANN genomes (Castelle et al. 2018, 2015) **(Figure 2A and Supplementary**  
162 **Dataset - Table S3)**. Another interesting module (module 23) **(Supplementary Figure 5)**,  
163 composed of ~100 protein families, is widely distributed in most archaeal genomes but was not  
164 identified in DPANN and surprisingly, not in the Thalamoarchaea. Module 23 includes functions  
165 involved in carbon metabolism, amino-acid synthesis, and many transporter families. For  
166 instance, we identified several families for subunits of the Mrp antiporter as widespread in  
167 Halobacteria, Methanogens and Thermococci, but they appear to be absent in DPANN and  
168 Thalamoarchaea. The Mrp antiporter functions as Na<sup>+</sup>/H<sup>+</sup> antiporter and also contributes to



169 sodium tolerance in Haloarchaea. Mrp has been reported to be involved in energy conservation  
170 in methanogens and in the metabolic system of hydrogen production in Thermococci.

171 The DPANN are an enigmatic set of lineages, the monophyly of which remains uncertain  
172 (Aouad et al. 2018). However, the protein family analysis clearly showed that these lineages group  
173 together and are distinct from other Archaea (**Figure 2B**). A detailed protein family analysis of  
174 groups within the DPANN is presented elsewhere (Castelle et al. in prep.).

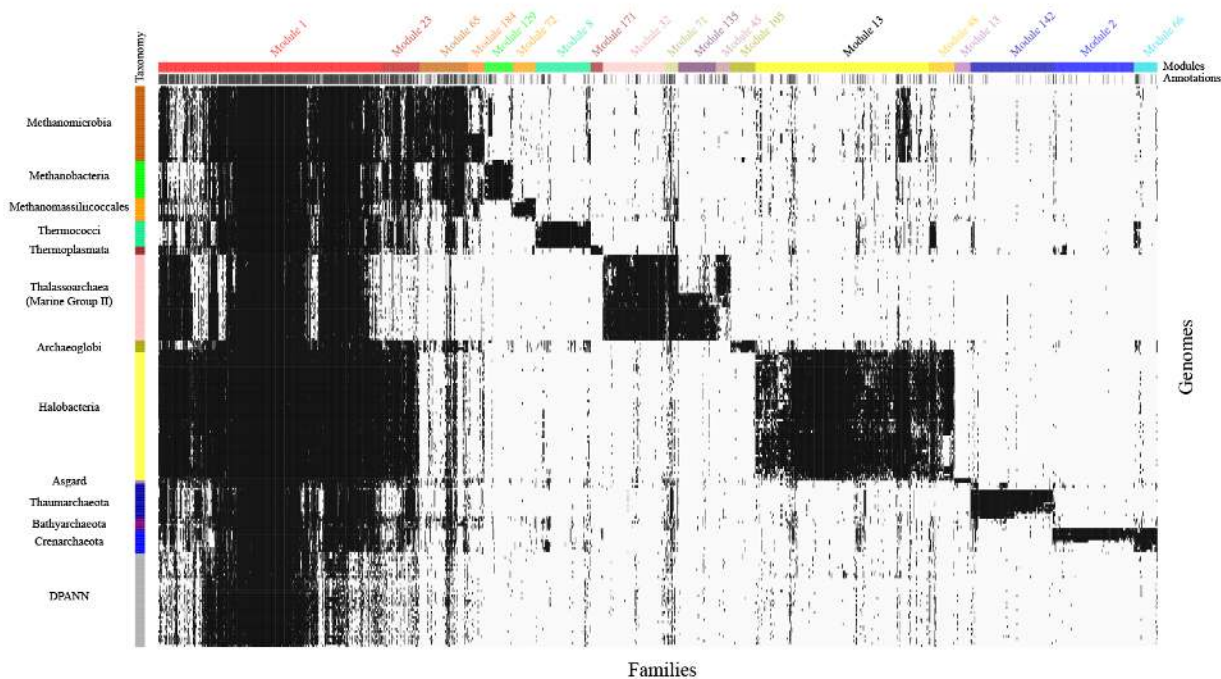
175

### 176 Major non-DPANN groups possess groups of conserved protein families.

177

178 We detected 96 modules that are restricted to non-DPANN lineages (**Supplementary**  
179 **Dataset - Table S2**). Only 9 of the 96 modules were found in multiple phyla and in 8 of these 9  
180 cases, the phyla that possess each module are phylogenetically unrelated (e.g., Crenarchaeota  
181 and Halobacteria). The 9th, module 44, is interesting in that it occurs in two phyla and those phyla  
182 are monophyletic (Thorarchaeota and Heimdallarchaeota of the Asgard superphylum). Thus, the  
183 vast majority of the non-DPANN modules (87) are restricted to a single phylum (**Supplementary**  
184 **Dataset - Table S2**) and, perhaps surprisingly given phylogenetic support for superphyla within  
185 Archaea, almost no modules are specific to superphyla.

186



187

188 **Figure 3. The distribution of the 2,632 families of the 19 modules discussed in this study.**

189 **Each column represents a protein family and each row represents a genome. Data are**

190 **clustered based on the presence (black)/ absence (white) profiles but also based on the**  
191 **taxonomy of the genomes and the module membership. The first colored top bar**  
192 **(annotations) shows the families with (black) / without (white) a predicted annotation**  
193 **whereas the second colored top bar (modules) indicates the module of each family. The**  
194 **colored side bar indicates the taxonomic assignment of each genome.**

195

196 Visualization of the distribution of protein families highlights the presence of modules that  
197 are not only lineage specific but are also well conserved within each lineage (**Figure 2A**). In fact,  
198 we identified such archaeal group-specific modules in 10 out of 11 non-DPANN with more than  
199 10 genomes (**Table 1, Figure 3 and Supplementary Figure 5**). For instance, there are two  
200 modules (modules 13 and 108) comprising 525 families that are fairly conserved in Halobacteria.  
201 On average, each of the 525 families appears in 65% of the halobacterial genomes, yet these  
202 families are mostly absent in non-halobacterial genomes (**Supplementary Figure 6**). These  
203 modules are slightly less conserved within each archaeal group than module 1 families  
204 (comprising core functions) (**Supplementary Figure 6**).

205

#### 206 **Do methanogens cluster together, despite their phylogenetic diversity?**

207 We identified one module of 128 protein families, module 65 (**Figure 3 and**  
208 **Supplementary Dataset - Table S3**), that is common to essentially all methanogens, despite the  
209 fact that methanogens are not monophyletic (**Figure 1**). This module contained *mcrA*  
210 (Fam05485), a key gene in methane production (Ermler et al. 1997) all the other subunits (BCDG)  
211 of methyl-coenzyme M reductase (Mcr), five subunits of the methyl-tetrahydromethanopterin  
212 (methyl-H4MPT): coenzyme M methyltransferase (Mtr), five hypothetical conserved proteins in  
213 methanogens (Borrel et al. 2014) and genes for transport of iron, magnesium, cobalt and nickel  
214 and for synthesis of key cofactors that are required for growth of methanogens. Details are  
215 provided in the **Supplementary Materials**.

216 Modules 72, 129 and 184 are enriched in subunits of the energy-converting hydrogenase  
217 A and B and in enzymes for the utilization of methanol (fam04064 and fam05405), methylamine  
218 (fam02336 and fam03937), dimethylamine (fam03076 and fam05873), and trimethylamine  
219 (fam04092 and fam21299), which are substrates for methanogenesis (Burke, Lo, and Krzycki  
220 1998) (for details, see the **Supplementary Materials**).

221 Interestingly, we recovered *mcr* subunits in lineages that are not considered as canonical  
222 methanogenic lineages (Evans et al. 2019). These include two genomes of Bathyarchaeota  
223 related to BA1 and BA2 (GCA\_002509245.1 and GCA\_001399805.1) (Evans et al. 2015), and

224 one Archaeoglobi genome related to JdFR-42 (GCA\_002010305) (Boyd et al. 2019; Wang et al.  
 225 2019). These genomes have been described as having divergent MCR genes. It is reassuring  
 226 that our method is sensitive enough to recover distant homology. Overall, the correspondence  
 227 between the distribution of protein families linked to methanogenesis and methanogens supports  
 228 the validity of our protein family delineation method (**Supplementary Figure 7**).

Modules	Lineage(s)	# Families	SignalP (%)	TMHMM (%)	Hypothetical families (%)	Hits to Bacteria (%)
1	Core genome	587	6	20	13	87
13,108	Halobacteria	525	14	36	82	34
66,2	Crenarchaeota	276	9	34	89	11
142	Thaumarchaeota	216	13	31	94	11
32,71	Marine Group II	199	19	55	77	32
8	Thermococci	146	12	32	84	24
65	Methanomicrobia	128	11	22	45	63
129	Methanobacteria	75	16	55	71	17
105	Archaeoglobi	65	3	40	94	12
72	Methanomassiliococcales	59	22	49	86	27
48	Asgard	42	17	36	79	17
171	Thermoplasmata	32	3	38	97	3

229  
 230 **Table 1. A list of the fourteen modules that are lineage specific but also well conserved**  
 231 **within eleven major archaeal lineages. A family was counted as having a signal peptide if**  
 232 **at least 25% of its protein sequences were predicted to have a signal peptide prediction**  
 233 **according to the SignalP software (Almagro Armenteros et al. 2019). A family was counted**  
 234 **as having a transmembrane helix if more than half of its protein sequences were predicted**  
 235 **to have a transmembrane helix according to the TMHMM software (Krogh et al. 2001).**  
 236 **Families were considered hypothetical if they have neither PFAM (Domain of Unknown**  
 237 **Function domains were excluded) nor KEGG annotations (see the supplementary dataset**  
 238 **- Table S3 for the full list of hypothetical families). Finally, a family was considered to have**  
 239 **bacterial homologs if the family matched with protein sequences of at least ten distinct**  
 240 **bacterial genomes (see Materials and Methods). The core module 1 is included as a**  
 241 **comparison.**

242

### 243 **Functions specific to Thalamoarchaea**

244 Modules 32 and 71, encompassing 199 families, were consistently associated with  
 245 genomes of Thalamoarchaea (**Figure 3 and Supplementary Dataset - Table S3**), which are  
 246 implicated in protein and saccharide degradation (Tully 2019) (for details, see the **Supplementary**  
 247 **Materials**). These modules contain protein degrading enzymes (several different classes of

248 peptidases and one oligotransporter) previously found in Thalamoarchaea (Tully 2019) and two  
249 new Thalamoarchaea-specific families of well-conserved peptidases. As reported by (Tully 2019),  
250 peptidase S15 (PF02129; fam03321) and peptidase M60-like (PF13402; fam05454) have a  
251 narrow distribution within Thalamoarchaea, and were not assigned to ones of the 96 modules.  
252 Interestingly, we identified modules specific to Thalamoarchaea subgroup a (MGIIa)  
253 (module\_135, containing 99 families) and Thalamoarchaea subgroup b (MGIIb) (module\_45,  
254 containing 39 families) with calcium-binding domains (**Supplementary Figure 8**). These proteins  
255 may be involved in signaling and regulation of protein-protein interactions in the cell (Michiels et  
256 al. 2002).

257

### 258 **Functions specific to Crenarchaeota**

259 The Crenarchaeota comprises thermophilic organisms that are divided into three main  
260 classes, the Thermoproteales, the Sulfolobales and the Desulfurococcales. Two distinct modules  
261 with distinct distributions were retrieved. Module 66 (61 families) is widespread in the three  
262 classes of Crenarchaeota whereas module 2 (215 families) is specific to the Sulfolobales class  
263 (**Figure 3 and Supplementary Dataset - Table S3**). Interestingly, the subunits of RNA  
264 polymerase (Korkhin et al. 2009), RpoG/Rpb8 (fam03177) are widespread in Crenarchaeota but  
265 Rpo13 (fam03159) seems restricted to the *Sulfolobales* class (Korkhin et al. 2009). The Rpo13  
266 protein family of Thermoproteales and Desulfurococcales may be highly divergent from the form  
267 described experimentally.

268 Comparison to PDB enabled annotation of three families with no PFAM and KEGG  
269 annotations as having functions related to the DNA replication machinery (**Supplementary**  
270 **Dataset - Table S4**). We were interested to find that this ubiquitous function is performed by  
271 specific protein families in Crenarchaeota, possibly reflecting adaptation to their high temperature  
272 habitats. One of these, PolB1-binding protein 2 (PBP2) (fam03141, PDB accession 5n35) (Yan  
273 et al. 2017), is a subunit of DNA polymerases B1 (PolB1) that are responsible for initial RNA  
274 primer extension with DNA, lagging and leading strand synthesis. The second is a single-stranded  
275 DNA-binding protein (DBP) ThermoDBP, which we also found to be conserved in Crenarchaeota  
276 and in Thermococci (fam03176, PDB accession 4psl) (Ghalei et al. 2014; Paytubi et al. 2012).  
277 Interestingly, however, the third is a Fe-S independent primase subunit PriX (fam03870, PDB  
278 accessions: 4wyh and 5of3) specific to Sulfolobales (**Supplementary Figure 9**). PriX is essential  
279 for the growth of *Sulfolobus* cells (Holzer et al. 2017; Liu et al. 2015). These observations point to  
280 fundamentally different transcription and replication mechanisms in the major groups within the  
281 Crenarchaeota.

282           Restricted to the Sulfolobales are also two multicopy thermostable acid protease  
283 thermopsin families (Lin and Tang 1990) (fam01298 and fam01602 in module 2). Fam01298 is  
284 also found in two genomes of Thermoproteales (**Supplementary Figure 9**). Extending a prior  
285 report that Crenarchaeota have anomalously large numbers of types I and III CRISPR-Cas  
286 systems (Vestergaard, Garrett, and Shah 2014), Crenarchaeota-specific module 66 contains four  
287 type I-A Cas families (one of which is the sulfolobales-specific CRISPR-associated protein csaX,  
288 fam07252) and four Cas families associated with type III systems (**Supplementary Figure 9**)  
289 (**Supplementary Dataset - Table S3**).

290

### 291 **Functions specific to Thaumarchaeota**

292           The phylum Thaumarchaeota mostly contains aerobic ammonia oxidizing archaea  
293 (Brochier-Armanet et al. 2008; Adam et al. 2017). Module 142, which contains 216 families, is  
294 specific to Thaumarchaeota (**Figure 3 and Supplementary Dataset - Table S3**). Although this  
295 module contains protein families for the three subunits of the ammonia monooxygenase, these  
296 three families are absent in genomes for two basal Thaumarchaeota lineages, as expected based  
297 on prior analyses (Adam et al. 2017) (**Supplementary Figure 10**). This module also contains a  
298 highly conserved hypothetical family (fam08021), referred to as AmoX (Bartossek et al. 2012),  
299 that is known to co-occur with the amoABC genomic cluster (**Supplementary Dataset - Table**  
300 **S5**). Importantly, essentially all other protein families in Module 142 currently lack functional  
301 annotations (**Supplementary Dataset - Table S3**).

302

### 303 **Functions specific to Thermococci**

304           The Thermococci comprises sulfur-reducing hyperthermophilic archaea (Palaeococcus,  
305 Thermococcus and Pyrococcus). Module 8 contains 146 families abundant in Thermococci and  
306 absent in other archaeal lineages (**Figure 3 and Supplementary Dataset - Table S3**). For  
307 example, 98% of the Thermococci genomes have a group 3b (NADP-reducing) [NiFe]  
308 hydrogenase. This hydrogenase, also known as sulfhydrogenase, is likely bidirectional (Schut et  
309 al. 2012). Only the subunit beta of the sulfur reductase (fam04571) is present in module 8.  
310 Subunits alpha (fam00341), delta (fam00630) and gamma (fam00435) are present in the core  
311 module (module 1), probably because they are homologs of other hydrogenases. We also  
312 detected hydrogen gas-evolving membrane-bound hydrogenases (MBH) in every Thermococci  
313 genome (fam03754 in module 8) (Yu et al. 2018; Schut et al. 2016) (**Supplementary Figure 11**).  
314 The MBH transfers electrons from ferredoxin to reduce protons to form H<sub>2</sub> gas (Sapra,  
315 Bagramyan, and Adams 2003). The Na<sup>+</sup>-translocating unit of the MBH enables H<sub>2</sub> gas evolution



316 by MBH to establish a Na<sup>+</sup> gradient for ATP synthesis near 100 °C in *Pyrococcus furiosus* (Yu et  
317 al. 2018). As with the sulfhydrogenase, only the subunit I of the MBH is present in module 8, other  
318 subunits of MBH are present in core modules 1 and 23 probably because MBH-type respiratory  
319 complexes are evolutionarily and functionally related to the Mrp H<sup>+</sup>/Na<sup>+</sup> antiporter system (Yu et  
320 al. 2018).

321 In the Thermococci-specific module 8 we detected the alpha and gamma subunits  
322 (represented by fam10869 and fam02435, respectively) of the Na<sup>+</sup>-pumping methylmalonyl-  
323 coenzyme A (CoA) decarboxylase that performs Na<sup>+</sup> extrusion at the expense of the free energy  
324 of decarboxylation reactions (Dimroth 1987; Buckel 2001). The beta and delta subunit, fam02317  
325 and fam00273, are present in the core module 1, again probably because they are homologs of  
326 proteins that perform different functions.

327 Interestingly, three families from module 8 are encoded adjacent in the Thermococci  
328 genomes (fam15060, fam07594 and fam05926) (**Supplementary Dataset - Table S6**). These  
329 are annotated as of fungal lactamase (renamed prokaryotic 5-oxoprolinase A, pxpA) and  
330 homologs of allophanate hydrolase subunits (renamed pxpB and pxpC) and are likely to form  
331 together an 5-oxoprolinase complex (Niehaus et al. 2017). While oxoproline is a major universal  
332 metabolite damage product and oxoproline disposal systems are common in all domains of life,  
333 the system encoded by these three families appears to be highly conserved in Thermococci  
334 (**Supplementary Figure 11**).

335 We found the ribosomal protein L41e (fam02171) (Yutin et al. 2012) in 83% of the  
336 genomes of Thermococci but sparsely distributed or absent in other lineages. It has previously  
337 been noted that the distribution of L41e in Archaea is uncertain (Lecompte et al. 2002).

338 Using PDB, we established annotations for three families in Thermococci-specific module  
339 8 that lacked PFAM or KEGG annotations (**Supplementary Dataset - Table S4**). The first  
340 appears to be a small protein that inhibits the proliferating cell nuclear antigen by breaking the  
341 DNA clamp in *Thermococcus kodakarensis* (fam09868) (Altieri et al. 2016). The second is the S  
342 component of an energy-coupling factor (ECF) transporter (fam02033) likely responsible for  
343 vitamin uptake (Zhang, Wang, and Shi 2010). The third (fam01133) is the Valosin-containing  
344 protein-like ATPase (VAT) that in *Thermoplasma acidophilum* functions in concert with the 20S  
345 proteasome by unfolding substrates and passing them on for degradation (Huang et al. 2016).  
346 Finally, three peptidases were detected in module 8 (fam01338, fam26972 and fam05052), thus  
347 may be specific to the Thermococci (**Supplementary Figure 11**).

348  
349



## 350 **Functions specific to Halobacteria**

351 We found that 525 families comprise the Halobacteria-specific modules 13 and 108.  
352 Module 108 is composed almost completely of hypothetical proteins (**Figure 3 and**  
353 **Supplementary Dataset - Table S3**).

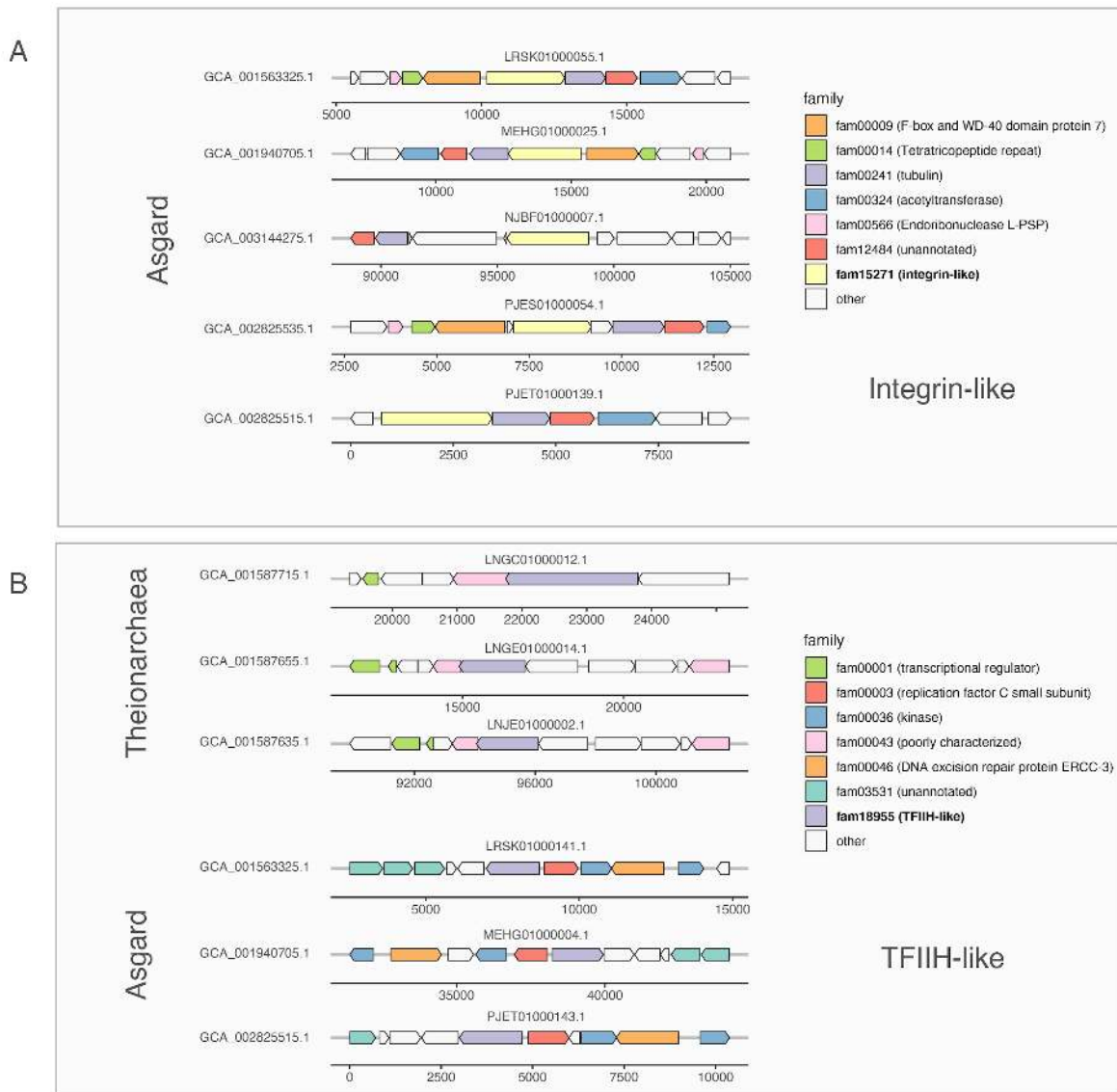
354 Module 13 contains the two subunits I (fam02395) and II (fam06634) of the high-affinity  
355 oxygen cytochrome *bd* oxidase (module 13) and was identified in half of the genomes  
356 (**Supplementary Fig. 12**). It also contains three families without KEGG and PFAM annotations,  
357 but close inspection using HMM-HMM comparison showed that they have distant homology with  
358 cytochrome-related proteins (**Supplementary Dataset - Table S4**). The first, fam02696, has  
359 distant homology with the catalytic subunit I of heme-copper oxygen reductases (fam00581) and  
360 the genes often colocalize with heme-copper oxygen reductases-related genes such as type C  
361 (*cbb<sub>3</sub>*) subunit I or the nitric oxide reductase subunit B (fam00581) (**Supplementary Dataset -**  
362 **Table S7**). The two other families are cytochrome c associated proteins (fam01001, cytochrome  
363 c biogenesis factor and fam02143, cytochrome C and quinol oxidase polypeptide I). Consistent  
364 with the presence of oxygen respiration-related families, a catalase-peroxidase gene is present  
365 in 90% (fam02210) of the halobacteria genomes (**Supplementary Fig. 12**). Module 13 also  
366 includes proteins for synthesis of proteinaceous gas vacuoles (fam03834, fam03740, fam02854  
367 and fam00889; identified in more than 45% of halobacterial genomes, **Supplementary Dataset**  
368 **- Table S3**) that regulate buoyancy of cells in aqueous environments (DasSarma and Arora 2006).  
369 The module also includes bacterioruberin 2'', 3''-hydratase (fam00736, CruF; identified in 97% of  
370 the halobacteria genomes). Adjacent in the Halobacteria genomes are two families found in the  
371 core module 1 (fam00008 and fam00115) and annotated as digeranylgeranyl glycerophospholipid  
372 reductase and UbiA prenyltransferases respectively (**Supplementary Dataset - Table S7**).  
373 Closer inspection of these three co-encoded enzymes in *Haloarcula japonica* DSM 6131  
374 (GCA\_000336635.1) showed they are identical with the bifunctional lycopene elongase and 1,2-  
375 hydratase (LyeJ, fam00115) and the carotenoid 3,4-desaturase (CrtD, fam00008) and the  
376 bacterioruberin 2'', 3''-hydratase (CruF, fam00736) genes described in *Haloarcula japonica* JCM  
377 7785<sup>T</sup> (Yang et al. 2015). Together, these three enzymes can generate C50 carotenoid  
378 bacterioruberin from lycopene in *Haloarcula japonica* (Yang et al. 2015). Our results showed that  
379 C50 carotenoid bacterioruberin is highly conserved in Halobacteria (**Supplementary Figure 12**).

380

## 381 **Functions specific to the six Asgard genomes.**

382 The module 48 contains 42 families that are specific and conserved in the six genomes of  
383 the superphylum Asgard (four genomes of Thorarchaeota and two genomes of

384 Heimdallarchaeota (**Figure 3**). Of these, 33 lack both KEGG and PFAM functional predictions  
 385 (**Supplementary Dataset - Table S3**). The Asgard archaea, which affiliate with eukaryotes in the  
 386 tree of life (Cox et al. 2008), encode many proteins that they share with eukaryotes (Hartman and  
 387 Fedorov 2002). We detected four eukaryotic signature protein families (ESPs) in module 48 that  
 388 were described in previous studies (**Supplementary Figure 13 and Supplementary Materials**)  
 389 (Zaremba-Niedzwiedzka et al. 2017; Spang et al. 2015; Akil and Robinson 2018).  
 390



391  
 392 **Figure 4. Schematic overview of integrin-like and TFIIH-like gene clusters identified in**  
 393 **archaea. A. Conserved gene clusters comprising archaeal integrin-like genes (fam15271)**

394 **identified in five Asgard genomes. B. Conserved gene clusters comprising archaeal TFIIH-**  
395 **like genes (fam18955) identified in three Theionarchaea and three Asgard genomes. A full**  
396 **gene synteny and genomic context of the genes neighboring the integrin-like (fam15271)**  
397 **and TFIIH-like (fam18955) genes is available in Supplementary Dataset - Table S8.**

398  
399 Interestingly, we found a family in module 48 (fam15271) that shows sequence similarity  
400 with the integrin beta 4. To the best we know, integrin genes were never described in archaea  
401 and fam15271 may represent a new ESP. The genes of fam15271 are always located next to  
402 tubulin genes (fam00241) in the five Asgard genomes (**Figure 4A and Supplementary Dataset**  
403 **- Table S8**). This is particularly interesting as recent studies have shed light on the crosstalk  
404 between integrin and the microtubule cytoskeleton (LaFlamme et al. 2018). Finally, one family in  
405 module 48 (fam18955) is annotated as the DNA excision repair protein ERCC-3 in three Asgard  
406 genomes and three Theionarchaea genomes. The genes neighboring the genes of fam18955  
407 differ between the two lineages (**Figure 4B and Supplementary Dataset - Table S8**) and the  
408 three Asgard sequences only share between 20 and 23% protein identity with the three  
409 Theionarchaea sequences. These differences may indicate two distinct functions for this family.  
410 Fam18955 shows distant homology with the protein RAD25 of *Saccharomyces cerevisiae*.  
411 RAD25 is a DNA helicase required for DNA repair and RNA polymerase II transcription in *S.*  
412 *cerevisiae* (Guzder et al. 1994). RAD25 is also one of the six subunits of the transcription factor  
413 IIH (TFIIH) in *S. cerevisiae* (Sung et al. 1996). Consistent with the role of RAD25 in *S. cerevisiae*,  
414 the genes of family18955 is found next to replication factor C small subunit genes in the three  
415 Asgard genomes (**Figure 4 and Supplementary Dataset - Table S8**).

416  
417 **Groups without lineage-specific metabolic signatures**

418 The Archaeoglobi and Thermoplasmata lineages are unusual in that they have modules  
419 specific to them (modules 105 and 171 respectively), but no specific capacities were identified  
420 only in these groups based on functional predictions (**Supplementary Dataset - Table S3**). These  
421 lineage-specific modules have the highest percentage of hypothetical families of any lineage-  
422 specific module (**Table 1**).

423 Bathyarchaeota is the only lineage having more than 10 genomes and that does not have  
424 a specific module of families that is widespread in the 19 Bathyarchaeota genomes  
425 (**Supplementary Dataset - Table S2**). This is intriguing as Bathyarchaeota are widespread  
426 across terrestrial marine ecosystems and are known to thrive in diverse chemical environments

427 (Kubo et al. 2012). Alternatively, Bathyarchaeota may be a superphylum. In this case, the lack of  
428 modules shared across a superphylum would reinforce the results noted above.

429

### 430 **Hypothetical proteins distinguish the major archaeal groups**

431 Even after augmenting functional predictions using PDB and EggNog databases, families  
432 with functional predictions represent a tiny proportion of the protein families that comprise the  
433 lineage-specific modules (**Table 1, Figure 5 and Supplementary Dataset - Table S3**). 1053 out  
434 of 1411 hypothetical families remain unannotated (**Supplementary Dataset - Table S4**). 358  
435 hypothetical families have small domain matches, but not enough information is available to  
436 predict functions. For example, many have domains with matches to zinc finger domains, but such  
437 domains occur in proteins with diverse functions (**Supplementary Dataset - Table S4**). We found  
438 that the hypothetical proteins are shorter than proteins from the core families of module 1  
439 (**Supplementary Figure 14**) and are more likely to have a transmembrane helix prediction and a  
440 signal peptide prediction (**Table 1**).

441 Previous studies highlighted the presence of numerous families of proteins with roles in  
442 metabolism that are of bacterial origin but occur only in specific archaeal phyla (Nelson-Sathi et  
443 al. 2015, 2012). Consequently, we compared all archaeal protein families against a database of  
444 bacterial genomes sampled from across the bacterial tree of life to determine the extent to which  
445 proteins acquired from bacteria contribute to the archaeal group specific modules (**see materials**  
446 **and methods**). From 3% (Thermoplasmata) to 34% (Halobacteria) of the protein families in  
447 modules that are archaeal group specific have homologs in  $\geq 10$  distinct bacterial genomes, with  
448 the exception of Methanomicrobia, where 63% of the protein families have bacterial homologs  
449 (**Table 1**). Thus, for almost all archaeal groups, the majority of the protein families that form  
450 modules that separate them from other archaeal groups did not evolve in (or were not acquired  
451 from) bacteria. Further, we conclude genes acquired from bacteria only account for a small  
452 fraction of the lineage-specific families.

453

454

455

456

457

458

459

460

## 461 Discussion

462 We constructed a set of protein families for Domain Archaea, each of which generally  
463 corresponds with a set of homologous proteins with the same predicted function (in cases where  
464 functions could be assigned). Protein families with functional predictions that are specific to  
465 certain archaeal lineages (e.g., genes involved in methanogenesis or ammonia oxidation) well  
466 predict functional traits specific to these lineages. These observations indicate that the protein  
467 family construction method is robust. The generated set of 10,866 protein families is provided as  
468 an important community research resource. The patterns of presence/absence of protein families  
469 across genomes highlight sets of co-occurring proteins (modules), and groupings of genomes  
470 based on these modules mostly recapitulate archaeal phylogeny.

471 What is most striking from our analyses is the prominence of families of hypothetical  
472 proteins in the sets of highly prevalent lineage-specific proteins. An important consideration is  
473 whether (i) divergence of the sequence of these proteins from proteins with known function simply  
474 precluded functional annotation or (ii) whether they are novel proteins that serve well known  
475 functions, or if (iii) they represent functions that are unique and evolved following the divergence  
476 of each lineage from other archaea. Our analyses were designed to avoid case (i) by relying on  
477 state-of-the-art HMM-based homology detection methods that appear to well-group proteins with  
478 shared functions (**Supplementary Figure 3**). However, the fact that we could identify some  
479 probable functions using protein modeling suggests that (i) is correct in at least a subset of cases.  
480 For instance, PriX (fam03870) has structural homology with PriL but no sequence similarity was  
481 detected between PriX and any other protein in our analysis. Both proteins are distinct  
482 components of the primase complex in *Sulfolobus solfataricus* suggesting that PriX evolved from  
483 PriL by duplication followed by subfunctionalization (Holzer et al. 2017; Liu et al. 2015). Lineage-  
484 specific hypothetical proteins that are actually homologs of known proteins but currently too  
485 divergent for functional assignment are interesting, as they may have been under pressure to  
486 evolve more rapidly than normal during lineage divergence. It is not possible to distinguish (ii)  
487 from (iii) with the data available. In general, the sets of relatively common archaeal proteins  
488 without functional assignments provides targets for future biochemical studies.

489 Overall, the prevalence of transmembrane helices and signal peptides in the hypothetical  
490 proteins in lineage-specific modules indicates that they are membrane associated or extracellular,  
491 thus possibly involved in cell-cell and cell-environment interactions (some may be transporters).  
492 Where the lineages are confined to specific environments (e.g., halophiles), lineage-specific  
493 protein families may have evolved to meet the requirements of that environment (case (i) or (iii)).  
494 It is important to note that some modules contain many protein families and probably represent

495 combinations of new functions that, at the present time, cannot be resolved. Regardless of the  
496 explanation, or combination of explanations, for the presence of large numbers of lineage-specific  
497 proteins, the results clearly indicate the importance of divergence or evolution of a specific subset  
498 of proteins during emergence of the major archaeal lineages.

499       Possibly also informative regarding archaeal evolution is the observation that, despite  
500 resolving a Domain-wide core module (module 1), we detected only one case where a clearly  
501 defined module is conserved at the superphylum level. It is important to note that, with additional  
502 genomes, the two newly recognized Asgard phyla may be reclassified into a single phylum,  
503 eliminating this exception. The apparent lack of protein family module support for superphyla may  
504 argue against the phyletic gradualism, in which one lineage gradually transforms into another,  
505 and favor the theory of cladogenesis, where a lineage splits into two distinct lineages (Gould and  
506 Eldredge 1977). We acknowledge that modules containing fewer than 20 protein families (the  
507 cutoff used to define modules) may be uniquely associated with superphyla, and that some  
508 potentially important archaeal lineages were not included in the current analysis due to lack of a  
509 sufficient number of high-quality genomes.

510       The observation that the patterns of presence/absence of shared protein families groups  
511 together archaea that historically have been assigned to the same lineage and separates them  
512 from other lineages, in combination with innumerable prior publications on archaeal physiology  
513 and taxonomy (Adam et al. 2017; Spang, Caceres, and Ettema 2017; Baker et al. 2020), supports  
514 the value of the current taxonomic classifications within Domain Archaea. Overall, the results  
515 reinforce the concept that early archaeal evolution rapidly generated the major lineages, the rise  
516 of which was linked to acquisition of a set of proteins (recognized here as modules) that were  
517 largely retained during subsequent evolution of each lineage.

518



## 519 **Methods**

### 520 Genome collection

521           569 unpublished genomes were added to the 2,618 genomes of Archaea downloaded  
522 from the NCBI genome database in September 2018.

523           132 genomes were obtained from metagenomes of sediment samples. Sediment samples  
524 were collected from the Guaymas Basin (27°N0.388, 111°W24.560, Gulf of California, Mexico)  
525 during three cruises at a depth of approximately 2000 m below the water surface. Sediment cores  
526 were collected during two Alvin dives, 4486 and 4573 in 2008 and 2009. Sites referred to as  
527 “Megamat” (genomes starting with “Meg”) and “Aceto Balsamico” (genomes starting with “AB” in  
528 name), Core sections between 0-18 cm from 4486 and from 0-33 cm 4573 and were processed  
529 for these analyses. Intact sediment cores were subsampled under N<sub>2</sub> gas, and immediately frozen  
530 at -80 °C on board. The background of sampling sites was described previously (Teske et al.  
531 2016). Samples were processed for DNA isolation from using the MoBio PowerMax soil kit  
532 (Qiagen) following the manufacturer’s protocol. Illumina library preparation and sequencing were  
533 performed using Hiseq 4000 at Michigan State University. Paired-end reads were interleaved  
534 using `interleav_fasta.py` ([https://github.com/jorvis/biocode/blob/master/fasta/interleave\\_fasta.py](https://github.com/jorvis/biocode/blob/master/fasta/interleave_fasta.py))  
535 and the interleaved sequences were trimmed using `Sickle` (<https://github.com/najoshi/sickle>) with  
536 the default settings. Metagenomic reads from each subsample were individually assembled using  
537 IDBA-UD with the following parameters: `--pre_correction --mink 65 --maxk 115 --step 10 --`  
538 `seed_kmer 55` (Peng et al. 2012). Metagenomic binning was performed on contigs with a  
539 minimum length of 2000 bp in individual assemblies using the binning tools MetaBAT (Kang et al.  
540 2015) and CONCOCT (Alneberg et al. 2014), and resulting bins were combined with using DAS  
541 Tool (Sieber et al., n.d.). CheckM lineage\_wf (v1.0.5) (Parks et al. 2015) was used to estimate  
542 the percentage of completeness and contamination of bins. Genomes with more than 50%  
543 completeness and 10% contamination were manually optimized based on GC content, sequence  
544 depth and coverage using `mmgenome` (Karst, Kirkegaard, and Albertsen, n.d.).

545           The remaining 447 genomes came from previous sequencing and binning efforts  
546 (genomes starting with “ggkbase”). In brief, 168 genomes were obtained from an aquifer adjacent  
547 to the Colorado River near the town of Rifle, Colorado, USA in 2011 (Anantharaman et al. 2016),  
548 50 genomes from the Crystal Geyser system in Utah, USA (Probst et al. 2014). For DNA  
549 processing and sequencing methods see (Probst et al. 2017; Anantharaman et al. 2016). Forty-  
550 one genomes were obtained from (Parks et al. 2017). Additionally, 188 genomes were obtained  
551 from groundwater samples from Genasci Dairy Farm, located in Modesto, California (CA) as  
552 described in (He et al., n.d.).

553

554 Genome completeness assessment and de-replication.

555 Genome completeness and contamination were estimated based on the presence of single-copy  
556 genes (SCGs) as described in (Olm et al. 2017; Anantharaman et al. 2016). Genome  
557 completeness was estimated using 38 SCGs. For non-DPANN archaea, genomes with more than  
558 26 SCGs (>70% completeness) and less than 4 duplicated copies of the SCGs (<10%  
559 contamination) were considered as draft-quality genomes. Due to the reduced nature of DPANN  
560 genomes (Castelle et al. 2015), genomes with more than 22 SCGs and less than 4 duplicated  
561 copies of the SCGs were considered as draft-quality genomes. Genomes were de-replicated  
562 using dRep (Olm et al. 2017) (version v2.0.5 with ANI > 95%). The most complete and less  
563 contaminated genome per cluster was used in downstream analyses.

564

565 Concatenated 14 ribosomal proteins phylogeny

566 A maximum-likelihood tree was calculated based on the concatenation of 14 ribosomal proteins  
567 (L2, L3, L4, L5, L6, L14, L15, L18, L22, L24, S3, S8, S17, and S19). Homologous protein  
568 sequences were aligned using MAFFT (version 7.390) (--auto option) (Kato and Standley 2016),  
569 and alignments refined to remove gapped regions using Trimal (version 1.4.22) (--gappyout  
570 option) (Capella-Gutiérrez, Silla-Martínez, and Gabaldón 2009). The protein alignments were  
571 concatenated, with a final alignment of 1,179 genomes and 2,388 positions. The tree was inferred  
572 using RAXML (Stamatakis 2014) (version 8.2.10) (as implemented on the CIPRES web server  
573 (Miller, Pfeiffer, and Schwartz 2010)), under the LG plus gamma model of evolution, and with the  
574 number of bootstraps automatically determined via the MRE-based bootstopping criterion. A total  
575 of 108 bootstrap replicates were conducted, from which 100 were randomly sampled to determine  
576 support values.

577

578 Protein clustering

579 Protein clustering into families was achieved using a two-step procedure as previously described  
580 in (Méheust et al. 2019). A first protein clustering was done using the fast and sensitive protein  
581 sequence searching software MMseqs2 (version  
582 9f493f538d28b1412a2d124614e9d6ee27a55f45) (Steinegger and Söding 2017). An all-vs-all  
583 search was performed using e-value: 0.001, sensitivity: 7.5 and cover: 0.5. A sequence similarity  
584 network was built based on the pairwise similarities and the greedy set cover algorithm from  
585 MMseqs2 was performed to define protein subclusters. The resulting subclusters were defined  
586 as subfamilies. In order to test for distant homology, we grouped subfamilies into protein families

587 using an HMM-HMM comparison procedure as follows: the proteins of each subfamily with at  
588 least two protein members were aligned using the result2msa parameter of mmseqs2, and, from  
589 the multiple sequence alignments, HMM profiles were built using the HHpred suite (version 3.0.3)  
590 (Soding 2005). The subfamilies were then compared to each other using hhblits (Remmert et al.  
591 2011) from the HHpred suite (with parameters -v 0 -p 50 -z 4 -Z 32000 -B 0 -b 0). For subfamilies  
592 with probability scores of  $\geq 95\%$  and coverage  $\geq 0.50$ , a similarity score (probability  $\times$  coverage)  
593 was used as weights of the input network in the final clustering using the Markov Clustering  
594 algorithm (Enright, Van Dongen, and Ouzounis 2002), with 2.0 as the inflation parameter. These  
595 clusters were defined as the protein families.

596

#### 597 Module definition and taxonomic assignment

598 Looking at the distribution of the protein families across the genomes, a clear modular  
599 organization emerged. Modules of families were defined using a cutoff of 0.95 on the dendrogram  
600 tree of the families. The dendrogram tree was obtained from a hierarchical clustering using the  
601 Jaccard distance that was calculated based on profiles of protein family presence/absence. The  
602 corresponding clusters define the modules.

603 A phyla distribution was assigned to each module using the method of (Méheust et al.  
604 2019). For each module, the median number of genomes per family ( $m$ ) was calculated. The  
605 genomes were ranked by the number of families they carry. The  $m$  genomes that carry the most  
606 of families were retained; their phyla distribution defines the taxonomic assignment of the module.

607

#### 608 Functional annotation

609 Protein sequences were functionally annotated based on the accession of their best Hmsearch  
610 match (version 3.1) (E-value cut-off 0.001) (Eddy 1998) against an HMM database constructed  
611 based on ortholog groups defined by the KEGG (Kanehisa et al. 2016) (downloaded on June 10,  
612 2015). Domains were predicted using the same hmsearch procedure against the Pfam  
613 database (version 31.0) (Punta et al. 2012). The domain architecture of each protein sequence  
614 was predicted using the DAMA software (version 1.0) (default parameters) (Bernardes et al.  
615 2016). SIGNALP (version 5.0) (parameters: -format short -org arch) (Almagro Armenteros et al.  
616 2019) was used to predict the putative cellular localization of the proteins. Prediction of  
617 transmembrane helices in proteins was performed using TMHMM (version 2.0) (default  
618 parameters) (Krogh et al. 2001). Protein sequences were also functionally annotated based on  
619 the accession of their best hmsearch match (version 3.1) (E-value cut-off  $1e-10$ ) against the  
620 PDB database (Rose et al. 2017) (downloaded in February 2020).

621

622 HMM-HMM Predictions

623 Subfamilies were used to perform HMM-HMM annotation against arCogs of EggNog (version 5.0)  
624 (Huerta-Cepas et al. 2019) using hhblits (Remmert et al. 2011) from the HHpred suite (with  
625 parameters -v 0 -p 50 -z 4 -Z 32000 -B 0 -b 0). Subfamilies were subsequently functionally  
626 annotated based on the EggNog accessions of their best probability score.

627

628 Sequence similarity analysis

629 The 75,737 subfamilies from the 10,866 families were searched against a bacterial database of  
630 2,552 bacterial genomes (Méheust et al. 2019) using hmmsearch (version 3.1) (E-value cut-off  
631 0.001) (Eddy 1998). Among them 46,261 subfamilies, comprising 8,300 families, have at least  
632 one hit against a bacterial genome.

633

634

635 **Data availability**

636 The newly reconstructed genomes have been deposited at NCBI under BioProject **XX (TBA)**. The  
637 genomes of the herein analysed archaea have been made publicly available on the ggkbase  
638 database (**TBA**). Detailed annotations of the families are provided in the **Supplementary Dataset**  
639 - **Table S3** accompanying this paper. Raw data files (phylogenetic tree and fasta sequences of  
640 the families) are made available via figshare under the following link: **TBA**.

641

642 **Author contributions**

643 R.M., C.J.C. and J.F.B. designed the study. R.M. and C.J.C. created the dataset. C.J.C performed  
644 the phylogenetic analysis. A.L.J. performed the bacterial analysis. R.M. performed the protein  
645 family, the module detection, the genome annotations and HMMs analyses. R.M., C.J.C. and  
646 J.F.B. wrote the manuscript. All authors read and approved the final manuscript.

647

648 **Competing interests**

649 J.F.B. is a founder of Metagenomi. The authors declare that they have no conflict of interest.

650

651 **Materials and Correspondence**

652 Correspondence and material requests should be addressed to [jbanfield@berkeley.edu](mailto:jbanfield@berkeley.edu) and  
653 [raphael.meheust@berkeley.edu](mailto:raphael.meheust@berkeley.edu).

654

655 **Acknowledgements**

656 We thank Dr. Brett Baker, Dr Kiley Seitz and Dr Xianzhe Gong for the permission to use the  
657 metagenomic datasets from the Guaymas Basin in this study. We thank Dr. Christine He for the  
658 permission to use the metagenomic dataset from Genasci. We acknowledge funding support from  
659 the Chan Zuckerberg Biohub and the Innovative Genomics Institute at UC Berkeley.

660

## 661 References

- 662 Adam, Panagiotis S., Guillaume Borrel, Céline Brochier-Armanet, and Simonetta Gribaldo.  
663 2017. "The Growing Tree of Archaea: New Perspectives on Their Diversity, Evolution and  
664 Ecology." *The ISME Journal* 11 (11): 2407–25.
- 665 Akil, Caner, and Robert C. Robinson. 2018. "Genomes of Asgard Archaea Encode Profilins  
666 That Regulate Actin." *Nature* 562 (7727): 439–43.
- 667 Almagro Armenteros, José Juan, Konstantinos D. Tsirigos, Casper Kaae Sønderby, Thomas  
668 Nordahl Petersen, Ole Winther, Søren Brunak, Gunnar von Heijne, and Henrik Nielsen.  
669 2019. "SignalP 5.0 Improves Signal Peptide Predictions Using Deep Neural Networks."  
670 *Nature Biotechnology* 37 (4): 420–23.
- 671 Alneberg, Johannes, Brynjar Smári Bjarnason, Ino de Bruijn, Melanie Schirmer, Joshua Quick,  
672 Umer Z. Ijaz, Leo Lahti, Nicholas J. Loman, Anders F. Andersson, and Christopher Quince.  
673 2014. "Binning Metagenomic Contigs by Coverage and Composition." *Nature Methods* 11  
674 (11): 1144–46.
- 675 Altieri, Amanda S., Jane E. Ladner, Zhuo Li, Howard Robinson, Zahur F. Sallman, John P.  
676 Marino, and Zvi Kelman. 2016. "A Small Protein Inhibits Proliferating Cell Nuclear Antigen  
677 by Breaking the DNA Clamp." *Nucleic Acids Research* 44 (20): 10015.
- 678 Anantharaman, Karthik, Christopher T. Brown, Laura A. Hug, Itai Sharon, Cindy J. Castelle,  
679 Alexander J. Probst, Brian C. Thomas, et al. 2016. "Thousands of Microbial Genomes Shed  
680 Light on Interconnected Biogeochemical Processes in an Aquifer System." *Nature*  
681 *Communications* 7 (1): 13219.
- 682 Aouad, Monique, Najwa Taib, Anne Oudart, Michel Lecocq, Manolo Gouy, and Céline Brochier-  
683 Armanet. 2018. "Extreme Halophilic Archaea Derive from Two Distinct Methanogen Class II  
684 Lineages." *Molecular Phylogenetics and Evolution* 127 (October): 46–54.
- 685 Baker, Brett J., Valerie De Anda, Kiley W. Seitz, Nina Dombrowski, Alyson E. Santoro, and  
686 Karen G. Lloyd. 2020. "Diversity, Ecology and Evolution of Archaea." *Nature Microbiology* 5  
687 (7): 887–900.
- 688 Bartossek, Rita, Anja Spang, Gerhard Weidler, Anders Lanzen, and Christa Schleper. 2012.  
689 "Metagenomic Analysis of Ammonia-Oxidizing Archaea Affiliated with the Soil Group."  
690 *Frontiers in Microbiology* 3 (June): 208.
- 691 Bernardes, J. S., F. R. J. Vieira, G. Zaverucha, and A. Carbone. 2016. "A Multi-Objective  
692 Optimization Approach Accurately Resolves Protein Domain Architectures." *Bioinformatics*  
693 32 (3): 345–53.
- 694 Borrel, Guillaume, Nicolas Parisot, Hugh M. B. Harris, Eric Peyretailade, Nadia Gaci, William  
695 Tottey, Olivier Bardot, et al. 2014. "Comparative Genomics Highlights the Unique Biology of  
696 Methanomassiliicoccales, a Thermoplasmatales-Related Seventh Order of Methanogenic  
697 Archaea That Encodes Pyrrolysine." *BMC Genomics* 15 (August): 679.
- 698 Boyd, Joel A., Sean P. Jungbluth, Andy O. Leu, Paul N. Evans, Ben J. Woodcroft, Grayson L.  
699 Chadwick, Victoria J. Orphan, Jan P. Amend, Michael S. Rappé, and Gene W. Tyson.  
700 2019. "Divergent Methyl-Coenzyme M Reductase Genes in a Deep-Subseafloor  
701 Archaeoglobi." *The ISME Journal*. <https://doi.org/10.1038/s41396-018-0343-2>.
- 702 Brochier-Armanet, Céline, Bastien Boussau, Simonetta Gribaldo, and Patrick Forterre. 2008.  
703 "Mesophilic Crenarchaeota: Proposal for a Third Archaeal Phylum, the Thaumarchaeota."  
704 *Nature Reviews. Microbiology* 6 (3): 245–52.
- 705 Buckel, W. 2001. "Sodium Ion-Translocating Decarboxylases." *Biochimica et Biophysica Acta*  
706 1505 (1): 15–27.
- 707 Burke, S. A., S. L. Lo, and J. A. Krzycki. 1998. "Clustered Genes Encoding the  
708 Methyltransferases of Methanogenesis from Monomethylamine." *Journal of Bacteriology*  
709 180 (13): 3432–40.



- 710 Capella-Gutiérrez, Salvador, José M. Silla-Martínez, and Toni Gabaldón. 2009. “trimAl: A Tool  
711 for Automated Alignment Trimming in Large-Scale Phylogenetic Analyses.” *Bioinformatics*  
712 25 (15): 1972–73.
- 713 Castelle, Cindy J., Christopher T. Brown, Karthik Anantharaman, Alexander J. Probst, Raven H.  
714 Huang, and Jillian F. Banfield. 2018. “Biosynthetic Capacity, Metabolic Variety and Unusual  
715 Biology in the CPR and DPANN Radiations.” *Nature Reviews. Microbiology* 16 (10): 629–  
716 45.
- 717 Castelle, Cindy J., Kelly C. Wrighton, Brian C. Thomas, Laura A. Hug, Christopher T. Brown,  
718 Michael J. Wilkins, Kyle R. Frischkorn, et al. 2015. “Genomic Expansion of Domain  
719 Archaea Highlights Roles for Organisms from New Phyla in Anaerobic Carbon Cycling.”  
720 *Current Biology: CB* 25 (6): 690–701.
- 721 Cox, Cyron J., Peter G. Foster, Robert P. Hirt, Simon R. Harris, and T. Martin Embley. 2008.  
722 “The Archaeobacterial Origin of Eukaryotes.” *Proceedings of the National Academy of*  
723 *Sciences of the United States of America* 105 (51): 20356–61.
- 724 DasSarma, Shiladitya, and Priya Arora. 2006. “Genetic Analysis of the Gas Vesicle Gene  
725 Cluster in Haloarchaea.” *FEMS Microbiology Letters*. [https://doi.org/10.1111/j.1574-](https://doi.org/10.1111/j.1574-6968.1997.tb10456.x)  
726 [6968.1997.tb10456.x](https://doi.org/10.1111/j.1574-6968.1997.tb10456.x).
- 727 Dimroth, P. 1987. “Sodium Ion Transport Decarboxylases and Other Aspects of Sodium Ion  
728 Cycling in Bacteria.” *Microbiological Reviews* 51 (3): 320–40.
- 729 Dombrowski, Nina, Jun-Hoe Lee, Tom A. Williams, Pierre Offre, and Anja Spang. 2019.  
730 “Genomic Diversity, Lifestyles and Evolutionary Origins of DPANN Archaea.” *FEMS*  
731 *Microbiology Letters* 366 (2). <https://doi.org/10.1093/femsle/fnz008>.
- 732 Eddy, S. R. 1998. “Profile Hidden Markov Models.” *Bioinformatics* 14 (9): 755–63.
- 733 Enright, A. J., S. Van Dongen, and C. A. Ouzounis. 2002. “An Efficient Algorithm for Large-  
734 Scale Detection of Protein Families.” *Nucleic Acids Research* 30 (7): 1575–84.
- 735 Ermler, U., W. Grabarse, S. Shima, M. Goubeaud, and R. K. Thauer. 1997. “Crystal Structure of  
736 Methyl-Coenzyme M Reductase: The Key Enzyme of Biological Methane Formation.”  
737 *Science* 278 (5342): 1457–62.
- 738 Evans, Paul N., Joel A. Boyd, Andy O. Leu, Ben J. Woodcroft, Donovan H. Parks, Philip  
739 Hugenholtz, and Gene W. Tyson. 2019. “An Evolving View of Methane Metabolism in the  
740 Archaea.” *Nature Reviews. Microbiology* 17 (4): 219–32.
- 741 Evans, Paul N., Donovan H. Parks, Grayson L. Chadwick, Steven J. Robbins, Victoria J.  
742 Orphan, Suzanne D. Golding, and Gene W. Tyson. 2015. “Methane Metabolism in the  
743 Archaeal Phylum Bathyarchaeota Revealed by Genome-Centric Metagenomics.” *Science*  
744 350 (6259): 434–38.
- 745 Garrett, Roger A., Gisle Vestergaard, and Shiraz A. Shah. 2011. “Archaeal CRISPR-Based  
746 Immune Systems: Exchangeable Functional Modules.” *Trends in Microbiology* 19 (11):  
747 549–56.
- 748 Ghalei, Homa, Holger von Moeller, Detlef Eppers, Daniel Sohmen, Daniel N. Wilson, Bernhard  
749 Loll, and Markus C. Wahl. 2014. “Entrapment of DNA in an Intersubunit Tunnel System of a  
750 Single-Stranded DNA-Binding Protein.” *Nucleic Acids Research*.  
751 <https://doi.org/10.1093/nar/gku259>.
- 752 Gould, Stephen Jay, and Niles Eldredge. 1977. “Punctuated Equilibria: The Tempo and Mode of  
753 Evolution Reconsidered.” *Paleobiology*. <https://doi.org/10.1017/s0094837300005224>.
- 754 Guzder, Sami N., Patrick Sung, Véronique Bailly, Louise Prakash, and Satya Prakash. 1994.  
755 “RAD25 Is a DNA Helicase Required for DNA Repair and RNA Polymerase II  
756 Transcription.” *Nature*. <https://doi.org/10.1038/369578a0>.
- 757 Hartman, Hyman, and Alexei Fedorov. 2002. “The Origin of the Eukaryotic Cell: A Genomic  
758 Investigation.” *Proceedings of the National Academy of Sciences of the United States of*  
759 *America* 99 (3): 1420–25.
- 760 He, Christine, Ray Keren, Michael Whittaker, Ibrahim F. Farag, Jennifer Doudna, Jamie H. D.

- 761 Cate, and Jillian Banfield. n.d. "Huge and Variable Diversity of Episymbiotic CPR Bacteria  
762 and DPANN Archaea in Groundwater Ecosystems."  
763 <https://doi.org/10.1101/2020.05.14.094862>.
- 764 Holzer, Sandro, Jiangyu Yan, Mairi L. Kilkenny, Stephen D. Bell, and Luca Pellegrini. 2017.  
765 "Primer Synthesis by a Eukaryotic-like Archaeal Primase Is Independent of Its Fe-S  
766 Cluster." *Nature Communications* 8 (1): 1718.
- 767 Huang, Rui, Zev A. Ripstein, Rafal Augustyniak, Michal Lazniewski, Krzysztof Ginalski, Lewis E.  
768 Kay, and John L. Rubinstein. 2016. "Unfolding the Mechanism of the AAA+ Unfoldase VAT  
769 by a Combined Cryo-EM, Solution NMR Study." *Proceedings of the National Academy of  
770 Sciences of the United States of America* 113 (29): E4190–99.
- 771 Hua, Zheng-Shuang, Yan-Ni Qu, Qiyun Zhu, En-Min Zhou, Yan-Ling Qi, Yi-Rui Yin, Yang-Zhi  
772 Rao, et al. 2018. "Genomic Inference of the Metabolism and Evolution of the Archaeal  
773 Phylum Aigarchaeota." *Nature Communications*. [https://doi.org/10.1038/s41467-018-  
774 05284-4](https://doi.org/10.1038/s41467-018-05284-4).
- 775 Huerta-Cepas, Jaime, Damian Szklarczyk, Davide Heller, Ana Hernández-Plaza, Sofia K.  
776 Forslund, Helen Cook, Daniel R. Mende, et al. 2019. "eggNOG 5.0: A Hierarchical,  
777 Functionally and Phylogenetically Annotated Orthology Resource Based on 5090  
778 Organisms and 2502 Viruses." *Nucleic Acids Research* 47 (D1): D309–14.
- 779 Kanehisa, Minoru, Yoko Sato, Masayuki Kawashima, Miho Furumichi, and Mao Tanabe. 2016.  
780 "KEGG as a Reference Resource for Gene and Protein Annotation." *Nucleic Acids  
781 Research* 44 (D1): D457–62.
- 782 Kang, Dongwan D., Jeff Froula, Rob Egan, and Zhong Wang. 2015. "MetaBAT, an Efficient Tool  
783 for Accurately Reconstructing Single Genomes from Complex Microbial Communities."  
784 *PeerJ* 3 (August): e1165.
- 785 Karst, Søren M., Rasmus H. Kirkegaard, and Mads Albertsen. n.d. "Mmgenome: A Toolbox for  
786 Reproducible Genome Extraction from Metagenomes." <https://doi.org/10.1101/059121>.
- 787 Katoh, Kazutaka, and Daron M. Standley. 2016. "A Simple Method to Control over-Alignment in  
788 the MAFFT Multiple Sequence Alignment Program." *Bioinformatics* 32 (13): 1933–42.
- 789 Korkhin, Yakov, Ulug M. Unligil, Otis Littlefield, Pamlea J. Nelson, David I. Stuart, Paul B. Sigler,  
790 Stephen D. Bell, and Nicola G. A. Abrescia. 2009. "Evolution of Complex RNA  
791 Polymerases: The Complete Archaeal RNA Polymerase Structure." *PLoS Biology* 7 (5):  
792 e1000102.
- 793 Krogh, A., B. Larsson, G. von Heijne, and E. L. Sonnhammer. 2001. "Predicting  
794 Transmembrane Protein Topology with a Hidden Markov Model: Application to Complete  
795 Genomes." *Journal of Molecular Biology* 305 (3): 567–80.
- 796 Kubo, Kyoko, Karen G. Lloyd, Jennifer F. Biddle, Rudolf Amann, Andreas Teske, and Katrin  
797 Knittel. 2012. "Archaea of the Miscellaneous Crenarchaeotal Group Are Abundant, Diverse  
798 and Widespread in Marine Sediments." *The ISME Journal*.  
799 <https://doi.org/10.1038/ismej.2012.37>.
- 800 LaFlamme, Susan E., Shomita Mathew-Steiner, Neetu Singh, Diane Colello-Borges, and  
801 Bethsaida Nieves. 2018. "Integrin and Microtubule Crosstalk in the Regulation of Cellular  
802 Processes." *Cellular and Molecular Life Sciences: CMLS* 75 (22): 4177–85.
- 803 Lecompte, Odile, Raymond Ripp, Jean-Claude Thierry, Dino Moras, and Olivier Poch. 2002.  
804 "Comparative Analysis of Ribosomal Proteins in Complete Genomes: An Example of  
805 Reductive Evolution at the Domain Scale." *Nucleic Acids Research* 30 (24): 5382–90.
- 806 Lin, X., and J. Tang. 1990. "Purification, Characterization, and Gene Cloning of Thermopsin, a  
807 Thermostable Acid Protease from *Sulfolobus Acidocaldarius*." *The Journal of Biological  
808 Chemistry* 265 (3): 1490–95.
- 809 Liu, Bing, Songying Ouyang, Kira S. Makarova, Qiu Xia, Yanping Zhu, Zhimeng Li, Li Guo,  
810 Eugene V. Koonin, Zhi-Jie Liu, and Li Huang. 2015. "A Primase Subunit Essential for  
811 Efficient Primer Synthesis by an Archaeal Eukaryotic-Type Primase." *Nature*

- 812 *Communications* 6 (June): 7300.
- 813 McKay, Luke J., Mensur Dlakić, Matthew W. Fields, Tom O. Delmont, A. Murat Eren, Zackary J.
- 814 Jay, Korinne B. Klingensmith, Douglas B. Rusch, and William P. Inskeep. 2019. “Co-
- 815 Occurring Genomic Capacity for Anaerobic Methane and Dissimilatory Sulfur Metabolisms
- 816 Discovered in the Korarchaeota.” *Nature Microbiology* 4 (4): 614–22.
- 817 Méheust, Raphaël, David Burstein, Cindy J. Castelle, and Jillian F. Banfield. 2019. “The
- 818 Distinction of CPR Bacteria from Other Bacteria Based on Protein Family Content.” *Nature*
- 819 *Communications* 10 (1): 4173.
- 820 Michiels, Jan, Chuanwu Xi, Jan Verhaert, and Jos Vanderleyden. 2002. “The Functions of
- 821 Ca(2+) in Bacteria: A Role for EF-Hand Proteins?” *Trends in Microbiology* 10 (2): 87–93.
- 822 Miller, M. A., W. Pfeiffer, and T. Schwartz. 2010. “Creating the CIPRES Science Gateway for
- 823 Inference of Large Phylogenetic Trees.” In *2010 Gateway Computing Environments*
- 824 *Workshop (GCE)*, 1–8.
- 825 Nelson-Sathi, Shijulal, Tal Dagan, Giddy Landan, Arnold Janssen, Mike Steel, James O.
- 826 McInerney, Uwe Deppenmeier, and William F. Martin. 2012. “Acquisition of 1,000
- 827 Eubacterial Genes Physiologically Transformed a Methanogen at the Origin of
- 828 Haloarchaea.” *Proceedings of the National Academy of Sciences of the United States of*
- 829 *America* 109 (50): 20537–42.
- 830 Nelson-Sathi, Shijulal, Filipa L. Sousa, Mayo Roettger, Nabor Lozada-Chávez, Thorsten
- 831 Thiergart, Arnold Janssen, David Bryant, et al. 2015. “Origins of Major Archaeal Clades
- 832 Correspond to Gene Acquisitions from Bacteria.” *Nature* 517 (7532): 77–80.
- 833 Niehaus, Thomas D., Mona Elbadawi-Sidhu, Valérie de Crécy-Lagard, Oliver Fiehn, and
- 834 Andrew D. Hanson. 2017. “Discovery of a Widespread Prokaryotic 5-Oxoprolinase That
- 835 Was Hiding in Plain Sight.” *The Journal of Biological Chemistry* 292 (39): 16360–67.
- 836 Offre, Pierre, Anja Spang, and Christa Schleper. 2013. “Archaea in Biogeochemical Cycles.”
- 837 *Annual Review of Microbiology* 67 (June): 437–57.
- 838 Olm, Matthew R., Christopher T. Brown, Brandon Brooks, and Jillian F. Banfield. 2017. “dRep: A
- 839 Tool for Fast and Accurate Genomic Comparisons That Enables Improved Genome
- 840 Recovery from Metagenomes through de-Replication.” *The ISME Journal* 11 (12): 2864–68.
- 841 Parks, Donovan H., Michael Imelfort, Connor T. Skennerton, Philip Hugenholtz, and Gene W.
- 842 Tyson. 2015. “CheckM: Assessing the Quality of Microbial Genomes Recovered from
- 843 Isolates, Single Cells, and Metagenomes.” *Genome Research* 25 (7): 1043–55.
- 844 Parks, Donovan H., Christian Rinke, Maria Chuvochina, Pierre-Alain Chaumeil, Ben J.
- 845 Woodcroft, Paul N. Evans, Philip Hugenholtz, and Gene W. Tyson. 2017. “Recovery of
- 846 Nearly 8,000 Metagenome-Assembled Genomes Substantially Expands the Tree of Life.”
- 847 *Nature Microbiology* 2 (11): 1533–42.
- 848 Paytubi, Sonia, Stephen A. McMahon, Shirley Graham, Huanting Liu, Catherine H. Botting, Kira
- 849 S. Makarova, Eugene V. Koonin, James H. Naismith, and Malcolm F. White. 2012.
- 850 “Displacement of the Canonical Single-Stranded DNA-Binding Protein in the
- 851 Thermoproteales.” *Proceedings of the National Academy of Sciences of the United States*
- 852 *of America* 109 (7): E398–405.
- 853 Peng, Y., H. C. M. Leung, S. M. Yiu, and F. Y. L. Chin. 2012. “IDBA-UD: A de Novo Assembler
- 854 for Single-Cell and Metagenomic Sequencing Data with Highly Uneven Depth.”
- 855 *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/bts174>.
- 856 Pester, Michael, Christa Schleper, and Michael Wagner. 2011. “The Thaumarchaeota: An
- 857 Emerging View of Their Phylogeny and Ecophysiology.” *Current Opinion in Microbiology*.
- 858 <https://doi.org/10.1016/j.mib.2011.04.007>.
- 859 Petitjean, Céline, Philippe Deschamps, Purificación López-García, and David Moreira. 2014.
- 860 “Rooting the Domain Archaea by Phylogenomic Analysis Supports the Foundation of the
- 861 New Kingdom Proteoarchaeota.” *Genome Biology and Evolution* 7 (1): 191–204.
- 862 Probst, Alexander J., Cindy J. Castelle, Andrea Singh, Christopher T. Brown, Karthik

- 863 Anantharaman, Itai Sharon, Laura A. Hug, et al. 2017. "Genomic Resolution of a Cold  
864 Subsurface Aquifer Community Provides Metabolic Insights for Novel Microbes Adapted to  
865 High CO<sub>2</sub> concentrations." *Environmental Microbiology*. <https://doi.org/10.1111/1462-2920.13362>.
- 867 Probst, Alexander J., Thomas Weinmaier, Kasie Raymann, Alexandra Perras, Joanne B.  
868 Emerson, Thomas Rattei, Gerhard Wanner, et al. 2014. "Biology of a Widespread  
869 Uncultivated Archaeon That Contributes to Carbon Fixation in the Subsurface." *Nature*  
870 *Communications* 5 (November): 5497.
- 871 Punta, Marco, Penny C. Coggill, Ruth Y. Eberhardt, Jaina Mistry, John Tate, Chris Boursnell,  
872 Ningze Pang, et al. 2012. "The Pfam Protein Families Database." *Nucleic Acids Research*  
873 40 (Database issue): D290–301.
- 874 Remmert, Michael, Andreas Biegert, Andreas Hauser, and Johannes Söding. 2011. "HHblits:  
875 Lightning-Fast Iterative Protein Sequence Searching by HMM-HMM Alignment." *Nature*  
876 *Methods* 9 (2): 173–75.
- 877 Rinke, Christian, Patrick Schwientek, Alexander Sczyrba, Natalia N. Ivanova, Iain J. Anderson,  
878 Jan-Fang Cheng, Aaron Darling, et al. 2013. "Insights into the Phylogeny and Coding  
879 Potential of Microbial Dark Matter." *Nature* 499 (7459): 431–37.
- 880 Rose, Peter W., Andreas Prlić, Ali Altunkaya, Chunxiao Bi, Anthony R. Bradley, Cole H.  
881 Christie, Luigi Di Costanzo, et al. 2017. "The RCSB Protein Data Bank: Integrative View of  
882 Protein, Gene and 3D Structural Information." *Nucleic Acids Research* 45 (D1): D271–81.
- 883 Sapa, Rajat, Karine Bagramyan, and Michael W. W. Adams. 2003. "A Simple Energy-  
884 Conserving System: Proton Reduction Coupled to Proton Translocation." *Proceedings of*  
885 *the National Academy of Sciences of the United States of America* 100 (13): 7545–50.
- 886 Schut, Gerrit J., Gina L. Lipscomb, Diep M. N. Nguyen, Robert M. Kelly, and Michael W. W.  
887 Adams. 2016. "Heterologous Production of an Energy-Conserving Carbon Monoxide  
888 Dehydrogenase Complex in the Hyperthermophile *Pyrococcus furiosus*." *Frontiers in*  
889 *Microbiology* 7 (January): 29.
- 890 Schut, Gerrit J., William J. Nixon, Gina L. Lipscomb, Robert A. Scott, and Michael W. W.  
891 Adams. 2012. "Mutational Analyses of the Enzymes Involved in the Metabolism of  
892 Hydrogen by the Hyperthermophilic Archaeon *Pyrococcus furiosus*." *Frontiers in*  
893 *Microbiology*. <https://doi.org/10.3389/fmicb.2012.00163>.
- 894 Sieber, Christian M. K., Alexander J. Probst, Allison Sharrar, Brian C. Thomas, Matthias Hess,  
895 Susannah G. Tringe, and Jillian F. Banfield. n.d. "Recovery of Genomes from  
896 Metagenomes via a Dereplication, Aggregation, and Scoring Strategy."  
897 <https://doi.org/10.1101/107789>.
- 898 Snel, B., P. Bork, and M. A. Huynen. 1999. "Genome Phylogeny Based on Gene Content."  
899 *Nature Genetics* 21 (1): 108–10.
- 900 Soding, J. 2005. "Protein Homology Detection by HMM-HMM Comparison." *Bioinformatics*.  
901 <https://doi.org/10.1093/bioinformatics/bti125>.
- 902 Spang, Anja, Eva F. Caceres, and Thijs J. G. Ettema. 2017. "Genomic Exploration of the  
903 Diversity, Ecology, and Evolution of the Archaeal Domain of Life." *Science* 357 (6351).  
904 <https://doi.org/10.1126/science.aaf3883>.
- 905 Spang, Anja, Jimmy H. Saw, Steffen L. Jørgensen, Katarzyna Zaremba-Niedzwiedzka, Joran  
906 Martijn, Anders E. Lind, Roel van Eijk, Christa Schleper, Lionel Guy, and Thijs J. G. Ettema.  
907 2015. "Complex Archaea That Bridge the Gap between Prokaryotes and Eukaryotes."  
908 *Nature* 521 (7551): 173–79.
- 909 Stamatakis, Alexandros. 2014. "RAxML Version 8: A Tool for Phylogenetic Analysis and Post-  
910 Analysis of Large Phylogenies." *Bioinformatics* 30 (9): 1312–13.
- 911 Steinegger, Martin, and Johannes Söding. 2017. "MMseqs2 Enables Sensitive Protein  
912 Sequence Searching for the Analysis of Massive Data Sets." *Nature Biotechnology* 35 (11):  
913 1026–28.



- 914 Sung, P., S. N. Guzder, L. Prakash, and S. Prakash. 1996. "Reconstitution of TFIIH and  
915 Requirement of Its DNA Helicase Subunits, Rad3 and Rad25, in the Incision Step of  
916 Nucleotide Excision Repair." *The Journal of Biological Chemistry* 271 (18): 10821–26.
- 917 Teske, Andreas, Dirk de Beer, Luke J. McKay, Margaret K. Tivey, Jennifer F. Biddle, Daniel  
918 Hoer, Karen G. Lloyd, et al. 2016. "The Guaymas Basin Hiking Guide to Hydrothermal  
919 Mounds, Chimneys, and Microbial Mats: Complex Seafloor Expressions of Subsurface  
920 Hydrothermal Circulation." *Frontiers in Microbiology* 7 (February): 75.
- 921 Tully, Benjamin J. 2019. "Metabolic Diversity within the Globally Abundant Marine Group II  
922 Euryarchaea Offers Insight into Ecological Patterns." *Nature Communications* 10 (1): 271.
- 923 Vestergaard, Gisle, Roger A. Garrett, and Shiraz A. Shah. 2014. "CRISPR Adaptive Immune  
924 Systems of Archaea." *RNA Biology*. <https://doi.org/10.4161/rna.27990>.
- 925 Wang, Yinzhao, Gunter Wegener, Jialin Hou, Fengping Wang, and Xiang Xiao. 2019.  
926 "Expanding Anaerobic Alkane Metabolism in the Domain of Archaea." *Nature Microbiology*  
927 4 (4): 595–602.
- 928 Woese, C. R., and G. E. Fox. 1977. "Phylogenetic Structure of the Prokaryotic Domain: The  
929 Primary Kingdoms." *Proceedings of the National Academy of Sciences of the United States*  
930 *of America* 74 (11): 5088–90.
- 931 Woese, C. R., R. Gupta, C. M. Hahn, W. Zillig, and J. Tu. 1984. "The Phylogenetic  
932 Relationships of Three Sulfur Dependent Archaeobacteria." *Systematic and Applied*  
933 *Microbiology* 5: 97–105.
- 934 Woese, C. R., O. Kandler, and M. L. Wheelis. 1990. "Towards a Natural System of Organisms:  
935 Proposal for the Domains Archaea, Bacteria, and Eucarya." *Proceedings of the National*  
936 *Academy of Sciences of the United States of America* 87 (12): 4576–79.
- 937 Yang, Ying, Rie Yatsunami, Ai Ando, Nobuhiro Miyoko, Toshiaki Fukui, Shinichi Takaichi, and  
938 Satoshi Nakamura. 2015. "Complete Biosynthetic Pathway of the C50 Carotenoid  
939 Bacterioruberin from Lycopene in the Extremely Halophilic Archaeon Haloarcula Japonica."  
940 *Journal of Bacteriology* 197 (9): 1614–23.
- 941 Yan, Jianguyu, Thomas R. Beattie, Adriana L. Rojas, Kelly Schermerhorn, Tamzin Gristwood,  
942 Jonathan C. Trinidad, Sonja V. Albers, et al. 2017. "Identification and Characterization of a  
943 Heterotrimeric Archaeal DNA Polymerase Holoenzyme." *Nature Communications* 8 (May):  
944 15075.
- 945 Yu, Hongjun, Chang-Hao Wu, Gerrit J. Schut, Dominik K. Haja, Gongpu Zhao, John W. Peters,  
946 Michael W. W. Adams, and Huilin Li. 2018. "Structure of an Ancient Respiratory System."  
947 *Cell* 173 (7): 1636–49.e16.
- 948 Yutin, Natalya, Pere Puigbò, Eugene V. Koonin, and Yuri I. Wolf. 2012. "Phylogenomics of  
949 Prokaryotic Ribosomal Proteins." *PLoS ONE*. <https://doi.org/10.1371/journal.pone.0036972>.
- 950 Zaremba-Niedzwiedzka, Katarzyna, Eva F. Caceres, Jimmy H. Saw, Disa Bäckström, Lina  
951 Juzokaite, Emmelien Vancaester, Kiley W. Seitz, et al. 2017. "Asgard Archaea Illuminate  
952 the Origin of Eukaryotic Cellular Complexity." *Nature* 541 (7637): 353–58.
- 953 Zhang, P., J. Wang, and Y. Shi. 2010. "Structure and Mechanism of the S Component of a  
954 Bacterial ECF Transporter." <https://doi.org/10.2210/pdb3p5n/pdb>.

955