

1 **Materials and Methods Supplement**

2

3 *Bin splitting, and additional genome quality assessment*

4

5 Two MAGs that classified as members of the Acidobacteria and Deinococcus-Thermus phyla
6 had high contamination estimates of 74% and 68%, respectively. After analyzing each MAG in
7 greater detail using both refineM [1] alongside manual TNF frequency and coverage plots
8 (Supplementary Figure 2a), we found that each could be split to create two additional MAGs.
9 The Acidobacteria MAG could be separated using read coverage and TNF into a low coverage
10 and high coverage MAG, after which the contamination was reduced to 5.8% and 3%. The
11 contaminated Deinococcus-Thermus MAG appeared to represent two distinct species, as the
12 refineM BLAST taxonomy assignments were largely split between *Thermus antranikianii* and
13 *Thermus islandicus*, corresponding to two different clouds on the TNF plot (Supplementary
14 Figure 2b). Separation of these sequences reduced the contamination estimates to less than 5%
15 for the two new MAGs. Furthermore, the TNF-based diversity of the *Thermus antranikianii*
16 MAG largely coincided with the diversity of the dominant Deinococcus-Thermus SAG
17 population (Supplementary Figure 2c), thus the *Thermus antranikianii* MAG was used for the
18 downstream Deinococcus-Thermus SAG to MAG comparisons.

19

20 *SSU rRNA collection and primer mismatch analysis*

21

22 16S rRNA genes were extracted from the assemblies using cmsearch of the Infernal package [2]
23 and the RF00177 covariance module of the 16S rRNA gene [3]. Only 16S rRNA genes greater

24 than 1 kb in length were used in subsequent analyses (n = 151 16S rRNA gene sequences). As
25 expected, lower quality genomes were less likely to contain a 16S gene (completeness estimate
26 average was 54% for genomes with a 16S gene, and 33% for those without). Once 16S rRNA
27 genes had been collected and filtered to 1 kb and above, primer binding was assessed with
28 PrimerProspector [4] using default parameters. Any sequence with an overall weighted score
29 greater than 1.0 suggested that the 16S rRNA gene would be missed by the tested primer set
30 (see Eloë-Fadrosh et al. 2016 for additional details [5]).

31

32 *Construction of concatenated marker gene phylogenies and 16S rRNA gene phylogeny*

33

34 Concatenated marker gene phylogenies were constructed by combining a dereplicated
35 reference set of genomes together with query genomes (e.g. Dewar Creek SAGs and MAGs).
36 Marker proteins were extracted from each genome using hmmsearch (version 3.1b2) and
37 alignments were constructed with MAFFT [6] using the mafft-linsi option. Alignments were
38 trimmed with trimAl 1.4 [7], removing sites when more than 90% of taxa contained a gap. For
39 the 16 ribosomal protein tree and 56 marker gene tree, genomes were removed if they
40 contained less than 50% of the markers in the set. The presence of all 3 subunits for the RNA
41 polymerase gene were required for a genome to be included in the 3 subunit RNA polymerase
42 phylogeny. Individual protein alignments were then concatenated to produce an alignment of
43 51,239 sites. Maximum likelihood phylogenies were constructed with IQ-TREE [8], using the
44 WAG substitution model and 1,000 bootstraps. The set of reference genomes was collected by
45 dereplicating the full set of IMG (Integrated Microbial Genomes) isolate genomes (64,005

46 genomes) [9] based on cd-hit [10] clustering of the RNA polymerase gene (*rpoB*) at 65%. This
47 produced a dereplicated set of unique family-level genomes, spanning all bacteria and archaea
48 (n=681).

49

50 Lineage-specific trees (Figure 5a and Supplemental Figure 4) were constructed in a similar
51 manner, however, only the UNI56 marker set was used. Outgroups were selected as the
52 nearest neighbor taxa from the full UNI56 archaea/bacteria tree. For these trees, the full set of
53 genomes were collected for each phylum from IMG/M (Integrated Microbial Genomes /
54 Metagenomes) [11], then dereplicated using the *rpoB* gene at different clustering levels,
55 ranging from 90 to 100%. The clustering level was varied by clade in order to produce roughly
56 50 references per phylum. The reference set for the Crenarchaeota was dereplicated at 80%
57 *rpoB* similarity, as this was a broader phylogenetic clade than the other sets of lineage specific
58 trees. The query genomes (Dewar creek SAGs and MAGs) were dereplicated at 100% RNA
59 Polymerase beta-subunit gene identity. The lineage-specific trees were constructed in the same
60 manner as outlined above.

61

62 The reference set for the 16S rRNA gene phylogeny was based on sequences extracted from the
63 681 reference genomes used in the multi-marker gene trees. Dewar Creek SAG and MAG query
64 sequences only included 16S rRNA gene sequences that were greater than 1kb in length (n =
65 151). The combined set of query and reference 16S rRNA genes was aligned using cmalign using
66 the -matchonly option, resulting in an alignment length of 1534 bp, and the tree was

67 constructed with IQ-TREE [8] under the general time-reversible evolutionary model with 1,000
68 bootstraps. All trees (16S rRNA gene and protein markers) were visualized with ggtree [12] in R.

69

70 *Relative abundance comparison between amplicon, SAG and MAG datasets*

71

72 Since community composition was compared across SAGs, MAGs and amplicon datasets, the
73 generation of abundance profiles from the three distinct approaches should be briefly
74 described. Relative abundances of amplicon groups were the result of 97% OTU clustering,
75 taxonomic assignment, and grouping at the phylum level. The SAG abundances were
76 straightforward, as taxon assignments were based on the UNI56 maker gene tree, then counts
77 were based on phylum level assignments. Taxonomy assignments of MAGs were also based on
78 the UNI56 marker gene tree and abundances were based on read mapping where reads from
79 the bulk metagenome were mapped to the collective set of MAGs using bbsplit from the
80 bbtools package [13], where a read could only be mapped once.

81

82 *Note on phylum level classifications*

83

84 *Candidatus* Kryptonia is described as a phylum within NCBI and the corresponding publication
85 [5] while the GTDB-Tk [14] places Kryptonia within the Bacteroidetes phylum, and *Candidatus*
86 Parcubacteria [15] within the Patescibacteria phylum. For the current work, we are using the
87 names from the original publications as both classifications are based on the commonly cited
88 16S rRNA phylum designations [16] and a concatenated ribosomal tree in the case of the

89 Patescibacteria [15], though we acknowledge the new names in GTDB-tk and note that specific
90 phylum names are not of critical importance to this manuscript as much of the study focusses
91 on comparisons between SAGs and paired MAGs, and further dissection within and between
92 dominant populations.

93

94 *Pairwise average nucleotide identity (ANI) and definition of species level clusters used in*
95 *downstream analyses*

96

97 Pairwise genomic ANI analysis was performed with fastANI [17]. Genome pairs were filtered to
98 include only those pairs with an alignment fraction $\geq 70\%$, which were then grouped into
99 clusters sharing $\geq 95\%$ ANI using mcl [18]. Pairwise ANIs $\geq 95\%$ were used to define species level
100 clusters [19]. These ANI clusters were used for downstream intra-species analyses including
101 gene family/orthologue clustering and population analyses.

102

103 *Gene annotations, gene content comparisons and gene family diversity assessment*

104

105 Genes were called and annotated using the Integrated Microbial Genomes (IMG) [11]. The
106 naming of contigs and genes followed JGI's in house nomenclature and can be cross-referenced
107 with the IMG webserver (img.jgi.doe.gov). For gene content comparisons, annotations of
108 individual genes were used in combination with gene family clustering using OrthoFinder 2.1.3
109 [20], and the ANI/species level genome collections as input. In addition to clustering genomes
110 into 95% ANI groups, completion cutoffs of 40% were used for ortholog clustering.

111

112 *Identification of SNPs*

113

114 In preparation for SNP calling, the highest quality SAG (SAG with the highest completeness
115 estimate and contamination less than 5%) from each 95% ANI group was identified and used as
116 the reference genome. Reads of all SAGs were mapped to the references and SNPs were called
117 using the MIDAS pipeline [21]. Briefly SNP calling was done by mapping all reads belonging to
118 genomes within an ANI group using bowtie2 (--very-sensitive, global alignment mode) and
119 reads with less than 95% similarity to the reference, average read quality of less than 30,
120 mapping quality less than 20, and base quality scores of less than 30 were discarded. For a SNP
121 to be counted, it had to have a minor allele frequency (MAF) of at least 10%.

122

123 *Whole genome phylogenies of the dominant populations*

124

125 Population level phylogenies based on the variant sites between SAGs were created by taking a
126 collection of within species genomes, identifying variant sites using NUCmer from the MUMmer
127 package [22], then producing a neighbor joining tree. Identification of strain level clusters was
128 performed via RhierBAPS [23] using the same NUCmer whole genome multiple sequence
129 alignment as input, by partitioning each genome sequence into the appropriate cluster based
130 on the allele frequencies within each cluster.

131

132 *Estimates of recombination*

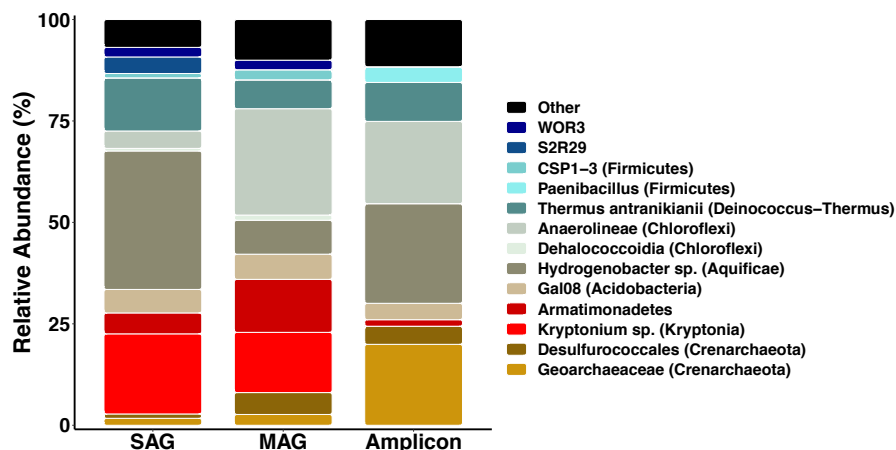
133

134 We generated SNP linkage disequilibrium (LD) profiles for the *Hydrogenobacter sp.*, *Kryptonium*
135 *sp.*, and *Thermus antranikianii* lineages, where the MIDAS constructed SNP depth and
136 frequency tables were used as input and converted to a SNP pair correlation matrix (R^2). LD
137 plots were created by reading in the SNP pair correlation matrix, creating a table of R^2 values by
138 distance, then plotting in ggplot2 [24]. The number of SNPs per kb was also calculated on a per
139 gene basis and mapped to their corresponding annotations (COG database used in figures).
140 SNPs per kb were calculated using the MIDAS script SNP_diversity.py
141 (<https://github.com/snayfach/MIDAS>) [21].

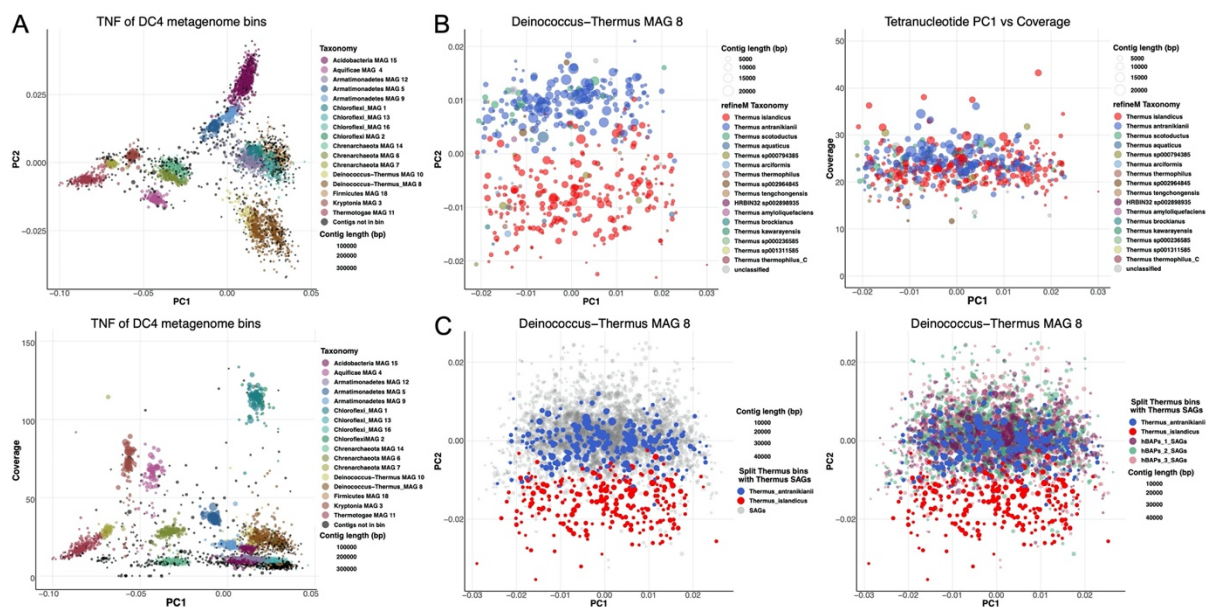
142

143

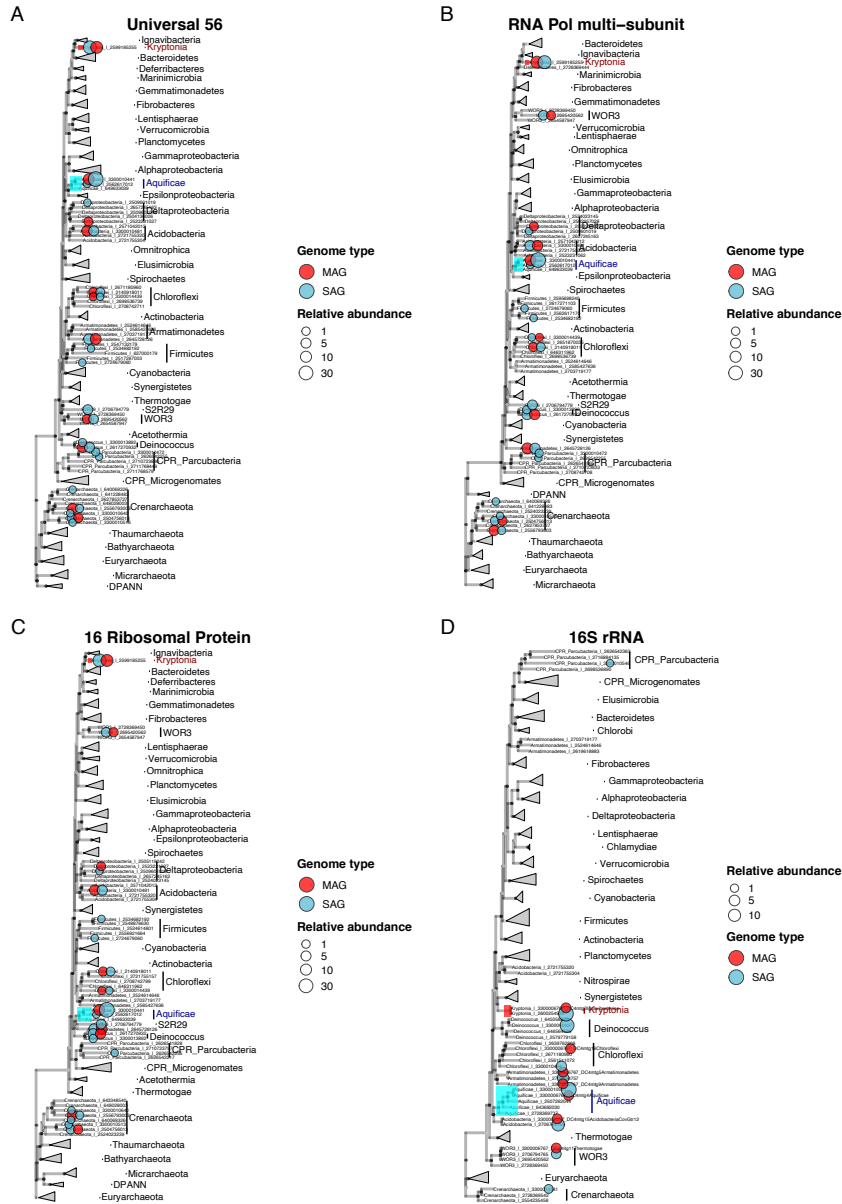
144 **Supplementary Figures and Tables**



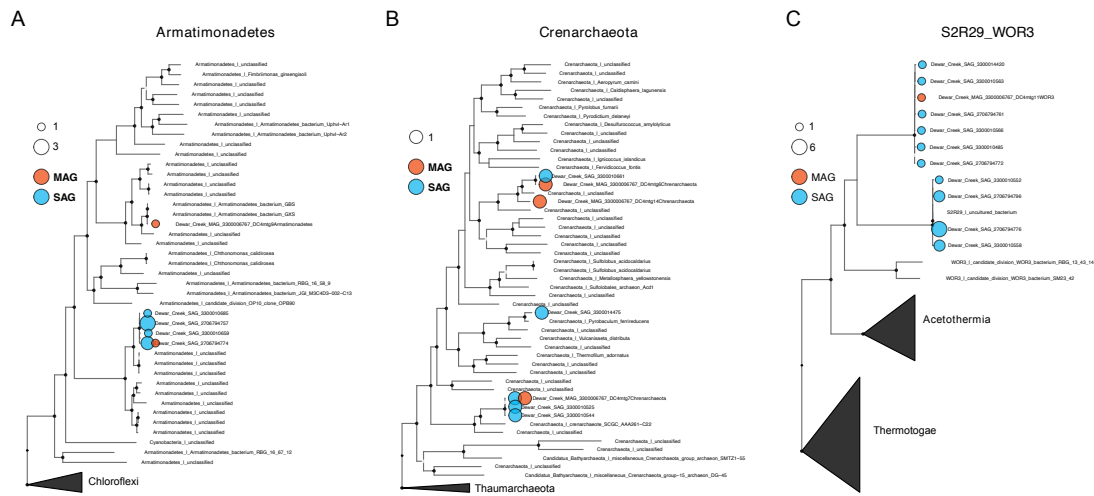
145
146 **Supplemental Figure 1.** The community composition of the single Dewar Creek sediment
147 sample using the three approaches. This figure is similar to Figure 1b, but with more
148 resolved taxonomic assignments.



149
150 **Supplemental Figure 2.** Example of bulk metagenome bin cleaning. **a** Tetranucleotide
151 frequency plot of all bins extracted from the bulk metagenome. All but one of the bins were
152 either high or medium quality MAGs based on the MIMAG/MISAG standards [25]. **b**
153 Demonstration of bin cleaning using the original highly contaminated Deinococcus-Thermus
154 MAG 8, colored by contig taxonomy assignments. The majority of contigs classified to either
155 *Thermus antranikianii* or *Thermus islandicus*. **b, left** shows two clear TNF compositional clouds
156 and **b, right** shows that these bins could not be separated based on coverage and PC1 alone. **c**
157 demonstrates that most Deinococcus-Thermus SAGs are most similar to the *Thermus*
158 *antranikianii* MAG. **Left** MAG contigs combined with SAG contigs. **Right** shows the same plot,
159 but where SAGs are colored by the HBAPS population clusters assigned to the *Thermus*
160 *antranikianii* SAGs.

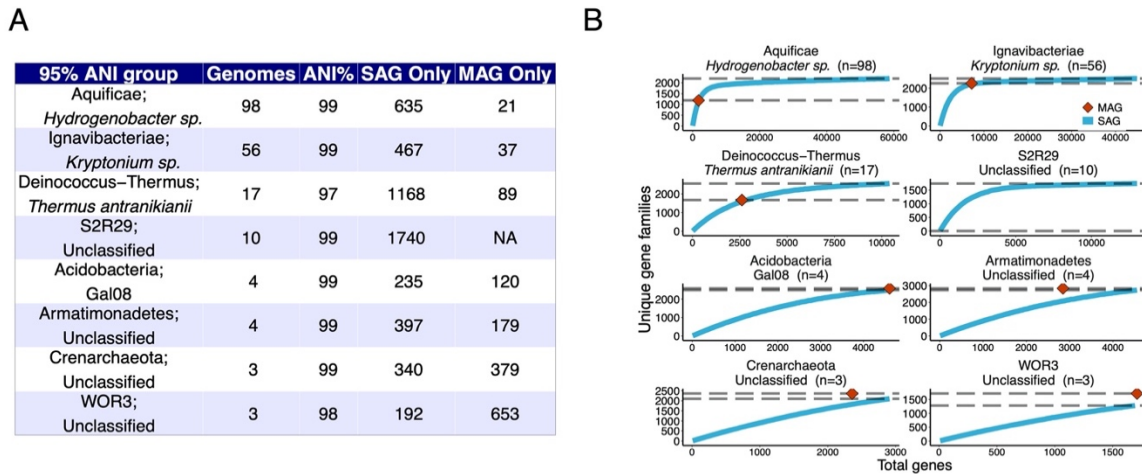


161
 162 **Supplementary Figure 3.** Phylogenetic trees of the Dewar Creek bacteria and archaea using
 163 different phylogenetic markers and marker sets including **a** the 56 universal markers used in
 164 main manuscript figure (Figure 2a), **b** the three subunits of the RNA Polymerase gene, **c** a set of
 165 16 conserved ribosomal proteins, and **d** the 16S rRNA genes derived from the same set of
 166 genomes used in the multi-protein phylogenies. Abundance counts within concatenated protein
 167 phylogenies (**a – c**) represent relative proportions within SAG and MAG datasets where MAG
 168 relative abundances are the result of bulk metagenome reads mapped to each MAG. Abundance
 169 counts within the 16S rRNA gene phylogeny are the result of 16S clustering at 87.5 % similarity
 170 (family level based on Yarza 16S rRNA standards [16]). Note: if any Dewar Creek lineages are
 171 missing from one phylum within one marker set, but present when using a different marker set,
 172 this means the genome either did not contain the markers or did not have enough markers to
 173 remain in the tree after quality filtering (cutoffs for the universal 56 marker genes and 16
 174 ribosomal proteins were set to contain at least 50% of the marker set).



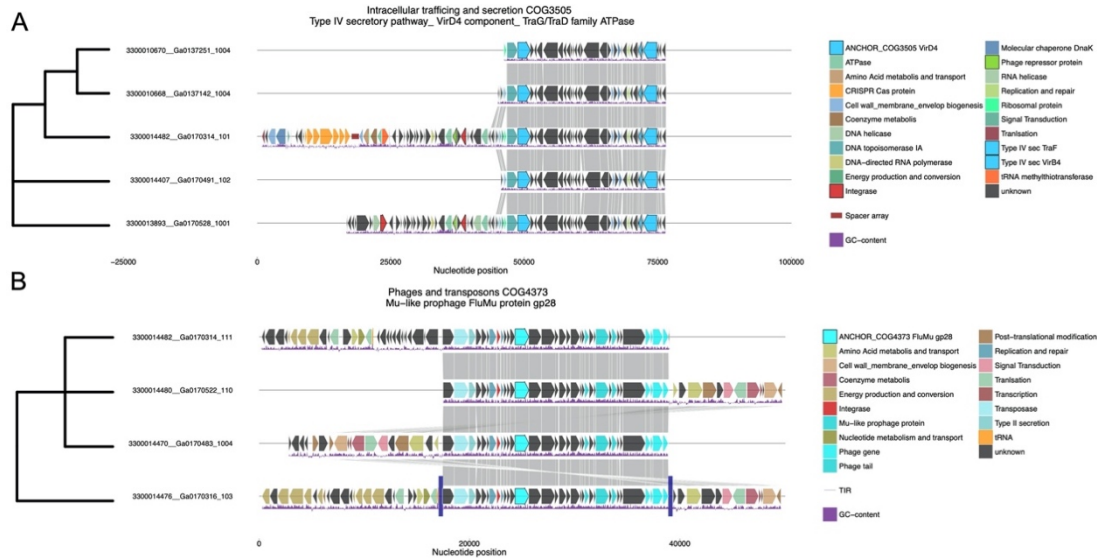
176
177
178
179
180
181
182

Supplemental Figure 4. Genome phylogenies of the abundant taxa excluding the Aquificae, *Candidatus* Kryptonia, and Deinococcus-Thermus phylogenies as these were shown in Figure 5a. Members of the **a** Armatimonadetes, **b** Crenarchaeota, and **c** S2R29 and WOR3 candidate phyla are displayed. The sizes of the Dewar Creek genome bubbles are based on 100% *rpoB* clustering. For example, when a bubble is larger than a value of 1, this means there are two or more identical genomes in that cluster based on the *rpoB* gene.

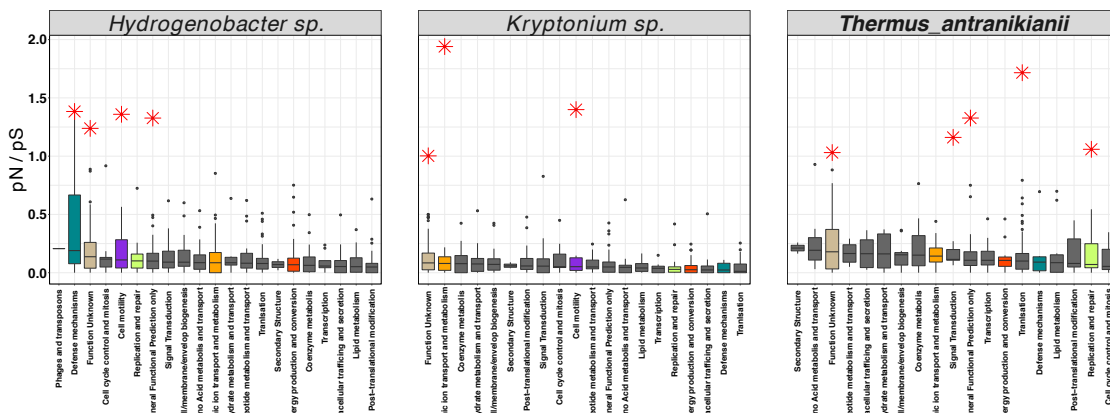


183
184
185
186
187
188
189
190

Supplementary Figure 5. **a** Table noting the average nucleotide identity (ANI) and the number of unique gene families observed in the SAG only or MAG only genomic subsets for each lineage. **b** Rarefaction curves highlighting the gene family diversity for SAGs (blue curves) and MAGs (red diamond).



191
 192 **Supplementary Figure 6. a** Observed synteny between 5 putative ICE plasmid sequences
 193 identified from two *Hydrogenobacter sp.* genomes. The “anchor” gene codes for a
 194 TraG/TraD ATPase, involved in T4SS transport. Sequences, 3300014482_Ga0170314_101
 195 and 3300013893_Ga0170528_1001 both contain integrases, phage repressor proteins,
 196 and the first has a tRNA-*Ala* adjacent to the integrase, a potential host integration site. **b**
 197 Synteny between 4 putative phage / prophage sequences. Note the tRNA, putative
 198 integration site adjacent to one of two Terminal inverted repeat sequences, designated by
 199 the blue vertical bars. The “ANCHOR” gene was the gene used to center both plots.
 200



201
 202 **Supplementary Figure 6.** pN/pS boxplots grouped by COG category. COG categories of interest
 203 are colored, and genes with pN/pS > 1 are noted by the red stars, as these are genes that may
 204 be under selection.

Population	Reference Completeness (%)	Genome size (Mb)	SNPs	% Polymorphic	NonSyn/kb	Syn/kb	Nuc div
<i>Hydrogenobacter sp.</i>	99	1.7	22810	1.36	4.4	42	2.0
<i>Kryptonium sp.</i>	96	2.6	23866	0.92	1.7	26	1.5
<i>Thermus antranikianii</i>	71	1.4	8571	0.61	7.2	50	1.8

205
 206 **Supplementary Table 1.** SNP statistics for each of the three analyzed populations.
 207
 208

209 **References**

- 210 1. Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, et al. Recovery
211 of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life.
212 *Nat Microbiol* 2017; 2: 1533–1542.
- 213 2. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches.
214 *Bioinformatics* 2013; 29: 2933.
- 215 3. Nawrocki EP, Burge SW, Bateman A, Daub J, Eberhardt RY, Eddy SR, et al. Rfam 12.0:
216 updates to the RNA families database. *Nucleic Acids Res* 2015; 43: D130.
- 217 4. Walters WA, Caporaso JG, Lauber CL, Berg-Lyons D, Fierer N, Knight R. PrimerProspector:
218 de novo design and taxonomic analysis of barcoded polymerase chain reaction primers.
219 *Bioinformatics* 2011; 27: 1159.
- 220 5. Eloe-Fadrosh EA, Ivanova NN, Woyke T, Kyrpides NC. Metagenomics uncovers gaps in
221 amplicon-based detection of microbial diversity. *Nat Microbiol* 2016; 1: 15032.
- 222 6. Katoh K, Standley DM. MAFFT Multiple Sequence Alignment Software Version 7:
223 Improvements in Performance and Usability. *Mol Biol Evol* 2013; 30: 772.
- 224 7. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. trimAl: a tool for automated
225 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 2009; 25: 1972.
- 226 8. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective
227 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol* 2015;
228 32: 268–74.
- 229 9. Markowitz VM, Chen I-MA, Chu K, Szeto E, Palaniappan K, Pillay M, et al. IMG/M 4
230 version of the integrated metagenome comparative analysis system. *Nucleic Acids Res*
231 2014; 42: D568–73.
- 232 10. Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein
233 or nucleotide sequences. *Bioinformatics* 2006; 22: 1658–1659.
- 234 11. Chen IMA, Chu K, Palaniappan K, Ratner A, Huang J, Huntemann M, et al. The IMG/M
235 data management and analysis system v.6.0: New tools and advanced capabilities.
236 *Nucleic Acids Res* 2021; 49: D751–D763.
- 237 12. Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. ggtree : an r package for visualization and
238 annotation of phylogenetic trees with their covariates and other associated data.
239 *Methods Ecol Evol* 2017; 8: 28–36.
- 240 13. Bushnell B, Rood J, Singer E. BBMerge – Accurate paired shotgun read merging via
241 overlap. *PLoS One* 2017; 12: e0185056.
- 242 14. Chaumeil PA, Mussig AJ, Hugenholtz P, Parks DH. GTDB-Tk: A toolkit to classify genomes
243 with the genome taxonomy database. *Bioinformatics* 2020; 36: 1925–1927.
- 244 15. Brown CT, Hug LA, Thomas BC, Sharon I, Castelle CJ, Singh A, et al. Unusual biology across
245 a group comprising more than 15% of domain Bacteria. *Nature* 2015; advance on.
- 246 16. Yarza P, Yilmaz P, Pruesse E, Oliver Glöckner F, Ludwig W, Schleifer K-H, et al. Uniting the
247 classification of cultured and uncultured bacteria and archaea using 16S rRNA gene
248 sequences. *Nat Publ Gr* 2014; 12.
- 249 17. Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High-throughput ANI
250 Analysis of 90K Prokaryotic Genomes Reveals Clear Species Boundaries. *bioRxiv* 2017;
251 225342.
- 252 18. Enright AJ, Dongen S Van, Ouzounis CA. An efficient algorithm for large-scale detection of

- 253 protein families. *Nucleic Acids Res* 2002; 30: 1575.
- 254 19. Klappenbach JA, Goris J, Vandamme P, Coenye T, Konstantinidis KT, Tiedje JM. DNA–DNA
255 hybridization values and their relationship to whole-genome sequence similarities. *Int J*
256 *Syst Evol Microbiol* 2007; 57: 81–91.
- 257 20. Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome
258 comparisons dramatically improves orthogroup inference accuracy. *Genome Biol* 2015;
259 16: 157.
- 260 21. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. An integrated metagenomics
261 pipeline for strain profiling reveals novel patterns of bacterial transmission and
262 biogeography. *Genome Res* 2016; 26: 1612–1625.
- 263 22. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and
264 open software for comparing large genomes. *Genome Biol* 2004; 5.
- 265 23. Tonkin-Hill G, Lees JA, Bentley SD, Frost SDW, Corander J. RhierBAPs: An R
266 implementation of the population clustering algorithm hierbaps [version 1; referees: 2
267 approved]. *Wellcome Open Res* 2018; 3.
- 268 24. Wickham H. ggplot2: Elegant Graphics for Data Analysis. 2016. Springer-Verlag New York.
- 269 25. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK, et al.
270 Minimum information about a single amplified genome (MISAG) and a metagenome-
271 assembled genome (MIMAG) of bacteria and archaea. *Nat Biotechnol* 2017; 35: 725–731.
272