1  **Ancestral absence of electron transport chains in Patescibacteria and DPANN**

2

3  Jacob P. Beam[1], Eric D. Becraft[1,15], Julia M. Brown[1], Frederik Schulz[2], Jessica K. Jarett[2,16],

4  Oliver Bezuidt[1], Nicole J. Poulton[1], Kayla Clark[1,17], Peter F. Dunfield[3], Nikolai V. Ravin[4], John

5  R. Spear[5], Brian P. Hedlund[6], Konstantinos A. Kormas[7], Stefan M. Sievert[8], Mostafa S.

6  Elshahed[9], Hazel A. Barton[10], Matthew B. Stott[11], Jonathan A. Eisen[12], Duane P. Moser[13], Tullis

7  C. Onstott[14], Tanja Woyke[2], and Ramunas Stepanauskas[1*]

8

9  [1]Bigelow Laboratory for Ocean Sciences, East Boothbay, Maine, USA, 04544

10  [2]Department of Energy Joint Genome Institute, Walnut Creek, California, USA, 94598

11  [3]University of Calgary, Calgary, AB, Canada, T2N 1N4

12  [4]Institute of Bioengineering, Research Center of Biotechnology of the Russian Academy of

13  Sciences, Moscow, 119071, Russia

14  [5]Civil and Environmental Engineering, Colorado School of Mines, Golden, Colorado 80401

15  [6]School of Life Sciences, University of Nevada Las Vegas, Las Vegas, NV, 89154, USA.

16  Nevada Institute of Personalized Medicine, University of Nevada Las Vegas, Las Vegas, NV,

17  USA, 89154

18  [7]Department of Ichthyology and Aquatic Environment, University of Thessaly, 38446, Volos,

19  Greece

20  [8]Biology Department, Woods Hole Oceanographic Institution, Woods Hole, MA, USA, 02543

21  [9]Department of Microbiology and Molecular Genetics, Oklahoma State University, Stillwater,

22  Oklahoma, USA

23  [10]Department of Biology, University of Akron, Akron, OH 44325

24    [11]School of Biological Sciences, University of Canterbury, Christchurch, 8041 New Zealand.

25    [12]University of California Davis, Davis, CA, USA, 95616

26    [13]Desert Research Institute, Las Vegas, Nevada, USA, 89119

27    [14]Princeton University, Princeton, New Jersey, USA, 08544

28    Current address:

29    [15]University of North Alabama, Florence, AL, USA, 35632

30    [16]AnimalBiome, Oakland, CA, USA, 94609

31    [17]U.S. Army Engineer Research and Development Center, US Army Corp of Engineers,

32    Vicksburg, Mississippi, 39180

33

34    *To whom correspondence should be addressed: rstepanauskas@bigelow.org

35

36

37

38

39

40

41

42

43

44

45

46

47  **Abstract**

48      Recent discoveries suggest that the candidate superphyla Patescibacteria and DPANN

49  constitute a large fraction of the phylogenetic diversity of Bacteria and Archaea. Their small

50  genomes and limited coding potential have been hypothesized to be ancestral adaptations to

51  obligate symbiotic lifestyles. To test this hypothesis, we performed cell-cell association,

52  genomic, and phylogenetic analyses on 4,829 individual cells of Bacteria and Archaea from 46

53  globally distributed surface and subsurface field samples. This confirmed the ubiquity and

54  abundance of Patescibacteria and DPANN in subsurface environments, the small size of their

55  genomes and cells, and the divergence of their gene content from other Bacteria and Archaea.

56  Our analyses suggest that most Patescibacteria and DPANN in the studied subsurface

57  environments do not form specific physical associations with other microorganisms. These data

58  also suggest that their unusual genomic features and prevalent auxotrophies may be a result of

59  minimal cellular energy transduction mechanisms that potentially precede the evolution of

60  respiration, thus relying solely on fermentation for energy conservation.

61

62

63

66

67

68

69

70 **Introduction**

71        Cultivation-independent research tools have revealed the coding potential of numerous,

72    deep branches of Bacteria and Archaea that were unknown until recently (Wrighton et al., 2012;

73    Rinke et al., 2013; Brown et al., 2015; Castelle et al., 2015) Among them, the candidate bacterial

74    superphylum Patescibacteria (also known as Candidate Phyla Radiation, CPR) and archaeal

75    superphylum DPANN have garnered particular attention, as they appear to constitute a large

76    fraction of microbial diversity in subsurface and various other environments (Brown et al., 2015;

77    Hug et al., 2016; Dombrowski et al., 2019). Patescibacteria and DPANN are characterized by

78    small genomes and cell sizes, and predicted minimal biosynthetic and metabolic potential

79    (Wrighton et al., 2012; Luef et al., 2015; Castelle and Banfield, 2018). They also appear to have

80    slow metabolism, as indicated by low per-cell ribosome counts (Luef et al., 2015) and slow

81    estimated genome replication rates (Brown et al., 2016). Host-dependent endo- and ecto-

82    symbioses have been observed in several Patescibacteria (Gong et al., 2014; He et al., 2015;

83    Cross et al., 2019) and the Nanoarchaeota and Nanohaloarchaeota phyla within DPANN (Huber

84    et al., 2002; Podar et al., 2013; Munson-McGee et al., 2015; Jarett et al., 2018; Hamm et al.,

85    2019). As a result, it has been posited that the unusual biological features of Patescibacteria and

86    DPANN reflect ancestral adaptations to symbiotic lifestyles (Castelle et al., 2018; Dombrowski

87    et al., 2019). However, direct evidence of symbiosis in Patescibacteria and DPANN is limited to

88    a small number of narrow phylogenetic groups inhabiting surface environments and, in the case

89    of Patescibacteria, dependent on eukaryotic hosts (Gong et al., 2014) or eukaryotic host systems

90    (He et al., 2015; Cross et al., 2019) (i.e., mammalian oral cavities), which suggests relatively

91    recent adaptations.

4

92      Here, we performed physical cell-cell association, genomic, and phylogenetic analyses on

93      4,829 of individual microbial cells from 46 globally distributed and environmentally diverse

94      locations to gain additional insights into the unusual biological features of Patescibacteria and

95      DPANN. Consistent with prior reports, we found these two superphyla abundant in many

96      subsurface environments, and also confirm their consistently small cell and genome sizes. Our

97      single cell genomic and biophysical observations do not support the prevailing view that

98      Patescibacteria and DPANN are dominated by symbionts (Castelle et al., 2018; Dombrowski et

99      al., 2019). Instead, based on the apparent lack of genes for complete electron transport systems,

100     we hypothesize that these two superphyla have not evolved the capacity for respiration and

101     therefore rely on fermentative metabolisms for energy conservation. Although complex

102     metabolic interdependencies are a rule rather than exception in natural microbiomes (Zengler and

103     Zaramela, 2018), the predicted fermentative energy conservation and limited biosynthetic

104     potential (Castelle et al., 2018; Dombrowski et al., 2019) of Patescibacteria and DPANN may

105     define a highly communal lifestyle of these two superphyla and provide explanation for the

106     extreme difficulty in obtaining them in pure culture.

107

108     **Materials and Methods**

109

110     Field sample collection

111     Field samples were collected from a global set of diverse environments that were found

112     to contain candidate phyla of Bacteria and Archaea in prior studies (Rinke et al., 2013; Thomas

113     et al., 2013; Moser et al., 2015; Becraft et al., 2017; Hershey et al., 2018; Sackett, 2018; Sackett

114     et al., 2018, 2019). Immediately after collection, samples were amended with sterile 5% glycerol

115 and 1 mM EDTA (final concentrations) and stored at -80 °C. Field sample metadata is located

116 with each individual SAG in Table S1.

117

118 <u>Single amplified genome (SAG) generation, sequencing, and *de novo* assembly</u>

119 SAG generation and sequencing were performed by Bigelow Laboratory for Ocean Sciences

120 Single Cell Genomics Center (SCGC) and U.S. Department of Energy Joint Genome Institute

121 (JGI) (Table S1). At SCGC, field samples were stained with SYTO-9 nucleic acids stain

122 (Thermo Fisher Scientific), separated using fluorescence-activated cell sorting (FACS), lysed

123 using a combination of freeze-thaw and alkaline treatment, and their genomic DNA was

124 amplified using WGA-X in a cleanroom, as previously described (Stepanauskas et al., 2017). For

125 sorting of cells with active oxidoreductases, the Beatrix field sample (plate AG-274) was pre-

126 incubated with the RedoxSensor Green stain (Thermo Fisher Scientific) following

127 manufacturer's instructions. During cell sorting, cell size estimates were performed using

128 calibrated index FACS (Stepanauskas et al., 2017). All SAGs generated at SCGC were subject to

129 Low Coverage Sequencing (LoCoS) using a modified Nextera library preparation protocol and

130 NextSeq 500 (Illumina) sequencing instrumentation (Stepanauskas et al., 2017). This resulted in

131 a variable number of 2x150bp reads per SAG, with an average of ~300k. The reads were de novo

132 assembled using a customized workflow utilizing SPAdes (Bankevich et al., 2012), as previously

133 described (Stepanauskas et al., 2017). The quality of the sequencing reads was assessed using

134 FastQC and the quality of the assembled genomes (contamination and completeness) was

135 assessed using checkM (Parks et al., 2015) and tetramer frequency analysis (Woyke et al., 2009).

136 This SAG generation, sequencing and assembly workflow was previously evaluated for

137 assembly errors using three bacterial benchmark cultures with diverse genome complexity and

138    GC content (%), indicating no non-target and undefined bases in the assemblies and average

139    frequencies of mis-assemblies, indels and mismatches per 100 kbp being 1.5, 3.0 and 5.0

140    (Stepanauskas et al., 2017). Functional annotation was first performed using Prokka (Seemann,

141    2014) with default Swiss-Prot databases supplied by the software. Prokka was run a second time

142    with a custom protein annotation database built from compiling Swiss-Prot (Bateman et al.,

143    2017) entries for Archaea and Bacteria. The uniquely barcoded sequencing libraries of SAGs

144    belonging to candidate divisions were combined, in equal proportions, into 48-library pools and

145    shipped to JGI for deeper sequencing with NextSeq 500 (Illumina) in 2x150 bp mode. Quality

146    filtering of raw reads was performed with BBTools v.37, read normalization with BBNorm, and

147    error correction with Tadpole (http://bbtools.jgi.doe.gov). The resulting reads were assembled

148    with SPAdes (Nurk et al., 2013) (v3.9.0, --phred-offset 33 –sc -k 22,55,95 –12), and 200 bp was

149    trimmed from the ends of assembled contigs, after which contigs with read coverage < 2 or < 2

150    kbp in length were discarded. Assemblies were annotated according to IMG standard protocols

151    (Huntemann et al., 2016; Chen et al., 2019). All SAGs are publicly available in IMG/M (Chen et

152    al., 2019), and can be found under their GOLD analysis project identifiers in Table S1.

153

154    Identification of heterogenous DNA sources

155        The 16S ribosomal RNA gene was identified in SAGs by searching them individually

156    using cmsearch, which is part of the infernal package (Nawrocki and Eddy, 2013), using the

157    bacterial 16S rRNA Rfam covariance model (rfam.xfam.org/family/RF00177). This method is

158    particularly helpful in predicting 16S rRNA genes in Patescibacteria and DPANN, which can

159    often have introns in their 16S rRNA genes (Brown, 2015). Taxonomic assignments to these 16S

160    rRNA genes were conducted using "classify.seqs" within mothur (Schloss et al., 2009) version

7

161    1.41.3 against the Silva 132 reference database and taxonomy file (Quast et al., 2013). The

162    resulting taxonomy file was used to search for SAGs that contained two 16S rRNA genes that

163    had different taxonomic phylum-level assignments and were marked as putative co-sorts; those

164    that did not have two 16S rRNA genes were marked as single sorts. The checkM (Parks et al.,

165    2015) contamination estimates were used to determine SAGs that had high values of potential

166    genome admixture (e.g., two different cellular origins). A Chi-squared test was performed in R

167    using the "chisq.test" function on potential co-sorted and single sorted SAGs, and Pearson's

168    residuals were retrieved from the output of this test and used to calculate the percent contribution

169    to each $X^2$ statistic, and plotted using the "corrplot" package in R Studio.

170

171    Genomes from prior studies

172        A total of 1,025 publicly available SAGs, metagenome bins, and isolate genomes (Table

173    S2) were used in this study from the Integrated Microbial Genomes and Microbiomes (IMG/M)

174    database (Chen et al., 2019) (genomes accessed April 2018). These genomes were selected by

175    clustering the RNA polymerase COG0086 protein sequence at 70% identity, and if there were

176    similar genomes at the 70% identity threshold, the one with the most complete set of 56 single

177    copy proteins was chosen as representative. Phylum-level classification and symbiotic lifestyle

178    assignments were exported from IMG/M. In cases were IMG/M lacked lifestyle assignments,

179    manual literature searches of organism names were used to determine whether they have

180    documented symbiotic relationships.

181

182

183

8

184 Concatenated single copy protein phylogeny

185    A set of 56 universal single copy marker proteins (Eloe-Fadrosh et al., 2016; Yu et al.,

186    2017) was used to build a phylogenetic tree for the newly generated SAGs and MAGs and a

187    representative set of bacteria and archaea based on publicly available microbial genomes in

188    IMG/M (Chen et al., 2019) (genomes accessed in April 2018). Marker proteins were identified

189    with hmmsearch (Eddy, 2011) version 3.1b2, using a specific Hidden Markov Model for each of

190    the markers. Genomes for which 5 or more different marker proteins could be identified were

191    included in the tree. For every marker protein, alignments were built with MAFFT (Nakamura et

192    al., 2018) v7.294b and subsequently trimmed with BMGE (Criscuolo and Gribaldo, 2010) v1.12

193    using BLOSUM30. Single protein alignments were concatenated and maximum likelihood

194    phylogenies inferred with FastTree2 (Price et al., 2010) using the options: -spr 4 -mlacc 2 -

195    slownni -lg (for archaea) and -spr 4 -mlacc 2 -slownni -lg (for bacteria).

196

197 Clusters of orthologous groups principal components analysis

198    Clusters of orthologous groups (COGs) were assigned to SAG (Table S1) and reference

199    genome (Table S2) predicted protein sequences using reverse position-specific blast (rpsblast)

200    (Altschul et al., 1997) with an e-value cutoff of 1e-5 and the cdd2cog script

201    (https://github.com/aleimba/bac-genomics-scripts/tree/master/cdd2cog). Genomes that were used

202    for the principal component analysis (PCA) had completeness estimates greater than or equal to

203    30%, and contained 16S rRNA genes for unambiguous phylum-level classification. Eigenvector

204    values were calculated in RStudio (RStudio Team, 2016) version 1.1.463 using the cmdscale

205    function from relative abundances of the different COG categories expressed as a percent out of

206    the total number of assigned COGs. PCA plots were visualized with ggplot2 (Wickham, 2016) in

207   RStudio (RStudio Team, 2016). A Wilcoxon test was performed in RStudio using the

208   "wilcox.test" function to determine statistical differences between principal components among

209   the different clusters discussed in the main text. The color scheme for these plots is based on the

210   Color Universal Design (https://jfly.uni-koeln.de/color/), and should be distinguishable by all

211   types of vision. This color scheme was used throughout all the figures in the manuscript.

212

213   Coding sequence density

214         Coding sequences (CDS) for SAGs and reference genomes were predicted using Prodigal

215   (Hyatt et al., 2010) version 2.6.3. The initial analysis of prokka CDS density revealed that

216   numerous SAGs and reference genomes had very low coding densities. Prokka utilizes the code

217   11 translation table by default, and many of these genomes could potentially use stop codons in

218   place of canonical codons (Wrighton et al., 2012; Rinke et al., 2013). We determined the correct

219   translation table to utilize for each genome by comparing the total CDS length from Code 11 and

220   Code 25 predictions, and if the Code 11 total CDS length was greater than the Code 25 total

221   CDS length, then the total length from Code 11 was used in the coding density calculation. If the

222   opposite was true, then the Code 25 total CDS length was used. The coding density was

223   calculated by dividing the total CDS sequence by the total assembly size.

224

225   Oxygen reductase identification

226         A published heme copper oxidase subunit I database (Sousa et al., 2011) from bacteria

227   and archaea was used as a database with blastp (Altschul et al., 1990) with an e-value cutoff of

228   1e-10 using the SAG and reference genomes as queries. The original database file had to be de-

229   replicated (i.e., removing 100% identical sequences) using the dedupe.sh script, which is part of

10

230     the BBMap package (https://github.com/BioInfoTools/BBMap). The sole crystal structure

231     sequence for the bd-ubiquinol oxidase subunit A from *Geobacillus thermodentrificans* (Safarian

232     et al., 2016) was used as a database for a blastp (Altschul et al., 1990) search using the SAGs and

233     reference genomes as queries with an e-value cutoff of 1e-10.

234

235     Oxygen reductase horizontal gene transfer

236          The protein sequences identified from the above section were retrieved from SAGs using

237     the grep function from the list of sequence file headers from the above analysis in the SeqKit

238     package (Shen et al., 2016). Reference protein sequences for Patescibacteria were retrieved via

239     the blastp server using the Patescibacteria SAG HCO sequences as queries and selecting for hits

240     only from sequences that were assigned to Patescibacteria and/or Candidate Phyla Radiation.

241     Other reference sequences for Patescibacteria were retrieved by manual literature searches from

242     relevant studies (Nelson and Stegen, 2015; León-Zayas et al., 2017; Castelle et al., 2018). The

243     search for Patescibacteria HCOs revealed that they only encoded for the low-affinity Type A

244     HCO, and all subsequent phylogenetic analyses focused solely on this HCO type. The multi-fasta

245     file containing all HCO sequences was filtered for sequences that were greater than 400 amino

246     acids in length, and aligned with mafft (Nakamura et al., 2018) using the "--auto" option and the

247     resulting alignment was trimmed with trimal (Capella-Gutiérrez et al., 2009) to remove gaps

248     using the "-gappyout" option. A maximum likelihood phylogenetic tree was created using

249     FastTree (Price et al., 2010) using the LG model of amino acid evolution. No DPANN genome

250     to date has had a positive identification of an HCO subunit I. The methodology for the HCO

251     phylogeny was repeated for the bd-ubiquinol oxygen reductases. Phylogenetic trees were

252     visualized and annotated using the online Interactive Tree of Life tool (Letunic and Bork, 2019).

253

254    Oxidoreductase annotation and abundance

255    Enzyme Commission 1 (EC1) class family proteins (i.e., oxidoreductases) were predicted from

256    the SAGs and reference genomes using the prokka "genome.tsv" annotation files. The total

257    number of predicted protein sequences annotated as EC1 for each genome was divided by the

258    total number of predicted protein sequences to provide the percent of protein encoding genes that

259    were predicted to be oxidoreductases. This allows for a direct comparison of all the genomes that

260    exhibited wide ranges in completeness estimates.

261

262    KEGG orthology assignment of electron transport chain proteins

263    The Kyoto Encyclopedia of Gene and Genomes (KEGG) orthology (KO) annotations were

264    assigned using KofamKOALA (Aramaki et al., 2019), which uses hmmsearch (Eddy, 2011)

265    against curated hidden Markov model (HMM) KO profiles. Only KO profiles related to energy

266    transduction oxidoreductases were used to search the genomes in this study, which were

267    extracted from Supplemental Table 1 in Jelen et al. (2016). Sequences were identified as positive

268    hits if their score was greater than or equal to 50% of the sequence threshold value as calculated

269    in KofamKOALA.

270

271    16S ribosomal RNA gene phylogeny

272        16S rRNA gene sequences predicted using cmsearch (Nawrocki and Eddy, 2013) were

273    filtered for sequences that were greater than or equal to 1200 bp using bioawk

274    (https://github.com/lh3/bioawk). Sequences that were 100% identical were removed using

275    dedupe.sh (https://github.com/BioInfoTools/BBMap). Sequences were then aligned using ssu-

276    align (Nawrocki, 2009), which produces two separate alignment files for Bacteria and Archaea.

277    Next, ambiguously aligned positions were removed using ssu-mask, and sequences were re-

278    checked to ensure that the masked alignment contained sequences that were greater than or equal

279    to 1200 bp. Sequences that did not meet these threshold requirements were removed from the

280    alignment file using ssu-mask with the "--seq-r" option and list of sequences to remove. The

281    Stockholm alignment file was converted to an aligned fasta file using ssu-mask with the "--

282    stk2afa" option. The masked and filtered alignment files for Bacteria and Archaea were used to

283    create phylogenetic trees using maximum likelihood reconstruction with FastTree (Price et al.,

284    2010) with the following parameters: "-nt -gtr -cat 20 -gamma". Both trees were visualized and

285    annotated using the Interactive Tree of Life (Letunic and Bork, 2019).

286

287    **Results and Discussion**

288

289    <u>Global presence of Patescibacteria and DPANN in subsurface environments</u>

290         To improve our understanding of the deep genealogy of Bacteria and Archaea, we

291    sequenced 4,829 single amplified genomes (SAGs; Table S1) of previously under-sampled

292    microbial lineages from 46 globally distributed field sites (Figure 1; Table S1). These sites were

293    chosen based on 16S rRNA gene amplicon screens that were enriched in bacterial and archaeal

294    candidate phyla. A maximum likelihood phylogenetic tree of concatenated single-copy proteins

295    (SCP) (Figure 2) positioned 22% and 4% of SAGs within Patescibacteria (n=492) and DPANN

296    archaea (n=81). The concatenated SCP phylogenetic tree revealed the separation of

297    Patescibacteria and DPANN from other Bacteria and Archaea, respectively, which corroborates

298    other phylogenetic reconstructions using diverse sets of single copy proteins and phylogenetic

13

299    tools (Rinke et al., 2013; Brown et al., 2015; Hug et al., 2016; Williams et al., 2017; Castelle et

300    al., 2018; Dombrowski et al., 2019). Patescibacteria comprised a median relative abundance of

301    13% (range=0-81%) and DPANN comprised a median abundance of 7.5% (range=0-23%) in 33

302    analyzed environmental sites, with elevated abundances in deep-sourced aquifer environments

303    (Figure 3). Most of the Patescibacteria and DPANN SAGs originated from 13 continental

304    subsurface sites in Africa, Asia, and North America (Table S1). These results confirm that

305    Patescibacteria and DPANN are globally abundant members of subsurface microbial

306    communities, expanding on the prior genomic studies that were predominantly based on a small

307    number of study locations in North America (Rinke et al., 2013; Luef et al., 2015; Castelle et al.,

308    2018).

309

310    <u>Evidence for physical cell-cell associations</u>

311         We searched for evidence of physical cell-cell associations—an implication of obligate

312    symbiosis—by identifying genomic sequences from multiple phylogenetically distinct organisms

313    within individual SAGs. First, we searched for multiple copies of conserved, single copy protein-

314    encoding genes using checkM (Parks et al., 2015), which is a commonly used tool to detect

315    genome contamination. This approach identified 1% of Patescibacteria SAGs (5/492), 1.2% of

316    DPANN SAGs (1/81), and 0.3% of SAGs from other phyla (5/1686) as containing DNA from

317    heterogeneous sources (Table S3). Next, we searched for non-identical, near-full-length (> 1,000

318    bp) 16S rRNA genes in individual SAG assemblies. Such cases accounted for 1.5% of

319    Patescibacteria (4/262), 0% DPANN (0/56), and 0.53% for other phyla (4/758) (Table S3). A

320    Chi-square test revealed that there was a significant relationship between phyla and potential co-

321    sorted SAGs from both checkM (p-value=1.2 x $10^{-13}$; $X^2$=224.2) and 16S rRNA gene analyses

322    (p-value<2.2 x $10^{-16}$; $X^2$=238.07), but the overall contribution of Patescibacteria and DPANN to

323    the significance of co-sorted SAGs was very low (<0.5%) relative to other phyla (Figure 4). Due

324    to the incomplete SAG assemblies (Table S1), these sequencing-based approaches may

325    underestimate the overall frequency of cell-cell associations in our data set. However, they

326    consistently show that putative cell-cell associations constitute only a minor fraction of all SAGs,

327    and that Patescibacteria and DPANN are not significantly enriched in such associations relative

328    to other phyla in the studied environments. Furthermore, all identified cases of heterogeneous

329    DNA in SAG assemblies were phylogenetically unique (Table S3), in contrast to the recurring

330    Nanoarchaeota-Crenarchaeota symbiotic associations found using the same techniques in hot

331    springs in prior studies (Munson-McGee et al., 2015; Jarett et al., 2018). Also, in mammalian

332    oral microbiomes, Saccharibacteria have been shown to be specifically associated with

333    Actinobacteria hosts (He et al., 2015; Cross et al., 2019). This suggests that the infrequent and

334    inconsistent presence of taxonomically heterogeneous DNA in SAGs most likely originated from

335    non-specific aggregation of multiple cells and/or attachment of extracellular DNA.

336         Based on a small number of transmission electron micrograph observations, it has been

337    suggested that Patescibacteria associations with other microorganisms may be fragile (Luef et al.,

338    2015). Thus, we cannot rule out the possibility that some Patescibacteria and DPANN cells were

339    attached to host cells *in situ* and became detached during sample collection and processing. To

340    reduce the risk of dispersing natural cell aggregates and associations, we performed only a gentle

341    mixing of the analyzed samples in preparation for cell sorting. In prior studies, similar techniques

342    successfully revealed host-symbiont associations in termite guts (Hongoh et al., 2008), marine

343    plankton (Martinez-Garcia et al., 2012) and hot springs (Jarett et al., 2018). This approach was

344    also used to determine symbiotic associations between anaerobic methane-oxidizing archaea and

345     their syntrophic partners in natural consortia from methane seeps (Hatzenpichler et al., 2016). It

346     is worthy to note that the Saccharibacteria-Actinobacteria symbiont-host relationship was only

347     disrupted by physical passage through a narrow-gauge needle multiple times (He et al., 2015).

348     Also, putatively co-sorted SAGs of Nanoarchaeota and Crenarchaeota from iron oxide microbial

349     mats were treated by a repeated physical disruption through multiple wash cycles and density

350     gradient centrifugation, from which co-sorted cells were obtained (Jarett et al., 2018). Thus,

351     although the techniques applied here may underestimate the overall counts of cell-cell

352     associations *in situ*, we found no evidence for Patescibacteria and DPANN to be enriched in such

353     associations relative to other phyla, and to form lineage-specific associations in the analyzed

354     environments.

355

356     <u>Cell diameters</u>

357     We employed calibrated index fluorescence-activated cell sorting (FACS) to determine

358     physical diameters of individual cells that were used in SAG generation (Stepanauskas et al.,

359     2017). This indicated that Patescibacteria (n=273) and DPANN (n=29) cells are extremely small

360     across their entire phylogenetic breadth, with median estimated diameters of 0.2 μm (Figure 5).

361     Several cases of larger, outlier diameter estimates may be due to attachment to other cells and

362     particles, cellular division, methodological artifacts, or true biological variation. The low

363     frequency of Patescibacteria and DPANN DNA recovery from larger particles (Table S1; Figure

364     5) provides further indication that most of these cells are not attached to other microorganisms.

365     Likewise, most of the SAGs with identified heterogeneous genome sources were larger than their

366     phylum median cell diameters (Table S3), which is consistent with their aggregation with other

367     cells.

368     To further investigate the composition of extremely small cells, we generated a

369     complementary library of SAGs from a single subsurface sample (AG-274; Table S1) with a

370     FACS gate targeting only ≤0.3 μm particles. Confirming our expectations, >90% of SAGs in this

371     cell diameter-specific library were composed of Patescibacteria and DPANN (Figure 3). The

372     obtained cellular size ranges are consistent with a prior report, which was based on transmission

373     electron micrographs from one field study site (Luef et al., 2015). These cell diameters

374     approximate the lower theoretical limits for cellular life (Maniloff et al., 1997).

375

376     General genome features

377     To identify functional coding potential differences of Patescibacteria and DPANN

378     compared to other Bacteria and Archaea, we performed a principal component analysis (PCA)

379     using the relative abundance of clusters of orthologous groups (COG) as input variables with

380     SAGs that had at least 30% completeness and a near full-length 16S rRNA gene (Figure 6). This

381     showed a clear separation of Patescibacteria and DPANN from other bacteria and archaea along

382     the first component (PC1) (Wilcoxon signed-rank test; p-value < 2.2 x 10$^{-16}$). Importantly, well-

383     described symbionts (Table S4) separated from both Patescibacteria along PC1 and DPANN

384     along PC2 (p-value = 2.57 x 10$^{-8}$ and 1.0 x 10$^{-7}$ for Patescibacteria and DPANN, respectively).

385     The only lineages that clustered with Patescibacteria and DPANN along PC1 and PC2 were

386     Dependentiae and Tenericutes, respectively.

387     The COG categories with the greatest negative effect on PC1, indicative of their relative

388     depletion in Patescibacteria and DPANN, included E (amino acid metabolism and transport), C

389     (energy production and conversion), P (inorganic ion transport and metabolism), and H

390     (coenzyme transport and metabolism). The COG categories with the greatest positive effect on

17

391   PC1, indicative of their relatively high fraction in genomes of Patescibacteria and DPANN,

392   included D (cell cycle control and mitosis) and O (post-translational modification, protein

393   turnover, chaperone functions). Archaea separated from bacteria along the second component

394   (PC2) (p-value $< 2.2 \times 10^{-16}$) primarily by their relative enrichment in COG categories B

395   (chromatin structure and dynamics), K (transcription), and S (unknown functions). This reflects

396   the major inter-domain differences in DNA packing and transcription, and the greater fraction of

397   archaeal genomes remaining uncharacterized, as compared to the genomes of Bacteria.

398      Genomes recently shaped by symbiosis often have low coding densities due to rapid gene

399   loss and pseudogene formation (McCutcheon and Moran, 2012). Inconsistent with this pattern,

400   we found the coding density of Patescibacteria and DPANN (median = 91%) to be typical of

401   Bacteria and Archaea (median = 90%), while well-characterized symbionts were separated by

402   their lower coding density (Figure 7a) (median = 0.87%, p-value = 0.035 and 0.028 compared to

403   Patescibacteria and DPANN). Although the reduced genome size of Patescibacteria and DPANN

404   has been viewed as an indication of a symbiotic lifestyle (Castelle et al., 2018), similar genome

405   sizes (1-2 Mbp) are typical among free-living, marine plankton (Swan et al., 2013; Giovannoni et

406   al., 2014). Furthermore, recent synthetic biology experimentation has pushed the minimal

407   genome size limit of a free-living microorganism to ~0.5 Mbp (Hutchison et al., 2016), far below

408   the predicted sizes of Patescibacteria and DPANN genomes. Collectively, these general genome

409   features of Patescibacteria and DPANN do not provide convincing evidence of an obligate

410   symbiotic lifestyle.

411      In this context, the observed gene content similarities between Patescibacteria and

412   Dependentiae, and between DPANN and most Tenericutes are intriguing (Figure 6).

413   Dependentiae is a candidate bacterial phylum that has been noted for its reduced coding

18

414    potential, including a depletion in electron transport chain components (McLean et al., 2013;

415    Yeoh et al., 2016). It has been speculated that these characteristics indicate a symbiotic lifestyle,

416    with energy acquired from hosts via ATP/ADP translocases, which has been confirmed

417    experimentally in a few Dependentiae members (Delafont et al., 2015; Pagnier et al., 2015; Deeg

418    et al., 2019). The well-characterized members of the bacterial phylum Tenericutes consist mostly

419    of obligate pathogens with reduced genomes (Moran and Wernegreen, 2000). Interestingly, most

420    Tenericutes are able to grow as free-living cells in rich media solely by fermentation (Tully et al.,

421    1977), and were originally hypothesized to represent ancient lineages of life due to their small

422    genome sizes and limited metabolisms (Morowitz, 1984). While we found all analyzed

423    Dependentiae and most Tenericutes deplete in oxidoreductases (Figures 7b; Figure 8), only

424    Tenericutes had a consistently low coding density (median = 71%) that is a characteristic of

425    recently evolved symbionts (McCutcheon and Moran, 2012) (Figure 7a). Thus, we hypothesize

426    that these two phyla cluster with the Patescibacteria and DPANN due to similar metabolic

427    features arrived at by convergent evolutionary processes.

428

429    Oxygen reductase genes

430         In search for an alternative explanation for the unique genealogy, genome content, and

431    cell sizes of Patescibacteria and DPANN, we examined their energy metabolic coding potential.

432    We found that only 0.6% of Patescibacteria SAGs (3/492) and none of the DPANN SAGs (0/81)

433    from these samples encoded for homologs of oxygen reductases ($O_2$red), as indicated by the

434    presence of oxygen-binding subunit I of either the heme-copper oxidase (HCO) or bd-ubiquinol

435    (bd) oxidase families. The incomplete genome recovery from individual SAGs cannot explain

436    this pattern, because the 492 Patescibacteria SAGs and 81 DPANN SAGs correspond to a

437    cumulative assembly of 162 and 27 randomly sampled, complete genomes. Furthermore, a

438    phylogenetic analysis revealed that all three oxygen reductases from Patescibacteria SAGs form

439    a cluster with other Patescibacteria sequences (Brown et al., 2015; Nelson and Stegen, 2015;

440    León-Zayas et al., 2017) that is nested within a clade comprised of other phyla (Figure 9). We

441    infer these phylogenetic relationships as an indication of a relatively recent horizontal gene

442    transfer (HGT), likely from Proteobacteria and Firmicutes for the HCO and bd sequences,

443    respectively. Although we did not detect any homologs of oxygen reductases in DPANN SAGs

444    from our samples, the publicly available bd $O_2$red sequences from DPANN metagenome bins

445    and isolates formed a clade with Actinobacteria and Firmicutes, which we also infer as likely

446    products of relatively recent HGT events. The topology of these $O_2$red phylogenetic trees is

447    consistent with prior reports, which have also been interpreted as evidence for prevalent HGT of

448    oxygen reductase genes among other phyla (Brochier-Armanet et al., 2009; Gribaldo et al., 2009;

449    Borisov et al., 2011). This suggests that the absence of oxygen reductases in Patescibacteria and

450    DPANN is ancestral and not a result of gene loss due to adaptations to symbiosis, as previously

451    hypothesized (Castelle et al., 2018).

452

453    <u>Distribution of electron transport chain complexes</u>

454        Patescibacteria and DPANN were depleted in the entire family of oxidoreductase enzyme

455    genes compared to other bacteria and archaea, ($p$-value $< 2.2 \times 10^{-16}$) (Figures 7b, 8). This

456    depletion was also significant in relation to symbionts with their comparatively small genome

457    sizes ($p$-value $< 0.05$). Oxidoreductases are key components of both aerobic and anaerobic

458    respiratory pathways (Jelen et al., 2016), so underrepresentation of them would suggest reduced

459    functionality of these energy transduction mechanisms. Accordingly, none of the Patescibacteria

460    and DPANN genomes were found to encode a complete ETC consisting of all four complexes

461    (Figure 8). Putative homologs of at least two of the four ETC complexes were found only in 3%

462    and 11% of Patescibacteria and DPANN genomes, respectively. We found putative homologs of

463    genes encoding individual complexes I, II, III, and IV in 0%, 2%, 3%, and 14% of

464    Patescibacteria genomes. The corresponding numbers for DPANN were 7%, 4%, 0%, and 21%.

465    Some of these computationally predicted genes are only distantly related to experimentally

466    verified homologs and therefore may constitute false positives. These findings are consistent

467    with the lack of complete ETC reports in prior studies of Patescibacteria genomes (Brown et al.,

468    2015), with the sole exception of a tentative nitric oxide respiration operon found in a single

469    metagenome bin (Castelle et al., 2017). The sparse and scattered distribution of the putative ETC

470    gene homologs in Patescibacteria and DPANN (Figure 8) suggest horizontal gene transfer

471    origins rather than ancestral inheritance. This is consistent with the phylogenetic reconstructions

472    of other energy transducing genes identified in Patescibacteria, which also suggest evolutionary

473    origins from horizontal gene transfer (Jaffe et al., 2019). Collectively, our observations indicate

474    that the absence of complete electron transport chains in Patescibacteria and DPANN is an

475    ancestral feature, which we propose is more parsimonious than multiple gene loss events due to

476    obligate symbiosis (Brown et al., 2015; Hug et al., 2016; Castelle et al., 2018; Dombrowski et

477    al., 2019; Méheust et al., 2019).

478

479    <u>Respiration activity</u>

480          To experimentally test for the presence of active oxidoreductases in a subsurface

481    microbial community, we employed the fluorogenic oxidoreductase probe RedoxSensor Green

482    on a deep groundwater sample from South Dakota. This revealed a wide range in fluorescence

21

483    intensity in phylogenetically diverse cells, with none of the Patescibacteria cells exceeding the

484    fluorescence of particles in a heat-killed, negative control (Figure 10). To the best of our

485    knowledge, RedoxSensor Green has not been tested extensively on diverse microbial lineages,

486    therefore these results should be considered tentative. Nonetheless, both genome content and *in*

487    *situ* physiology analyses indicate the absence of respiration in Patescibacteria and DPANN,

488    which corroborates earlier reports of these lineages containing few, if any, components of energy

489    transducing pathways other than fermentation (Castelle et al., 2018).

490

491    16S rRNA gene phylogeny

492        The placement of Patescibacteria and DPANN in the tree of life is widely debated (Hug

493    et al., 2016; Williams et al., 2017; Dombrowski et al., 2019). Most current phylogenetic

494    inferences are based on concatenated single-copy proteins (CSCP), which has the advantage of

495    higher phylogenetic resolution, as compared to phylogenies of individual genes (Rinke et al.,

496    2013). However, the unknown genetic change at heterogeneously evolving sites and large

497    sequence divergence may limit the accuracy of such trees (Pace, 2009; Dombrowski et al., 2019).

498    To complement the CSCP-resolved genealogy (Figure 2), we performed a large-scale

499    phylogenetic analysis of the well-established 16S rRNA gene (Woese, 2002) (length > 1,200 bp)

500    separately for Bacteria and Archaea. The obtained phylogenetic inference (Figure 8) supported

501    the separation of Patescibacteria and DPANN from other bacterial and archaeal lineages, in

502    agreement with the phylogenies based on CSCP genes (Castelle et al., 2018) (Figure 2) and a

503    recent large scale bacterial 16S rRNA gene tree (Schulz et al., 2017). Importantly, we did not

504    observe grouping of Patescibacteria with fast-evolving lineages (e.g., obligate insect symbionts

505    and Tenericutes) that could be due to long branch attraction in the 16S rRNA gene phylogeny.

22

506    This suggests that the divergent branching of Patescibacteria and DPANN is probably not a

507    result of recent, accelerated divergence.

508

509    **Concluding Remarks**

510        Using the collective evidence from cell-cell association, coding potential and

511    phylogenomic analyses, we propose a new explanation of the unusual biological features of

512    Patescibacteria and DPANN. Although the Patescibacteria and DPANN contain symbionts

513    (Huber et al., 2002; Podar et al., 2013; Gong et al., 2014; He et al., 2015; Munson-McGee et al.,

514    2015; Jarett et al., 2018; Cross et al., 2019; Hamm et al., 2019) and auxotrophies (Castelle et al.,

515    2018; Dombrowski et al., 2019), we believe that there is not sufficient evidence to conclude that

516    an ancestral adaptation to symbiosis has led to the reduction of their cell sizes and coding

517    potential (Castelle et al., 2018; Dombrowski et al., 2019; Méheust et al., 2019). Instead, our data

518    indicate that most Patescibacteria and DPANN do not form symbiotic cell-cell associations in

519    subsurface environments, and that their divergent coding potential, small genomes, and small

520    cell sizes may be a result of a primitive energy metabolism that relies solely on substrate-level

521    phosphorylation (fermentation), potentially preceding the evolution of electron transport

522    phosphorylation (respiration). Auxotrophies are very common among microorganisms, and

523    represent a wide range of dependencies for exogenous cellular components (Zengler and

524    Zaramela, 2018). Patescibacteria and DPANN may be on the extreme end of the spectrum in

525    their dependence on other community members, perhaps a reflection of an ancient evolutionary

526    strategy to limit cellular biosynthetic energy requirements, as energetic allocation is a major

527    driver of genome evolution in bacteria and archaea (Lynch and Marinov, 2015).

528

23

529    The authors declare no conflict of interest.

530

531    **Acknowledgements**

532

558

## Author contributions

560    JPB led data analyses and manuscript preparation. RS developed the concept and managed the

561    project, with contributions by TW, TCO, DM, JAE, JPB and EDB. EDB, JMB, FS, JKJ, OB, KC

562    contributed to data analyses. NJP performed cell sorting and size calibration at Bigelow

563    Laboratory. TCO, DPM, PD, NVR, JRS, BPH, KAK, SMS, MSE, HAB and MBS oversaw field

564    sample collection. All authors contributed to data interpretation and manuscript preparation.

565

## References

567    Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. (1990). Basic local

568        alignment search tool. *J. Mol. Biol.* 215, 403–410. doi:10.1016/S0022-2836(05)80360-2.

569    Altschul, S. F., Madden, T. L., Schäffer, A. A., Zhang, J., Zhang, Z., Miller, W., et al. (1997).

570        Gapped BLAST and PSI-BLAST: A new generation of protein database search programs.

571        *Nucleic Acids Res.* 25, 3389–3402. doi:10.1093/nar/25.17.3389.

572    Aramaki, T., Blanc-Mathieu, R., Endo, H., Ohkubo, K., Kanehisa, M., Goto, S., et al. (2019).

573        KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive score

574        threshold. *bioRxiv*, 602110. doi:10.1101/602110.

575    Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., et al.

576        (2012). SPAdes: A new genome assembly algorithm and its applications to single-cell

577        sequencing. *J. Comput. Biol.* 19, 455–477. doi:10.1089/cmb.2012.0021.

578    Bateman, A., Martin, M. J., O'Donovan, C., Magrane, M., Alpi, E., Antunes, R., et al. (2017).

579        UniProt: The universal protein knowledgebase. *Nucleic Acids Res.* 45, D158–D169.

580        doi:10.1093/nar/gkw1099.

581    Becraft, E. D., Woyke, T., Jarett, J., Ivanova, N., Godoy-Vitorino, F., Poulton, N., et al. (2017).

582        Rokubacteria: Genomic giants among the uncultured bacterial phyla. *Front. Microbiol.* 8,

583        1–12. doi:10.3389/fmicb.2017.02264.

584    Borisov, V. B., Gennis, R. B., Hemp, J., and Verkhovsky, M. I. (2011). The cytochrome bd

585        respiratory oxygen reductases. *Biochim. Biophys. Acta - Bioenerg.* 1807, 1398–1413.

586        doi:10.1016/j.bbabio.2011.06.016.

587    Brochier-Armanet, C., Talla, E., and Gribaldo, S. (2009). The multiple evolutionary histories of

588        dioxygen reductases: Implications for the origin and evolution of aerobic respiration. *Mol.*

589        *Biol. Evol.* 26, 285–297. doi:10.1093/molbev/msn246.

590    Brown, C. T., Hug, L. A., Thomas, B. C., Sharon, I., Castelle, C. J., Singh, A., et al. (2015).

591        Unusual biology across a group comprising more than 15% of domain Bacteria. *Nature* 523,

592        208–211. doi:10.1038/nature14486.

593    Brown, C. T., Olm, M. R., Thomas, B. C., and Banfield, J. F. (2016). Measurement of bacterial

594        replication rates in microbial communities. *Nat. Biotechnol.* 34, 1256–1263.

595        doi:10.1038/nbt.3704.

596    Capella-Gutiérrez, S., Silla-Martínez, J. M., and Gabaldón, T. (2009). trimAl: A tool for

597        automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25,

598   1972–1973. doi:10.1093/bioinformatics/btp348.

599 Castelle, C. J., and Banfield, J. F. (2018). Major New Microbial Groups Expand Diversity and

600   Alter our Understanding of the Tree of Life. *Cell* 172, 1181–1197.

601   doi:10.1016/j.cell.2018.02.016.

602 Castelle, C. J., Brown, C. T., Anantharaman, K., Probst, A. J., Huang, R. H., and Banfield, J. F.

603   (2018). Biosynthetic capacity, metabolic variety and unusual biology in the CPR and

604   DPANN radiations. *Nat. Rev. Microbiol.* 16, 629–645. doi:10.1038/s41579-018-0076-2.

605 Castelle, C. J., Brown, C. T., Thomas, B. C., Williams, K. H., and Banfield, J. F. (2017).

606   Unusual respiratory capacity and nitrogen metabolism in a Parcubacterium (OD1) of the

607   Candidate Phyla Radiation. *Sci. Rep.* 7, 1–12. doi:10.1038/srep40101.

608 Castelle, C. J., Wrighton, K. C., Thomas, B. C., Hug, L. A., Brown, C. T., Wilkins, M. J., et al.

609   (2015). Genomic expansion of domain archaea highlights roles for organisms from new

610   phyla in anaerobic carbon cycling. *Curr. Biol.* 25, 690–701. doi:10.1016/j.cub.2015.01.014.

611 Chen, I. M. A., Chu, K., Palaniappan, K., Pillay, M., Ratner, A., Huang, J., et al. (2019). IMG/M

612   v.5.0: An integrated data management and comparative analysis system for microbial

613   genomes and microbiomes. *Nucleic Acids Res.* 47, D666–D677. doi:10.1093/nar/gky901.

614 Criscuolo, A., and Gribaldo, S. (2010). BMGE (Block Mapping and Gathering with Entropy): A

615   new software for selection of phylogenetic informative regions from multiple sequence

616   alignments. *BMC Evol. Biol.* 10. doi:10.1186/1471-2148-10-210.

617 Cross, K. L., Campbell, J. H., Balachandran, M., Campbell, A. G., Cooper, S. J., Griffen, A., et

618   al. (2019). Targeted isolation and cultivation of uncultivated bacteria by reverse genomics.

619   *Nat. Biotechnol.* doi:10.1038/s41587-019-0260-6.

620 Deeg, C. M., Zimmer, M. M., George, E. E., Husnik, F., Keeling, P. J., and Suttle, C. A. (2019).

621     Chromulinavorax destructans, a pathogen of microzooplankton that provides a window into

622     the enigmatic candidate phylum dependentiae. *PLoS Pathog.* 15, 1–18.

623     doi:10.1371/journal.ppat.1007801.

624  Delafont, V., Samba-Louaka, A., Bouchon, D., Moulin, L., and Héchard, Y. (2015). Shedding

625     light on microbial dark matter: A TM6 bacterium as natural endosymbiont of a free-living

626     amoeba. *Environ. Microbiol. Rep.* 7, 970–978. doi:10.1111/1758-2229.12343.

627  Dombrowski, N., Lee, J. H., Williams, T. A., Offre, P., and Spang, A. (2019). Genomic

628     diversity, lifestyles and evolutionary origins of DPANN archaea. *FEMS Microbiol. Lett.*

629     366, 1–12. doi:10.1093/femsle/fnz008.

630  Eddy, S. R. (2011). Accelerated profile HMM searches. *PLoS Comput. Biol.* 7.

631     doi:10.1371/journal.pcbi.1002195.

632  Eloe-Fadrosh, E. A., Paez-Espino, D., Jarett, J., Dunfield, P. F., Hedlund, B. P., Dekas, A. E., et

633     al. (2016). Global metagenomic survey reveals a new bacterial candidate phylum in

634     geothermal springs. *Nat. Commun.* 7, 1–10. doi:10.1038/ncomms10476.

635  Giovannoni, S. J., Cameron Thrash, J., and Temperton, B. (2014). Implications of streamlining

636     theory for microbial ecology. *ISME J.* 8, 1553–1565. doi:10.1038/ismej.2014.60.

637  Gong, J., Qing, Y., Guo, X., and Warren, A. (2014). "Candidatus Sonnebornia yantaiensis", a

638     member of candidate division OD1, as intracellular bacteria of the ciliated protist

639     Paramecium bursaria (Ciliophora, Oligohymenophorea). *Syst. Appl. Microbiol.* 37, 35–41.

640     doi:10.1016/j.syapm.2013.08.007.

641  Gribaldo, S., Talla, E., and Brochier-Armanet, C. (2009). Evolution of the haem copper oxidases

642     superfamily: a rooting tale. *Trends Biochem. Sci.* 34, 375–381.

643     doi:10.1016/j.tibs.2009.04.002.

644     Hamm, J. N., Erdmann, S., Eloe-Fadrosh, E. A., Angeloni, A., Zhong, L., Brownlee, C., et al.

645         (2019). Unexpected host dependency of Antarctic Nanohaloarchaeota. *Proc. Natl. Acad.*

646         *Sci. U. S. A.* 116, 14661–14670. doi:10.1073/pnas.1905179116.

647     Hatzenpichler, R., Connon, S. A., Goudeau, D., Malmstrom, R. R., Woyke, T., and Orphan, V. J.

648         (2016). Visualizing in situ translational activity for identifying and sorting slow-growing

649         archaeal - bacterial consortia. *Proc. Natl. Acad. Sci. U. S. A.* 113, E4069–E4078.

650         doi:10.1073/pnas.1603757113.

651     He, X., McLean, J. S., Edlund, A., Yooseph, S., Hall, A. P., Liu, S. Y., et al. (2015). Cultivation

652         of a human-associated TM7 phylotype reveals a reduced genome and epibiotic parasitic

653         lifestyle. *Proc. Natl. Acad. Sci. U. S. A.* 112, 244–249. doi:10.1073/pnas.1419038112.

654     Hershey, O. S., Kallmeyer, J., Wallace, A., Barton, M. D., and Barton, H. A. (2018). High

655         microbial diversity despite extremely low biomass in a deep karst aquifer. *Front. Microbiol.*

656         9, 1–13. doi:10.3389/fmicb.2018.02823.

657     Hongoh, Y., Sharma, V. K., Prakash, T., Noda, S., Taylor, T. D., Kudo, T., et al. (2008).

658         Complete genome of the uncultured Termite Group 1 bacteria in a single host protist cell.

659         *Proc. Natl. Acad. Sci. U. S. A.* 105, 5555–5560. doi:10.1073/pnas.0801389105.

660     Huber, H., Hohn, M. J., Rachel, R., Fuchs, T., Wimmer, V. C., and Stetter, K. O. (2002). A new

661         phylum of Archaea represented by a nanosized hyperthermophilic symbiont. *Nature* 417,

662         63–67. doi:10.1038/417063a.

663     Hug, L. A., Baker, B. J., Anantharaman, K., Brown, C. T., Probst, A. J., Castelle, C. J., et al.

664         (2016). A new view of the tree of life. *Nat. Microbiol.* 1, 1–6.

665         doi:10.1038/nmicrobiol.2016.48.

666     Huntemann, M., Ivanova, N. N., Mavromatis, K., Tripp, H. J., Paez-Espino, D., Tennessen, K., et

667      al. (2016). The standard operating procedure of the DOE-JGI Metagenome Annotation

668      Pipeline (MAP v.4). *Stand. Genomic Sci.* 11, 1–5. doi:10.1186/s40793-016-0138-x.

669    Hutchison, C. A., Chuang, R. Y., Noskov, V. N., Assad-Garcia, N., Deerinck, T. J., Ellisman, M.

670      H., et al. (2016). Design and synthesis of a minimal bacterial genome. *Science (80-. ).* 351.

671      doi:10.1126/science.aad6253.

672    Hyatt, D., Chen, G. L., LoCascio, P. F., Land, M. L., Larimer, F. W., and Hauser, L. J. (2010).

673      Prodigal: Prokaryotic gene recognition and translation initiation site identification. *BMC*

674      *Bioinformatics* 11. doi:10.1186/1471-2105-11-119.

675    Jaffe, A. L., Castelle, C. J., Carnevali, P. B. M., Gribaldo, S., and Banfield, J. F. (2019). The rise

676      of diversity in metabolic platforms across the Candidate Phyla Radiation. *bioRxiv*, 1–29.

677    Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T., and Aluru, S. (2018). High

678      throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat.*

679      *Commun.* 9, 1–8. doi:10.1038/s41467-018-07641-9.

680    Jarett, J. K., Nayfach, S., Podar, M., Inskeep, W., Ivanova, N. N., Munson-Mcgee, J., et al.

681      (2018). Single-cell genomics of co-sorted Nanoarchaeota suggests novel putative host

682      associations and diversification of proteins involved in symbiosis 06 Biological Sciences

683      0604 Genetics. *Microbiome* 6, 1–14. doi:10.1186/s40168-018-0539-8.

684    Jelen, B. I., Giovannelli, D., and Falkowski, P. G. (2016). The Role of Microbial Electron

685      Transfer in the Coevolution of the Biosphere and Geosphere. *Annu. Rev. Microbiol.* 70, 45–

686      62. doi:10.1146/annurev-micro-102215-095521.

687    León-Zayas, R., Peoples, L., Biddle, J. F., Podell, S., Novotny, M., Cameron, J., et al. (2017).

688      The metabolic potential of the single cell genomes obtained from the Challenger Deep,

689      Mariana Trench within the candidate superphylum Parcubacteria (OD1). *Environ.*

690    *Microbiol.* 19, 2769–2784. doi:10.1111/1462-2920.13789.

691    Letunic, I., and Bork, P. (2019). Interactive Tree Of Life (iTOL) v4: recent updates and new

692        developments. *Nucleic Acids Res.* 47, W256–W259. doi:10.1093/nar/gkz239.

693    Luef, B., Frischkorn, K. R., Wrighton, K. C., Holman, H. Y. N., Birarda, G., Thomas, B. C., et

694        al. (2015). Diverse uncultivated ultra-small bacterial cells in groundwater. *Nat. Commun.* 6,

695        1–8. doi:10.1038/ncomms7372.

696    Lynch, M., and Marinov, G. K. (2015). The bioenergetic costs of a gene. *Proc. Natl. Acad. Sci.*

697        *U. S. A.* 112, 15690–15695. doi:10.1073/pnas.1514974112.

698    Maniloff, J., Nealson, K. H., Psenner, R., Loferer, M., and Folk, R. L. (1997). Nannobacteria:

699        Size Limits and Evidence. *Science (80-. ).* 276, 1773–1776.

700    Martinez-Garcia, M., Brazel, D., Poulton, N. J., Swan, B. K., Gomez, M. L., Masland, D., et al.

701        (2012). Unveiling in situ interactions between marine protists and bacteria through single

702        cell sequencing. *ISME J.* 6, 703–707. doi:10.1038/ismej.2011.126.

703    McCutcheon, J. P., and Moran, N. A. (2012). Extreme genome reduction in symbiotic bacteria.

704        *Nat. Rev. Microbiol.* 10, 13–26. doi:10.1038/nrmicro2670.

705    McLean, J. S., Lombardo, M. J., Badger, J. H., Edlund, A., Novotny, M., Yee-Greenbaum, J., et

706        al. (2013). Candidate phylum TM6 genome recovered from a hospital sink biofilm provides

707        genomic insights into this uncultivated phylum. *Proc. Natl. Acad. Sci. U. S. A.* 110.

708        doi:10.1073/pnas.1219809110.

709    Méheust, R., Burstein, D., Castelle, C. J., and Banfield, J. F. (2019). The distinction of CPR

710        bacteria from other bacteria based on protein family content. *Nat. Commun.* 10.

711        doi:10.1038/s41467-019-12171-z.

712    Moran, N. A., and Wernegreen, J. J. (2000). Lifestyle evolution in symbiotic bacteria: Insights

31

713      from genomics. *Trends Ecol. Evol.* 15, 321–326. doi:10.1016/S0169-5347(00)01902-9.

714     Morowitz, H. (1984). The Completeness of Molecular Biology. *Isr. J. Med. Sci.* 20, 750–753.

715     Moser, D. P., Hamilton-Brehm, S. D., Fisher, J. C., Bruckner, J. C., Kruger, B., and Sackett, J.

716       (2015). Radiochemically-supported microbial communities: a potential mechanism for

717       biocolloid production of importance to actinide transport.

718     Munson-McGee, J. H., Field, E. K., Bateson, M., Rooney, C., Stepanauskas, R., and Young, J.

719       (2015). Distribution across Yellowstone National Park Hot Springs. *Appl. Environ.*

720       *Microbiol.* 81, 7860–7868. doi:10.1128/AEM.01539-15.Editor.

721     Nakamura, T., Yamada, K. D., Tomii, K., and Katoh, K. (2018). Parallelization of MAFFT for

722       large-scale multiple sequence alignments. *Bioinformatics* 34, 2490–2492.

723       doi:10.1093/bioinformatics/bty121.

724     Nawrocki, E. P. (2009). Structural RNA Homology Search and Alignment Using Covariance

725       Models.

726     Nawrocki, E. P., and Eddy, S. R. (2013). Infernal 1.1: 100-fold faster RNA homology searches.

727       *Bioinformatics* 29, 2933–2935. doi:10.1093/bioinformatics/btt509.

728     Nelson, W. C., and Stegen, J. C. (2015). The reduced genomes of Parcubacteria (OD1) contain

729       signatures of a symbiotic lifestyle. *Front. Microbiol.* 6, 1–14.

730       doi:10.3389/fmicb.2015.00713.

731     Nurk, S., Bankevich, A., Antipov, D., Gurevich, A. A., Korobeynikov, A., Lapidus, A., et al.

732       (2013). Assembling single-cell genomes and mini-metagenomes from chimeric MDA

733       products. *J. Comput. Biol.* 20, 714–737. doi:10.1089/cmb.2013.0084.

734     Pace, N. R. (2009). Mapping the Tree of Life: Progress and Prospects. *Microbiol. Mol. Biol. Rev.*

735       73, 565–576. doi:10.1128/mmbr.00033-09.

736     Pagnier, I., Yutin, N., Croce, O., Makarova, K. S., Wolf, Y. I., Benamar, S., et al. (2015). Babela

737          massiliensis, a representative of a widespread bacterial phylum with unusual adaptations to

738          parasitism in amoebae. *Biol. Direct* 10, 1–17. doi:10.1186/s13062-015-0043-z.

739     Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P., and Tyson, G. W. (2015).

740          CheckM: Assessing the quality of microbial genomes recovered from isolates, single cells,

741          and metagenomes. *Genome Res.* 25, 1043–1055. doi:10.1101/gr.186072.114.

742     Podar, M., Makarova, K. S., Graham, D. E., Wolf, Y. I., Koonin, E. V., and Reysenbach, A. L.

743          (2013). Insights into archaeal evolution and symbiosis from the genomes of a nanoarchaeon

744          and its inferred crenarchaeal host from Obsidian Pool, Yellowstone National Park. *Biol.*

745          *Direct* 8, 1–20. doi:10.1186/1745-6150-8-9.

746     Price, M. N., Dehal, P. S., and Arkin, A. P. (2010). FastTree 2 - Approximately maximum-

747          likelihood trees for large alignments. *PLoS One* 5. doi:10.1371/journal.pone.0009490.

748     Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., et al. (2013). The SILVA

749          ribosomal RNA gene database project: Improved data processing and web-based tools.

750          *Nucleic Acids Res.* 41, 590–596. doi:10.1093/nar/gks1219.

751     Rinke, C., Schwientek, P., Sczyrba, A., Ivanova, N. N., Anderson, I. J., Cheng, J. F., et al.

752          (2013). Insights into the phylogeny and coding potential of microbial dark matter. *Nature*

753          499, 431–437. doi:10.1038/nature12352.

754     Sackett, J. D. (2018). Prokaryotic diversity and aqueous geochemistry of subsurface

755          environments of the Death Valley Regional Flow System.

756     Sackett, J. D., Huerta, D. C., Kruger, B. R., Hamilton-Brehm, S. D., and Moser, D. P. (2018). A

757          comparative study of prokaryotic diversity and physicochemical characteristics of devils

758          hole and the ash meadows fish conservation facility, a constructed analog. *PLoS One* 13, 1–

759   21. doi:10.1371/journal.pone.0194404.

760   Sackett, J. D., Kruger, B. R., Becraft, E. D., Jarett, J. K., Stepanauskas, R., Woyke, T., et al.

761   (2019). Four Draft Single-Cell Genome Sequences of Novel, Nearly Identical

762   Kiritimatiellaeota Strains Isolated from the Continental Deep Subsurface . *Microbiol.*

763   *Resour. Announc.* 8, 1–4. doi:10.1128/mra.01249-18.

764   Safarian, S., Rajendran, C., Müller, H., Preu, J., Langer, J. D., Ovchinnikov, S., et al. (2016).

765   Structure of a bd oxidase indicates similar mechanisms for membraneintegrated oxygen

766   reductases. *Science (80-. )*. 352, 583–586. doi:10.1126/science.aaf2477.

767   Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B., et al.

768   (2009). Introducing mothur: Open-source, platform-independent, community-supported

769   software for describing and comparing microbial communities. *Appl. Environ. Microbiol.*

770   75, 7537–7541. doi:10.1128/AEM.01541-09.

771   Schulz, F., Eloe-Fadrosh, E. A., Bowers, R. M., Jarett, J., Nielsen, T., Ivanova, N. N., et al.

772   (2017). Towards a balanced view of the bacterial tree of life. *Microbiome* 5, 140.

773   doi:10.1186/s40168-017-0360-9.

774   Seemann, T. (2014). Prokka: Rapid prokaryotic genome annotation. *Bioinformatics* 30, 2068–

775   2069. doi:10.1093/bioinformatics/btu153.

776   Shen, W., Le, S., Li, Y., and Hu, F. (2016). SeqKit: A cross-platform and ultrafast toolkit for

777   FASTA/Q file manipulation. *PLoS One* 11, 1–10. doi:10.1371/journal.pone.0163962.

778   Sousa, F. L., Alves, R. J., Pereira-Leal, J. B., Teixeira, M., and Pereira, M. M. (2011). A

779   bioinformatics classifier and database for Heme-Copper oxygen reductases. *PLoS One* 6, 1–

780   9. doi:10.1371/journal.pone.0019117.

781   Stepanauskas, R., Fergusson, E. A., Brown, J., Poulton, N. J., Tupper, B., Labonté, J. M., et al.

34

782      (2017). Improved genome recovery and integrated cell-size analyses of individual

783      uncultured microbial cells and viral particles. *Nat. Commun.* 8. doi:10.1038/s41467-017-

784      00128-z.

785      Swan, B. K., Tupper, B., Sczyrba, A., Lauro, F. M., Martinez-Garcia, M., Gonźalez, J. M., et al.

786      (2013). Prevalent genome streamlining and latitudinal divergence of planktonic bacteria in

787      the surface ocean. *Proc. Natl. Acad. Sci. U. S. A.* 110, 11463–11468.

788      doi:10.1073/pnas.1304246110.

789      Team, Rs. (2016). RStudio: Integrated Development for R.

790      Thomas, J. M., Moser, D. P., Fisher, J. C., Reihle, J., Wheatley, A., Hershey, R. L., et al. (2013).

791      Using Water Chemistry, Isotopes and Microbiology to Evaluate Groundwater Sources,

792      Flow Paths and Geochemical Reactions in the Death Valley Flow System, USA. *Procedia*

793      *Earth Planet. Sci.* 7, 842–845. doi:10.1016/j.proeps.2013.03.033.

794      Tully, J. G., Whitcomb, R. F., Clark, F. H., and Williamson, D. L. (1977). Pathogenic

795      mycoplasmas: Cultivation and vertebrate pathogenicity of a new spiroplasma. *Science (80-.*

796      *).* 195, 892–894. doi:10.1126/science.841314.

797      Wickham, H. (2016). ggplot2: Elegant Graphics for Data Analysis. Available at:

798      https://ggplot2.tidyverse.org.

799      Williams, T. A., Szöllosi, G. J., Spang, A., Foster, P. G., Heaps, S. E., Boussau, B., et al. (2017).

800      Integrative modeling of gene and genome evolution roots the archaeal tree of life. *Proc.*

801      *Natl. Acad. Sci. U. S. A.* 114, E4602–E4611. doi:10.1073/pnas.1618463114.

802      Woese, C. R. (2002). On the evolution of cells. *Proc. Natl. Acad. Sci. U. S. A.* 99, 8742–8747.

803      doi:10.1073/pnas.132266999.

804      Woyke, T., Xie, G., Copeland, A., González, J. M., Han, C., Kiss, H., et al. (2009). Assembling

805       the marine metagenome, one cell at a time. *PLoS One* 4. doi:10.1371/journal.pone.0005299.
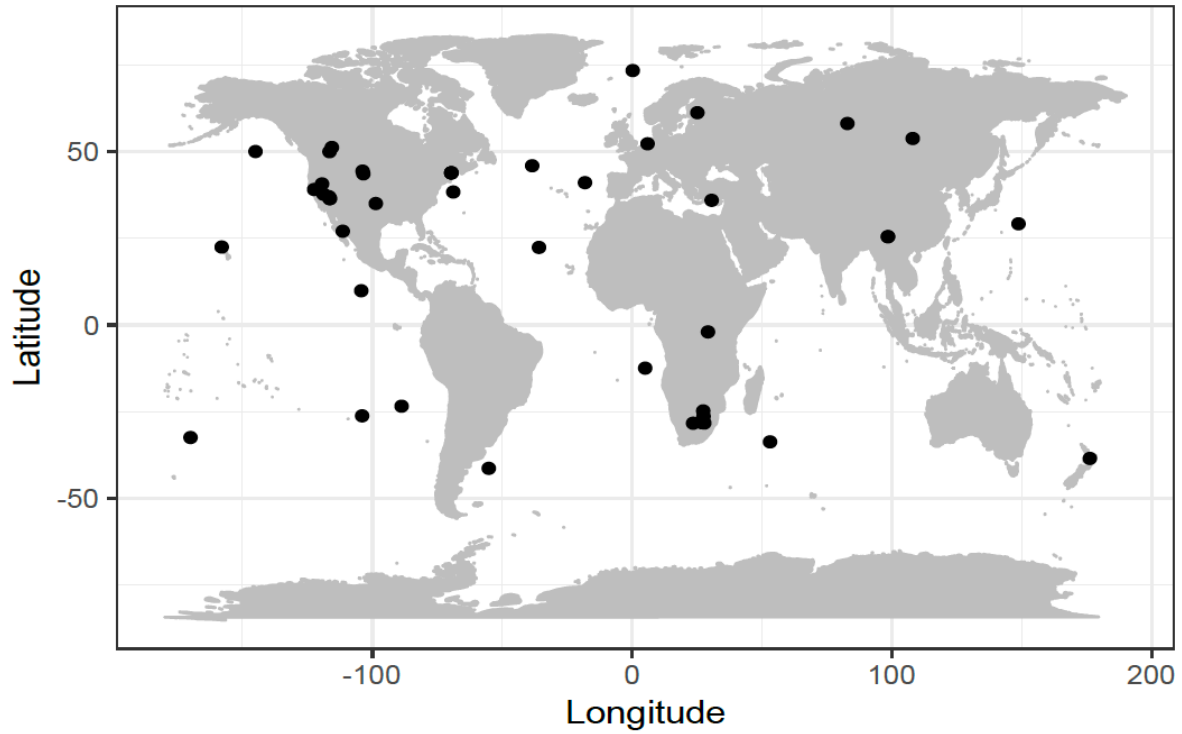
806    Wrighton, K. C., Thomas, B. C., Sharon, I., Miller, C. S., Castelle, C. J., VerBerkmoes, N. C., et

807       al. (2012). Fermentation, hydrogen, and sulfur metabolism in multiple uncultivated bacterial

808       phyla. *Science (80-. ).* 337, 1661–1665. doi:10.1126/science.1224041.

809    Yeoh, Y. K., Sekiguchi, Y., Parks, D. H., and Hugenholtz, P. (2016). Comparative genomics of

810       candidate phylum tm6 suggests that parasitism is widespread and ancestral in this lineage.

811       *Mol. Biol. Evol.* 33, 915–927. doi:10.1093/molbev/msv281.

812    Yu, F. B., Blainey, P. C., Schulz, F., Woyke, T., Horowitz, M. A., and Quake, S. R. (2017).

813       Microfluidic-based mini-metagenomics enables discovery of novel microbial lineages from

814       complex environmental samples. *Elife* 6, 1–20. doi:10.7554/eLife.26580.

815    Zengler, K., and Zaramela, L. S. (2018). The social network of microorganisms - How

816       auxotrophies shape complex communities. *Nat. Rev. Microbiol.* 16, 383–390.

817       doi:10.1038/s41579-018-0004-5.

818

**Figure 1.** Geographic locations of sample collection sites.

845
846 **Figure 2.** Maximum likelihood concatenated phylogenetic tree of single copy proteins (n=5)
847 from Bacteria (a) and Archaea (b). All SAGs from this study are highlighted red. Patescibacteria
848 and DPANN are highlighted with grey and labeled by individual proposed phyla within the
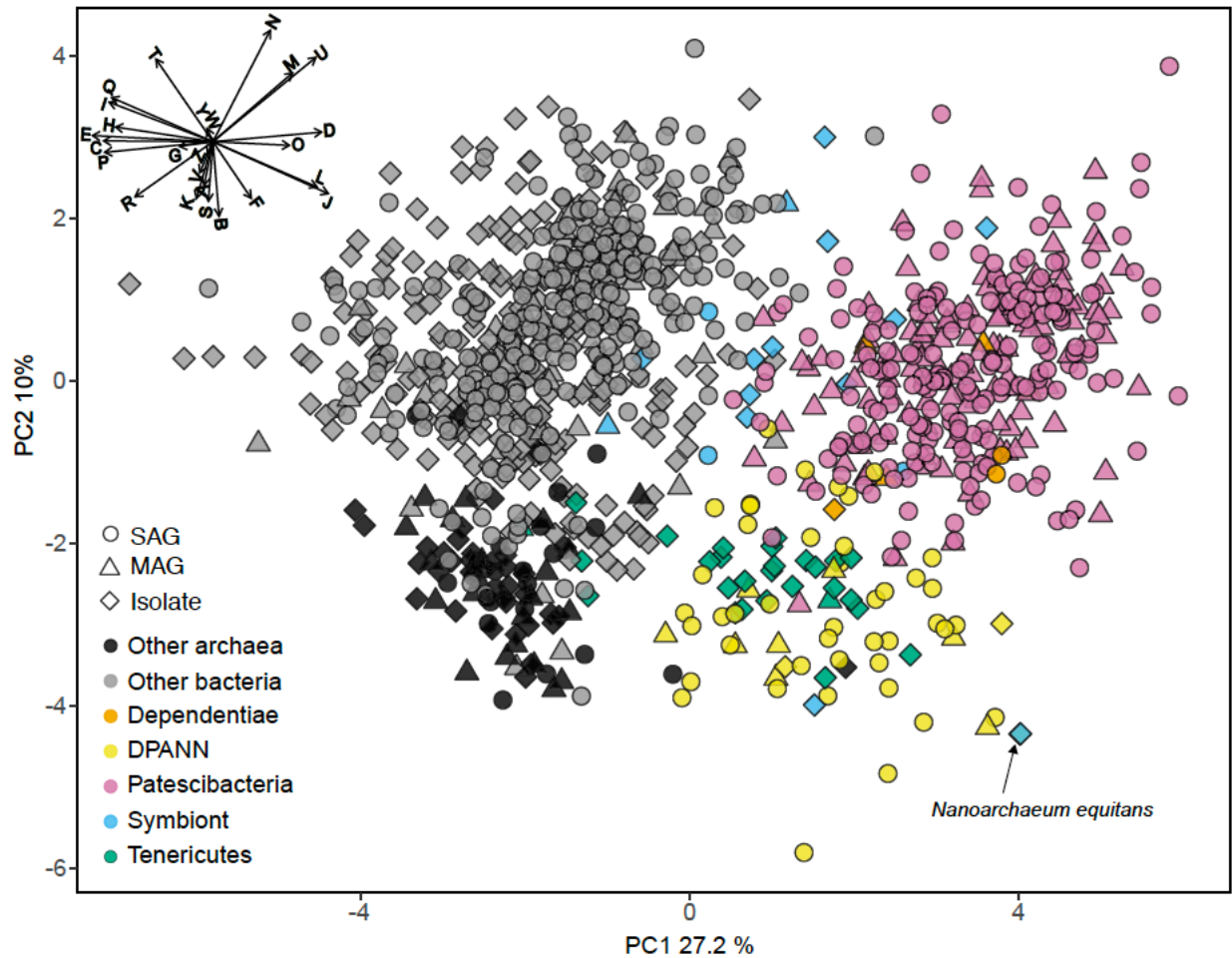849 superphylum (Rinke et al., 2013; Brown et al., 2015).

**Figure 3.** Relative abundance of Patescibacteria and DPANN from 34 geographically diverse samples determined from randomly sequenced LoCoS SAGs. The plate identifiers can be cross-referenced with specific SAGs and geographic sites in Table S1. The AG-274 sample contains small cells from a water-filled rock fracture at 1,340 m depth below surface in the Beatrix gold mine in South Africa.

**Figure 4.** Plot of the percent contribution of individual phyla to the Chi-square statistic from checkM (**A**) and 16S rRNA gene (**B**) co-sorting analyses. Classification of phyla in (**A**) from concatenated phylogenetic tree in Figure 2 (Table S1) and from 16S taxonomy in (**B**).

**Figure 5.** Phylum-resolved cell diameters. Solid black bars indicate medians; boxes represent the interquartile ranges (IQR) of the $1_{st}$ (Q1) and $3_{rd}$ (Q3) quartiles; whiskers denote the minimum (Q1 - 1.5*IQR) and maximum (Q3 + 1.5*IQR) values; outliers outside of the whiskers are marked by black dots. Orange indicates Bacteria and green indicates Archaea. A pairwise ranked-sum Wilcoxon test confirmed that the median diameter of Patescibacteria (highlighted in magenta) was smaller than most other phyla (27/36 phyla with p-values < 0.05; Table S5). The median diameter of DPANN (highlighted in yellow) was not significantly different from other archaea (1/36 phyla with p-values < 0.05; Table S5), likely due to the large variability in DPANN cell diameters. Individual cell diameters are available in Table S1 and pairwise p-values are located in Table S5.
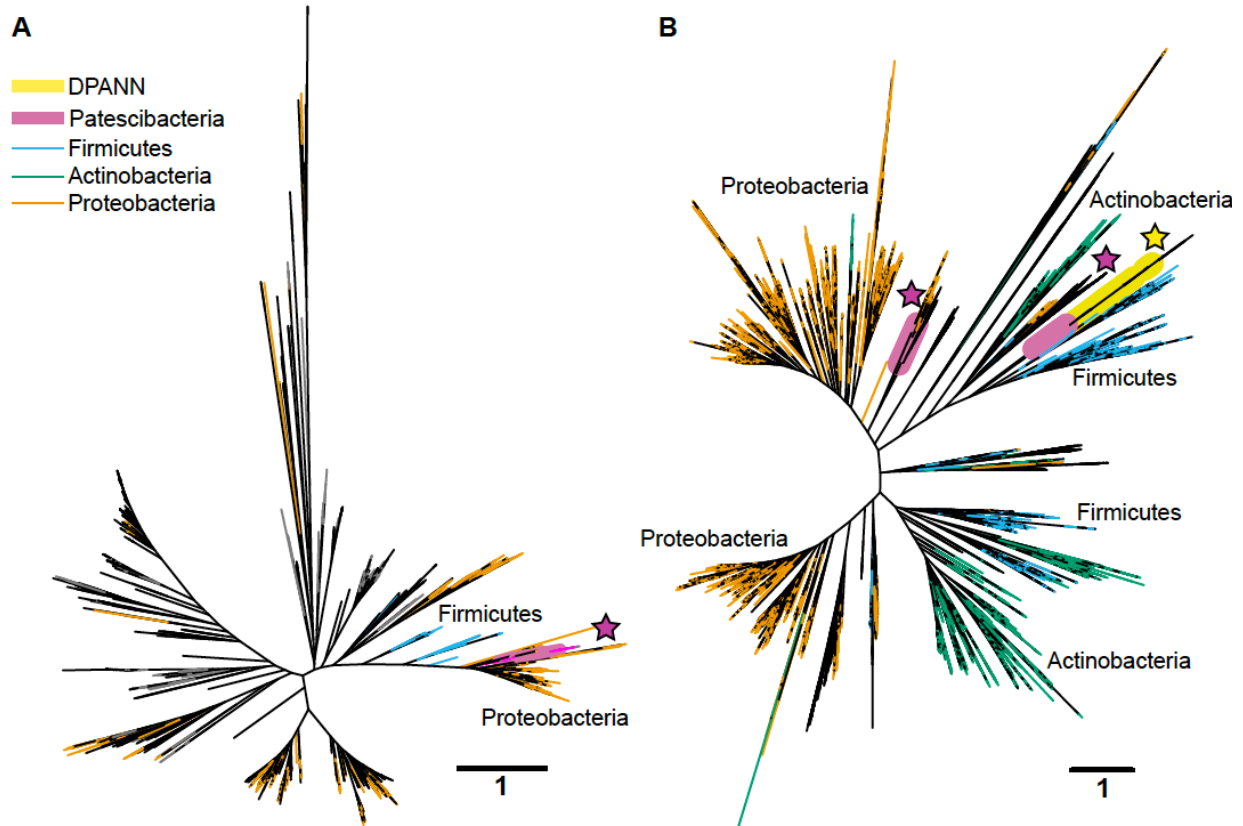
**Figure 6.** Principal components analysis (PCA) of the relative abundance of clusters of orthologous groups (COG) categories as the input variables. SAGs from this study (Table S1) and other studies (Table S2) with >30% completeness and had a near-full-length 16S rRNA gene and were included in the phylogenetic tree in Figure 8 (n=1,092). The vector plot in the upper left corner shows the COG categories that contributed to the most separation of the genomes: **Information Storage and Processing** Translation, ribosomal structure and biogenesis (J), RNA processing and modification (A), Transcription (K), Replication ,recombination and repair (L), Chromatin structure and dynamics (B); **Cellular Processes and Signaling** Cell cycle control, cell division, chromosome partitioning (D), Nuclear structure (Y), Defense mechanisms (V), Signal transduction mechanisms (T), Cell wall/membrane/envelope biogenesis (M), Cell motility (N), Cytoskeleton (Z), Extracellular structures (W), Intracellular trafficking, secretion, and vesicular transport (U), Posttranslational modification, protein turnover, chaperones (O); **Metabolism** Energy production and conversion (C), Carbohydrate transport and metabolism (G), Amino acid transport and metabolism (E), Nucleotide transport and metabolism (F), Coenzyme transport and metabolism (H), Lipid transport and metabolism (I), Inorganic ion transport and metabolism (P), Secondary metabolites biosynthesis, transport and catabolism (Q),; **Poorly Characterized** General function prediction only (R), Function unknown (S). SAG, single amplified genome; MAG, metagenome assembled genome. Symbiont genomes are listed in Table S4. Note position of *Nanoarchaeum equitans* with black arrow.
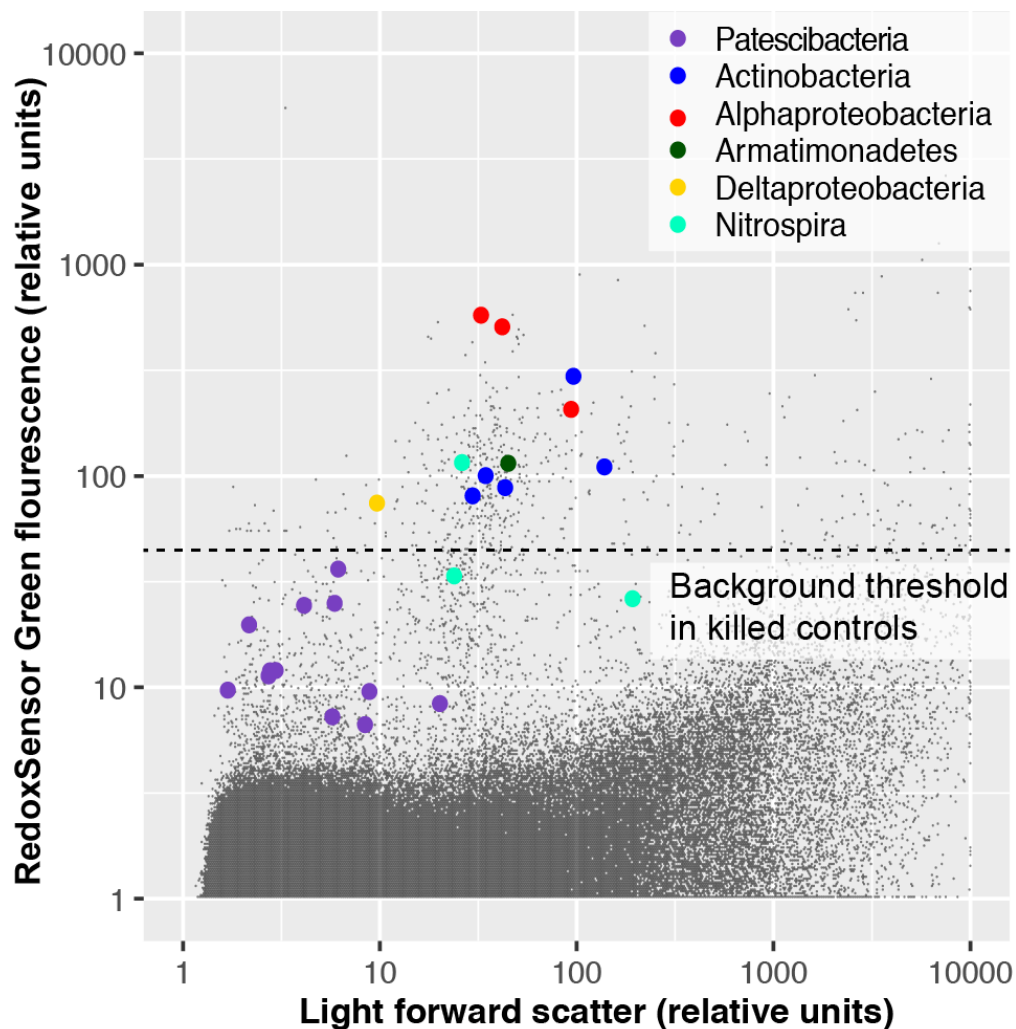
888  **Figure 7.** Relationship plots between estimated genome size, coding density, oxidoreductase
889  count, and cell diameter among SAGs (Table S1) and other genome sequences (Table S2) that
890  were greater than or equal to 30% complete, and were included in the 16S rRNA gene tree in
891  Figure 8. Symbiont genomes are listed in Table S4.

**Figure 8.** Maximum likelihood phylogeny of near full-length (>1,200 bp) 16S rRNA genes from Bacteria and Archaea, annotated with the distribution of counts for electron transport chain complexes, oxygen reductases, and oxidoreductase (Enzyme Commission 1; EC1) relative abundances from SAGs in this study (Table S1) and previously reported genome sequences (Table S2). The four innermost rings depict the counts of the electron transport chain complexes: I (NADH dehydrogenase subunits), II (succinate dehydrogenase subunits), III (cytochrome c reductase subunits), and IV (oxygen, nitrate, sulfate, iron, arsenate, and selenate reductase subunits). The outermost ring shows the relative abundance of oxidoreductases (Ox) for each genome assembly as a gradient from low (blue) to high (yellow). The peripheral stacked bar charts show the counts of oxygen reductases from both the heme copper oxidase and bd-ubiquinol oxidase oxygen reductase ($O_2$red) families grouped as high (orange) or low (sky blue) affinity for oxygen (note scale bar differences between bacterial and archaeal trees). Patescibacteria are highlighted in magenta and DPANN are highlighted in yellow. Other bacterial and archaeal phyla are highlighted in alternating white and grey.

44

**Figure 9.** Maximum likelihood phylogenetic trees of the oxygen-binding subunit I from the heme copper oxidase (HCO) type A (a) and the A subunit from the bd-ubiquinol (b) oxygen reductase families. Patescibacteria and DPANN sequences are marked with magenta and yellow stars, respectively. The Patescibacteria HCO type A sequences (a) are nested within a larger clade containing mostly Proteobacteria (orange), and the Patescibacteria and DPANN bd-ubiquinol sequences (b) are nested within Proteobacteria (orange) and Firmicutes (blue) dominated clades. The scale bar represents the estimated number of substitutions per site

913
914 **Figure 10.** Oxidoreductase activity in subsurface (~300 m below surface) microbial cells from
915 Homestake Mine (Lead, South Dakota, USA) measured by RedoxSensor Green (RSG;
916 ThermoFisher).
917
918
919
920
921
922
923
924
925
926
927
928
929
930

**Supplemental Table Captions**

**Table S1.** Deep-sequenced and LoCoS SAGs from this study with genomic statistics and associated environmental metadata. Data are ordered with the following column headers:
1=genome; single amplified genome (SAG identifier)
2=gold.analysis.id; Gold analysis identifier (used to search genome in IMG/M)
3=phylum
4=assembly.completeness; checkM completeness estimates
5= contamination; checkM estimated genome contamination
6=assembly.size; SAG assembly size
7=est.genome.size; estimated genome size
8=coding.density
9=ec1.count; counts of oxidoreductases from SAG assembly
10=est.ec1.count; estimated counts of oxidoreductases from predicted genome size of n Mbp
11=16s.copy.number; number of predicted 16S rRNA genes
12=cell.diameter; estimated cell diameter
13=sequencing.center
14=sample.collection.site; name of site where samples were collected
15=sample.type
16=date.collected; sample collection date dd/mm/yy
17=latitude
18=longitude
19=depth
20=dissolved.oxygen (micromoles/L)
21=ph
22=salinity (practical salinity units, psu)
23=temperature (degrees Celsius
24=h2s; dissolved hydrogen sulfide (millimoles/L)
NA=not applicable

**Table S2.** Genomes from other studies with associated genomic statistical information (accessed from IMG/M on April 2018).

**Table S3.** Potential co-sorted SAGs from deep-sequenced and LoCoS datasets. SAGs can be cross-referenced for specific information with Table S1.

**Table S4.** Symbiont genome assemblies and taxonomic names used in Figures 6 and 7.

**Table S5.** Pairwise Wilcoxon's test p-values on all phyla versus phyla cell diameter estimations in Figure 5.