1 **Supplementary Materials for Magnuson et al. "Active lithoautotrophic and methane-**

2 **oxidizing microbial community in an anoxic, sub-zero, and hypersaline High Arctic**

3 **spring".**

4

5 **Extended Materials and Methods**

6 i. Site description and sample collection

7 Lost Hammer (LH) spring discharges through a precipitated mineral salt tufa as described in

8 previous publications (1-5) (Figure S1). LH emits gases composed of methane (50%), nitrogen

9 (35%), carbon dioxide (10%), and trace hydrogen, helium, and short-chain alkanes (1). The

10 spring sediments and water contain high concentrations of sulfate (100 000 mg/kg) as well as

11 sulfide (<50 mg/kg), ammonia (2.55 mg/kg), nitrate/nitrite (2.87 mg/kg), and iron (13 000

12 mg/kg) (1). Physical and geochemical parameters in the outlet (Table S1) have remained highly

13 stable since 2005, allowing comparison among samples collected in different years.

14

15 For this study, sediment samples (top ~10 cm) were collected in July 2017 and July 2019 with an

16 ethanol-sterilized scoop, and stored in sterile Falcon tubes filled to maximum to avoid aerobic

17 headspace. Sediment from July 2017 was used for metagenomic sequencing, and sediment from

18 July 2019 was used for RNA and SAG sequencing. In parallel, sediment for RNA extractions

19 was mixed with Zymo Research DNA/RNA Shield (Irvine, CA, USA). Samples were kept at

20 <5°C during transportation to Montreal, after which they were stored at -20°C for DNA and

21 RNA extraction and at -5°C (unfrozen) for SAG sequencing. Physical and geochemical

22 parameters of the overlying water were measured *in situ* with a YSI Professional Plus

23 Multiparameter instrument (Yellow Springs, OH, USA) and a PyroScience Piccolo2 oxygen

24  meter (Aachen, Germany). Ortho-phosphate and ammonia were measured *in situ* with

25  CHEMetrics Inc. (Midland, VA, USA) test kits.

26

27  ii. DNA extraction, metagenome sequencing, and metagenome data analyses

28  DNA was extracted from two 5 g portions of sediment with a Qiagen DNeasy PowerMax Soil

29  Kit (Hilden, Germany). The resulting DNA from each sediment sample was concentrated with a

30  Thermo Fisher Scientific SpeedVac Vacuum Concentrator (Waltham, MA, USA) and sequenced

31  on a HiSeq2500 (2x126 base reads) (Illumina, San Diego, CA, USA) at The Centre for Applied

32  Genomics (Toronto, ON, Canada). Low-quality reads and bases were trimmed with

33  Trimmomatic (v0.38, settings LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15) (6).

34  Remaining reads were classified with Kaiju (v1.7.3, default settings, nr_euk database) (7) and

35  phyloFlash (v3.4) (8). The reads from the two sediment samples were co-assembled with

36  Megahit (v1.1.3, setting meta-sensitive) (9) as well as assembled separately with metaSPAdes

37  (v3.13.0, default settings) (10). Reads were mapped to each assembly with BBMap (v38.26,

38  minid=0.95) and contigs longer than 5000 bp were binned with MetaBAT (v2.12.1) (11). Bin

39  completeness and contamination was estimated with CheckM (v1.0.12) (12). The Megahit co-

40  assembly and resulting bins were selected for downstream analysis based on sequencing statistics

41  as determined by metaQUAST (v5.0.1) (13) and number of high- and medium-quality bins.

42  Sequencing and assembly statistics can be found in Tables S2 and S3. The metagenome was

43  annotated with the Joint Genome Institute's IMG/M system using KEGG, COG, and pfam

44  databases (14, 15). Bins were classified with the Genome Taxonomy Database Toolkit (v1.3.0,

45  reference data R05-RS95) (16). Phylogenomic trees were constructed with Anvi'o (v6.2) using

46  the Bacteria_71 collection of single copy genes (17). Amino acid sequences for all genes were

47    concatenated, with a total alignment length of 24 451 bp, and approximately-maximum-

48    likelihood trees were constructed in FastTree with default settings within Anvi'o, with midpoint

49    rooting. Additional analyses were as follows: FeGenie was used to identify iron-related genes

50    (18); *amoA* and *pmoA* were distinguished by HMMer using FunGene Hidden Markov Models

51    (19); hydrogenases were classified with hydDB (20). Reductive and oxidative DsrAB were

52    classified as follows: DsrAB amino acid sequences were aligned against reference sequences

53    from Muller et al. (21) using MUSCLE with default settings (22). Maximum likelihood

54    phylogenetic trees were constructed in CLC Genomics Workbench (v. 12.0.3) using the WAG

55    protein substitution model and 1000 bootstraps (Figure S9). DsrAB were classified as reductive

56    or oxidative based on phylogenetic clustering and the presence of accessory proteins as in

57    Anantharaman et al. (23).

58

59    iii. Single cell sorting, genome amplification, sequencing, and data analyses

60    Sediment was shipped on ice to the Single Cell Genomics Center at the Bigelow Laboratory for

61    Ocean Sciences (East Boothbay, ME, USA) for SYTO-9 fluorescence-activated single-cell

62    sorting (FACS), DNA extraction, and genome amplification with WGA-X (as described in

63    Stepanauskas et al. (24)). The 16S rRNA gene was PCR amplified from the genomes (primers

64    27F (25) and 1492R (26)) and sequenced at the Centre de Recherche at Université Laval

65    (Quebec City, QC, Canada) on an Applied Biosystems 3730xl DNA Analyzer (Foster City, CA,

66    USA). Obtained 16S rRNA sequences were annotated by BLAST with the SILVA rRNA

67    database (release 138) (27). For whole genome sequencing, libraries were prepared with a

68    Nextera XT DNA Library Prep Kit and sequenced on a MiSeq (Illumina) with MiSeq Reagent

69    Kit v3 (600 cycles, 2x300 base reads). Low quality reads and bases were trimmed with BBDuk

70    (minimum Phred quality score 15, minimum length 30 bp) and contaminant human reads were

71    removed with DeconSeq (28). Genomes were assembled with SPAdes (v3.13.1, settings --sc --

72    careful) (29) and screened for contamination using JGI's Kmer Frequency Analysis tool and by

73    read classification with Kaiju (v1.7.3) (7). Average nucleotide identity of the SAGs against other

74    SAGs and metagenome bins was calculated with FastANI (v1.3) (30). Genome annotation was

75    done as described for the metagenome.

76

77    iv. mRNA extraction, sequencing, and analysis

78    RNA was extracted in triplicate with a Zymo Research ZymoBIOMICS DNA/RNA Miniprep

79    Kit from approximately 3 g sediment per extraction. Extracted samples were then treated with

80    the Invitrogen Turbo DNA-free kit (Carlsbad, CA, USA) to remove contaminating DNA. The

81    treated samples were then pooled and concentrated with a New England Biolabs Monarch RNA

82    Cleanup Kit (Ipswitch, MA, USA). Ribosomal RNA was depleted with a New England BioLabs

83    NEBNext rRNA Depletion Kit (Bacteria) and a sequencing library was prepared with a New

84    England BioLabs Ultra II RNA Library Prep Kit. The generated cDNA library was sequenced at

85    The Center for Applied Genomics at the Hospital for Sick Children (Toronto, ON, Canada) on a

86    NovaSeq 6000 (Illumina) with an S Prime 100-cycle flow cell (2x100 base reads). Reads were

87    trimmed with Trimmomatic (settings LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15

88    MINLEN 50) and rRNA reads were removed with SortMeRNA (v4.2.0) (31). Human

89    contaminant reads were removed with BBMap using the BBTools RemoveHuman masked

90    reference genome. Reads were classified with Kaiju (v1.7.3), and reads classified within

91    common sequencing contaminant groups as in Sheik et al. (2018) (32) were removed (~81,500

92    reads removed representing 77 genera; the genera with the most reads removed were

93    *Pseudomonas* (~12,000 reads), *Streptococcus* (~7000 reads), and *Ralstonia* (~6000 reads)).

94    Remaining reads were aligned to metagenome contigs and SAG scaffolds with bowtie2

95    (v2.3.5.1, with setting -very-sensitive-local) (33). Reads aligned to CDS regions were counted

96    with HTSeq (v0.12.4) (34) and transcripts per million reads (tpm) was calculated to normalize

97    expression values for each gene. For comparison of stress response genes in SAGs to those in

98    related genomes, protein-coding genes from reference genomes were queried against LH SAG

99    protein-coding genes with BLASTp using a cut-off e-value of 1e-15.

100

101   v. Sequencing data availability

102   Sequencing reads, metagenome, MAGs, and SAGs were deposited in NCBI under BioProject

103   PRJNA699472. JGI annotations of the metagenome and SAGs are available under GOLD Study

104   ID Gs0135943. SAGs were also deposited individually in JGI (Table S4).

105

106   vi. Gibbs free energy calculations

107   Gibbs free energy values were calculated using the method reported in Jones et al. (2018) (35)

108   with the following parameters based on measurements reported in Table S1: 278.1° K, pH 6.4,

109   7.39 M ionic strength, $1.02 \times 10^{-7}$ M $CH_4$, $3.81 \times 10^{-4}$ M $NH_3$, $8.61 \times 10^{-6}$ M $H_2S$, $1.57 \times 10^{-5}$ M

110   $O_2$, $2.31 \times 10^{-5}$ M $NO_3^-$, 1.04 M $SO_4^{2-}$, $4.11 \times 10^{-5}$ M $CO_2$, $3.16 \times 10^{-7}$ M $H^+$, $3.12 \times 10^{-5}$ M $NO_2^-$,

111   and 0.23 M $Fe^{3+}$.

112

113   **Supplementary Results and Discussion**

114   *Discrepancy between recovered MAGs and SAGs*

115     Using the criteria of average nucleotide identity >95%, only one SAG was found to correspond

116     to MAGs. Typically, more overlap between the two datasets would be expected, particularly for

117     those taxa abundant in the microbial community (36); conversely, the SAGs appear to be

118     enriched for taxa at low abundance in the metagenome. In the following sections we will discuss

119     possible reasons as to why this occurred.

120

121     Firstly, multiple selection steps occurred during generation of the SAGs that may have excluded

122     or did exclude genomes from taxa abundant in the metagenome and MAGs. While 365 events

123     (cells) were sorted during single cell sorting, only the 95 wells with the greatest estimated

124     genome amplification (based on fluorescence during the amplification reaction) were selected for

125     further analysis. This may have resulted in a biased selection of cells: for example, selecting for

126     those genomes most responsive to the MDA reaction or those most resistant to environmental

127     fluxes that may have occurred during transportation or sorting. Additionally, after an initial

128     round of 16S rRNA Sanger sequencing, some wells with highly similar (>98%) 16S rRNA

129     sequences were excluded from subsequent genomic sequencing to maximize sequencing

130     coverage on the remaining wells. This included several *Halomonas* and *Desulfobulbaceae*

131     genomes represented by MAGs, therefore potentially reducing the number of SAGs that may

132     have corresponded to MAGs. As a result of these filtering steps, SAGs corresponding to MAGs

133     may have been excluded and low-abundance taxa including archaeal genomes may have been

134     enriched.

135

136     Secondly, 20% of reads map to the high- and medium-quality MAGs, due in part to stringent

137     quality control during binning (i.e. only contigs >5000 bp were binned). As a result, while the

138    MAGs broadly represent the taxonomic groups present in the metagenome, they represent only a

139    portion of the sum diversity. Therefore, this restrictive binning process in combination with the

140    relatively high number of SAGs corresponding to taxa at low abundance in the metagenome may

141    also have contributed to the discrepancy.

142

143    Thirdly, the samples collected for metagenomic analysis and for SAG sequencing were collected

144    in different years (2017 vs. 2019, respectively). This was necessitated by the small amount of

145    sediment removed during sampling to avoid disturbing the site in combination with the low

146    biomass of the sediment, requiring relatively high amounts of sample for processing. While the

147    physical and geochemical parameters of the spring have remained stable for the ~15 years that

148    we have studied the spring, including the measurements taken in 2017 and 2019, we can't

149    exclude the possibility of minor changes in the spring that could have affected the microbial

150    community sampled. Additionally, due to the low biomass and high salinity of the sediment, it is

151    difficult to extract DNA. Biases introduced by the differing processing steps during MAG and

152    SAG sequencing (for example, DNA extraction vs. separation of cells from sediment) may

153    therefore have had outsize effects.

154

155    To summarize, we suggest several factors potentially contributing to the discrepancy between

156    MAGs and SAGs: 1) Selection during generation of SAGs, 2) Stringent binning criteria, and 3)

157    Differences in input samples and biases introduced by differing experimental procedures. We

158    conclude by noting that although there is a discrepancy between the MAGs and SAGs, there is

159    overlap between the two samples additional to what was discussed in the manuscript. In addition

160    to the corresponding SAG and MAG noted in the manuscript, two additional SAGs had 16S

161    rRNA sequences >98% identical to unbinned 16S rRNA sequences, and nearly all taxonomic

162    groups represented by the SAGs were also present in the metagenome reads and 16S sequences

163    as classified by kaiju and PhyloFlash (with the sole exception of Iainarchaeota reads).

164    Additionally, comparison of 16S rRNA sequences in the SAGs to previous 16S rRNA

165    sequencing from over ten years ago (1) identified common sequences (>98%) between the two

166    datasets, including sequences for *Halomonas*, ANME-1, and *Iainarchaeota*, suggesting that the

167    discrepancy is more likely due to the potential technical factors discussed above rather than

168    significant changes in the microbial community. The relative abundances of the microbial

169    community represented in the metagenome are consistent with those observed in previous

170    CARD-FISH and 16S rRNA and metagenomic sequencing (1, 2). We therefore suggest that the

171    metagenome and MAGs more accurately represent the taxa abundant in the microbial

172    community, whereas the SAGs are disproportionately enriched in low-abundance bacteria and

173    archaea due to the potential factors discussed above.
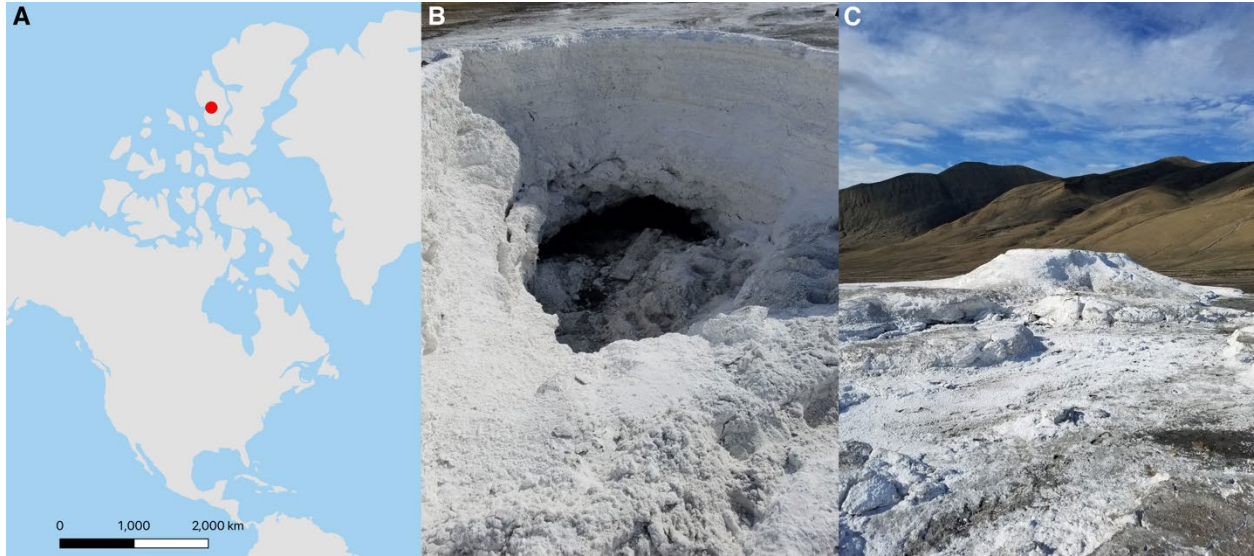
174

175

176

177

178

179

180

181

182

183

184    **Supplementary Figures and Tables**



186    **Figure S1.** A. Location of the LH spring on Axel Heiberg Island in the Canadian High Arctic.

187    Map generated in QGIS with the Natural Earth dataset. B. View looking into the LH salt tufa to

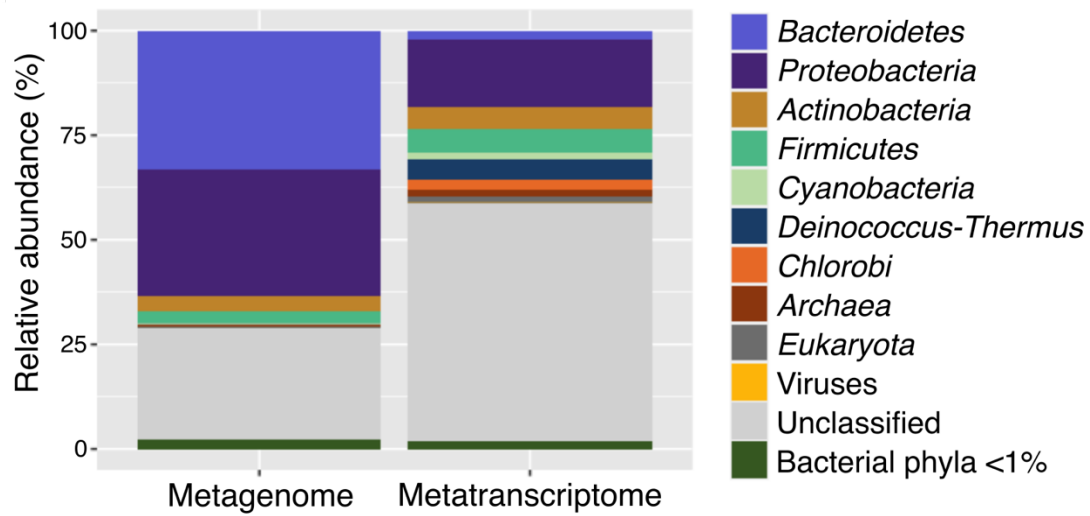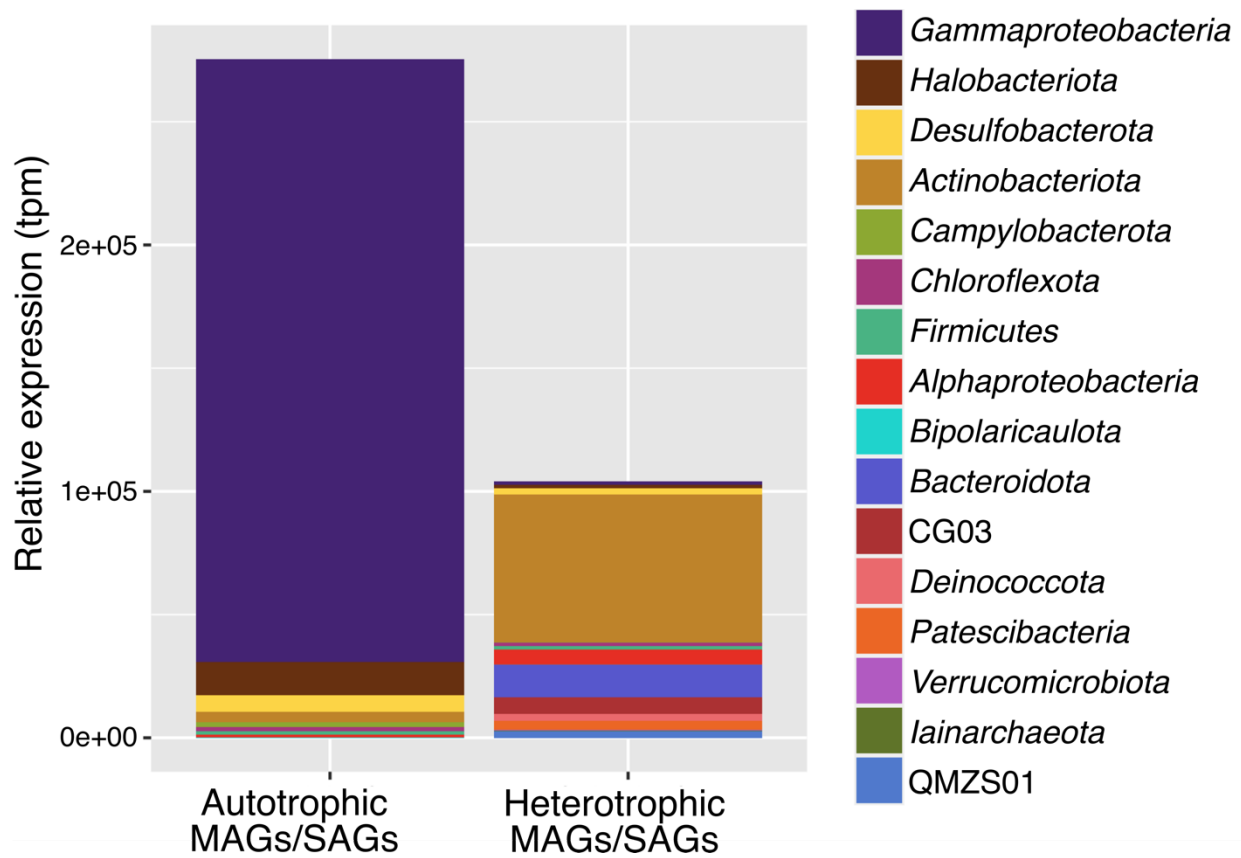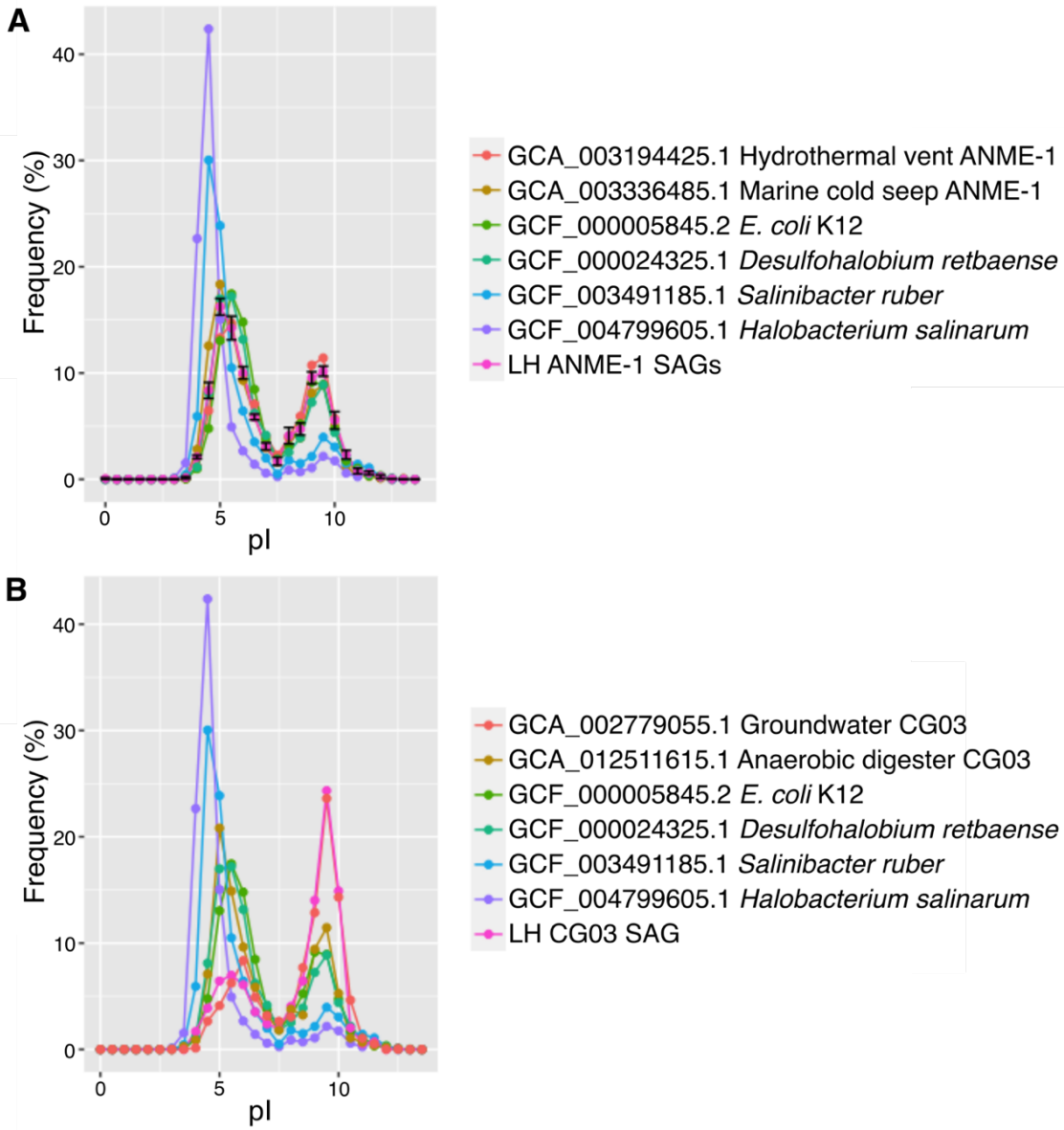188    the spring source (July 2019). C. View of the LH tufa and outflow channels (July 2019).

194

**Figure S2.** Taxonomic classification of metagenome and metatranscriptome reads. Reads were

classified by Kaiju using the NCBI non-redundant database. Metagenome reads are an average of

the relative abundance in duplicate samples.

198

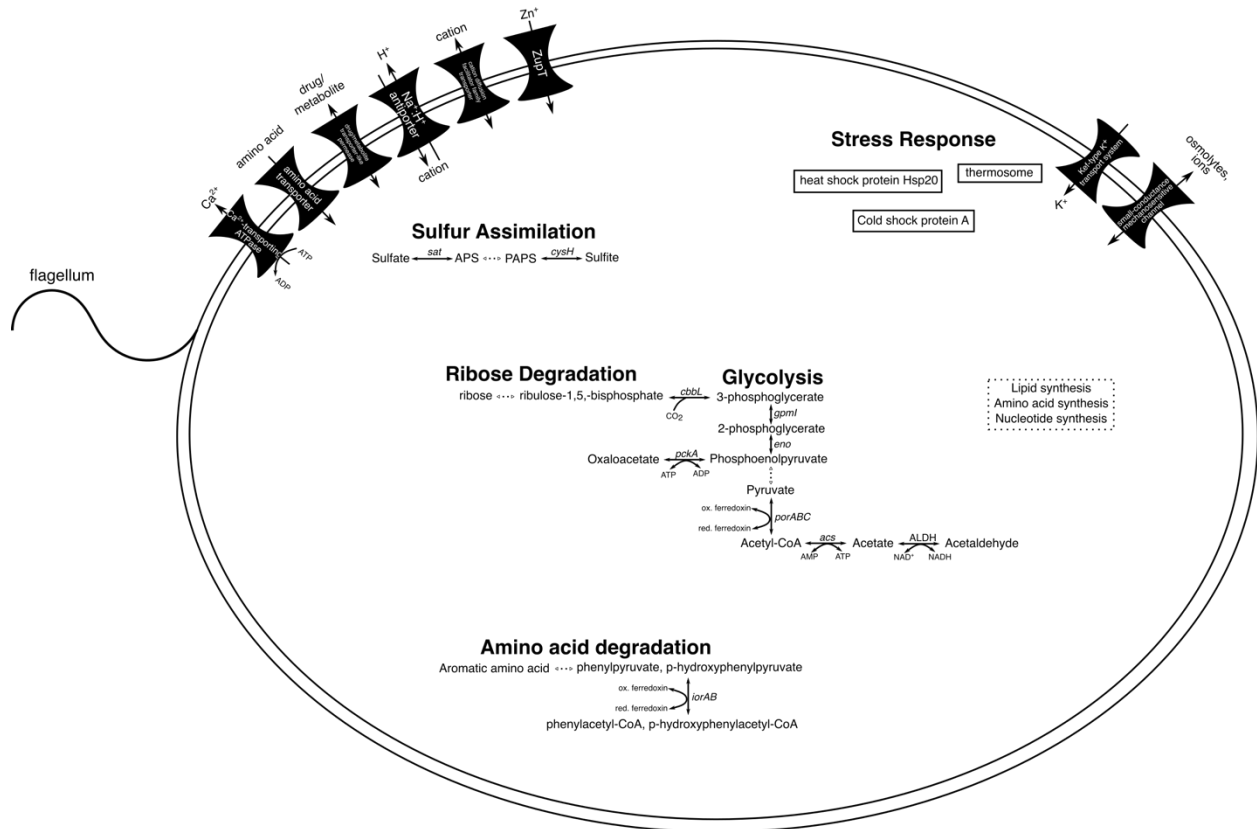**Figure S3.** Relative expression (tpm) attributed to MAGs and SAGs containing $CO_2$-fixation

genes (autotrophic MAGs/SAGs) and those without (heterotrophic MAGs/SAGs). Unbinned

genes with mapped transcripts were omitted from this analysis.

**Figure S4.** Comparison of predicted protein isoelectric points for A. ANME-1 SAGs S10-S14 and B. CG03 SAG S2. Protein isoelectric points were calculated with the ExPASy Compute pI/MW online tool. Genomes for comparison (based on similar analysis in Nigro et al. (37)) include microorganisms known to accumulate salts through the "salting-in" osmoregulation strategy (*Salinibacter ruber* and *Halobacterium salinarum*) and those that do not accumulate salts (*E. coli* K12 and *Desulfohalobium retbaense*).

**Figure S5.** Metabolic reconstruction of *Iainarchaeia* SAG S9 (670822 bp, 52.9% completeness, 0% contamination). Solid lines represent genes present within the genome, and dashed lines indicate steps or pathways absent from the genome. Table S15 contains a complete list of annotated genes represented in this figure.

216

**Figure S6.** Metabolic reconstruction of QMZS01 SAGs S4-S8 (between 563912-779153 bp,

57.7-76.2% completeness, 0.8-2.2% contamination). Solid lines represent genes present within

the genome, and dashed lines indicate steps or pathways absent from the genome. Table S16

contains a complete list of annotated genes represented in this figure.
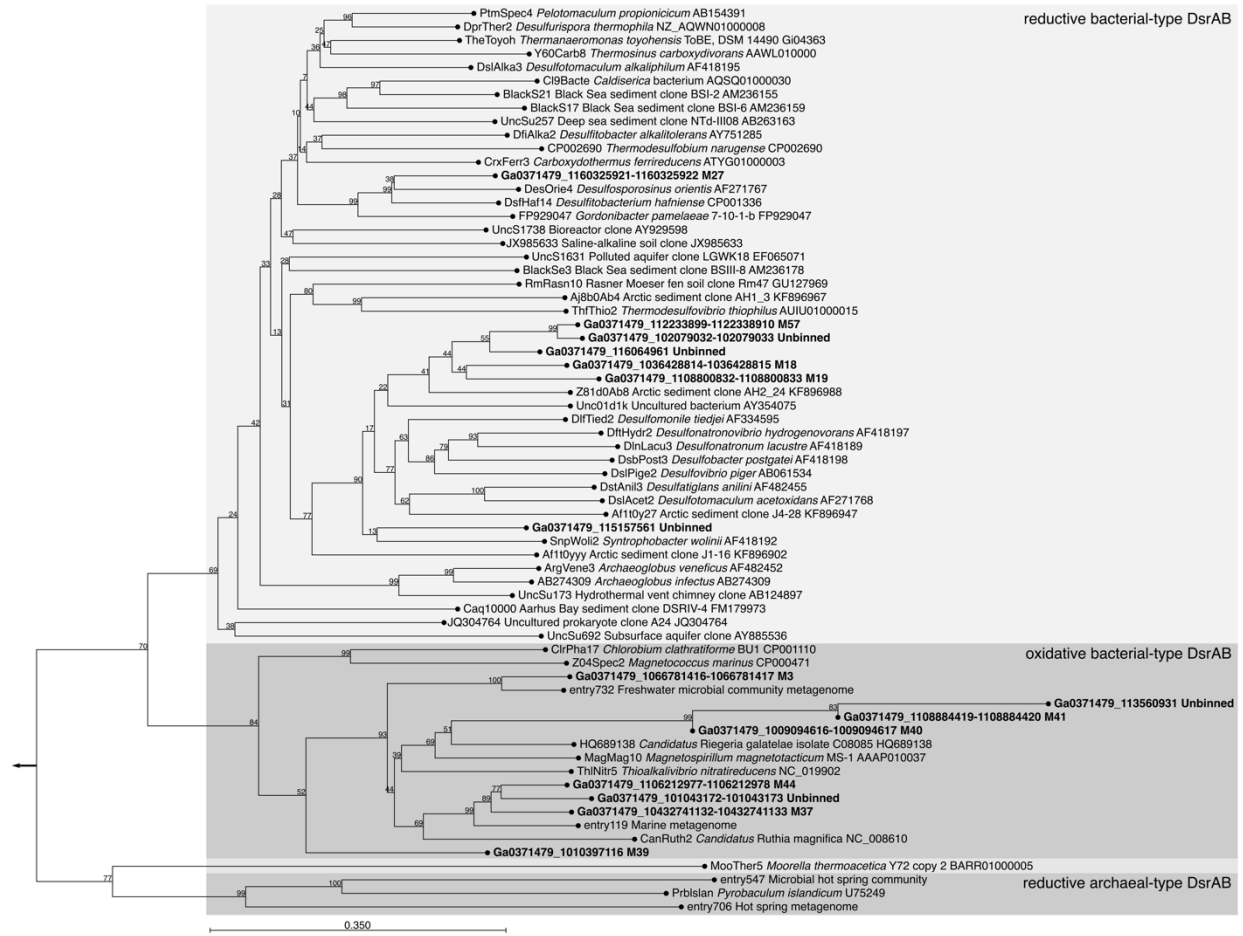
229



**Figure S7.** Maximum likelihood phylogenetic tree of RbcL/CbbL sequences. The tree was

constructed in CLC Genomics Workbench with 1000 bootstraps and WAG substitution model

using reference sequences from NCBI. A RubisCO-like protein sequence from

*Rhodopseudomonas palustris* TIE-1 (ACF00976.1) was used as an outgroup (direction indicated

by arrow).

**Figure S8.** Metabolic reconstruction of candidate phylum CG03 SAG S2 (1039757 bp, 62.8% completeness, 0% contamination). Solid lines represent genes present within the genome, and dashed lines indicate steps or pathways absent from the genome. Table S17 contains a complete list of annotated genes represented in this figure.

244

**Figure S9.** Maximum likelihood phylogenetic tree of DsrAB. The tree was constructed in CLC

Genomics Workbench with 1000 bootstraps and WAG substitution model using reference

sequences from Muller et al. (2015). Sequence names include the gene IDs and MAG number

where applicable. A sequence from eggNOG group COG2221 (*Campylobacter ureolyticus*

JFJK01000015_gene476) was used as an outgroup (direction indicated by arrow).

250

251

252

253

254

255   **Table S1 (xlsx).** Physical and geochemical parameters of Lost Hammer water and sediment.

256   Located in Supplementary Tables xlsx file.

257

258   **Table S2.** Metagenome and metatranscriptome sequencing statistics.

| DNA/RNA | Sample (NCBI ID) | Sequencing platform | Read length | Raw read pairs | Read pairs after quality control |
|---------|------------------|---------------------|-------------|----------------|----------------------------------|
| DNA | SRR13628066 | Illumina HiSeq2500 | 2x126 | 116844950 | 116746762 |
| DNA | SRR13628065 | Illumina HiSeq2500 | 2x126 | 107914225 | 107647735 |
| RNA | SRR13633097 | Illumina NovaSeq6000 | 2x100 | 9009952 | 4233149 |

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274    **Table S3.** Metagenomic assembly statistics.

| Assembler | Megahit v1.1.3 |
|---|---|
| **# contigs** | 1620755 |
| **# contigs >= 1000 bp** | 276578 |
| **# contigs >= 5000 bp** | 36050 |
| **# contigs >= 10000 bp** | 14841 |
| **# contigs >= 25000 bp** | 3946 |
| **# contigs >= 50000 bp** | 1251 |
| **Total length (bp)** | 1599970111 |
| **GC %** | 51.63 |
| **N50** | 2968 |
| **N75** | 1032 |
| **L50** | 67307 |
| **L75** | 264669 |
| **% reads mapping to contigs (average of two samples)** | 89.7 |
| **% reads mapping to contigs > 5000 bp** | 67.9 |
| **Average coverage of metatranscriptome against the metagenome (reads/base)** | $0.093 \pm 59$ |
| **Average coverage of metatranscriptome against CDS with mapped reads (reads/base)** | $13 \pm 604$ |

275
276
277
278
279
280
281
282
283

284    **Table S4 (xlsx).** SAG supplemental information. Located in Supplementary Tables xlsx file.

285

286    **Table S5 (xlsx).** MAG supplemental information. Located in Supplementary Tables xlsx file.

287

288    **Table S6 (xlsx).** Relative expression for KEGG, COG, and pfam IDs. Located in Supplementary

289    Tables xlsx file.

290

291    **Table S7 (xlsx).** Taxonomic classification of genes of interest with mapped transcripts. Located

292    in Supplementary Tables xlsx file.

293

294    **Table S8 (xlsx).** Gene presence and expression in medium-quality MAGs. Located in

295    Supplementary Tables xlsx file.

296

297    **Table S9 (xlsx).** Complete list of genes with mapped transcripts. Located in Supplementary

298    Tables xlsx file.

299

300

301

302

303

304

305

306

307 **Table S10.** Gibbs free energy of redox pairs present in LH. Values were calculated using the

308 methodology in Jones et al. (2018) (35) (additional details in Supplementary Methods). These

309 values should be considered preliminary estimates as some parameter concentrations are based

310 on single-year measurements.

| Redox pair | $\triangle$G reaction (kJ/mol electron$^{-1}$) |
|---|---|
| $H_2/O_2$ | -111.9 |
| $H_2/NO_3^-$ (to $NH_3$) | -61.3 |
| $N_2/NO_3^-$ (to $NO_2^-$) | -68.1 |
| $H_2/SO_4^{2-}$ | -3.0 |
| $H_2/CO_2$ | -9.1 |
| $H_2S/O_2$ | -109.0 |
| $H_2S/NO_3^-$ | -93.8 |
| $Fe^{2+}/O_2$ | -16.0 |
| $NH_3/O_2$ | -52.9 |
| $NH_3/NO_3^-$ | -90.8 |
| $NH_3/SO_4^{2-}$ | -56.5 |
| $CH_4/O_2$ | -205.7 |
| $CH_4/NO_3^-$ | -32.4 |
| $CH_4/SO_4^{2-}$ | 6.1 |

311
312
313
314
315
316
317
318

319     **Table S11 (xlsx).** Stress response gene comparison to related genomes. Located in

320     Supplementary Tables xlsx file.

321

322     **Table S12 (xlsx).** Gene content of MAGs. Located in Supplementary Tables xlsx file.

323

324     **Table S13 (xlsx).** Gene content of SAGs. Located in Supplementary Tables xlsx file.

325

326     **Table S14 (xlsx).** ANME-1 composite genome gene content. Located in Supplementary Tables

327     xlsx file.

328

329     **Table S15 (xlsx).** *Iainarchaeia* sp. S9 gene content. Located in Supplementary Tables xlsx file.

330

331     **Table S16 (xlsx).** QMZS01 composite genome gene content. Located in Supplementary Tables

332     xlsx file.

333

334     **Table S17 (xlsx).** CG03 sp. S2 gene content. Located in Supplementary Tables xlsx file.

335

336     **Table S18 (xlsx).** Metagenome copy number and total relative expression of genes of interest.

337     Located in Supplementary Tables xlsx file.

338

339

340

341

342    **References**

343    1. Niederberger TD, Perreault NN, Tille S, Lollar BS, Lacrampe-Couloume G, Andersen D, et al.

344    Microbial characterization of a subzero, hypersaline methane seep in the Canadian High Arctic.

345    ISME J. 2010;4(10):1326-39.

346    2. Lay CY, Mykytczuk NC, Yergeau E, Lamarche-Gagnon G, Greer CW, Whyte LG. Defining

347    the functional potential and active community members of a sediment microbial community in a

348    high-Arctic hypersaline subzero spring. Appl Environ Microbiol. 2013;79(12):3637-48.

349    3. Lamarche-Gagnon G, Comery R, Greer CW, Whyte LG. Evidence of *in situ* microbial activity

350    and sulphidogenesis in perennially sub-0 degrees C and hypersaline sediments of a high Arctic

351    permafrost spring. Extremophiles. 2015;19(1):1-15.

352    4. Battler MM, Osinski GR, Banerjee NR. Mineralogy of saline perennial cold springs on Axel

353    Heiberg Island, Nunavut, Canada and implications for spring deposits on Mars. Icarus.

354    2013;224(2):364-81.

355    5. Ward MK, Pollard WH. A hydrohalite spring deposit in the Canadian high Arctic: a potential

356    Mars analogue. Earth Planet Sci Lett. 2018;504:126-38.

357    6. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data.

358    Bioinformatics. 2014;30(15):2114-20.

359    7. Menzel P, Ng KL, Krogh A. Fast and sensitive taxonomic classification for metagenomics

360    with Kaiju. Nat Commun. 2016;7(1):1-9.

361    8. Gruber-Vodicka HR, Seah BKB, Pruesse E. phyloFlash: rapid small-subunit rRNA profiling

362    and targeted assembly from metagenomes. mSystems. 2020;5(5):1-16.

363  9. Li D, Liu CM, Luo R, Sadakane K, Lam TW. MEGAHIT: an ultra-fast single-node solution

364  for large and complex metagenomics assembly via succinct de Bruijn graph. Bioinformatics.

365  2015;31(10):1674-6.

366  10. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. metaSPAdes: a new versatile

367  metagenomic assembler. Genome Res. 2017;27(5):824-34.

368  11. Kang DD, Froula J, Egan R, Wang Z. MetaBAT, an efficient tool for accurately

369  reconstructing single genomes from complex microbial communities. PeerJ. 2015;3:1-15.

370  12. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the

371  quality of microbial genomes recovered from isolates, single cells, and metagenomes. Genome

372  Res. 2015;25(7):1043-55.

373  13. Mikheenko A, Prjibelski A, Saveliev V, Antipov D, Gurevich A. Versatile genome assembly

374  evaluation with QUAST-LG. Bioinformatics. 2018;34(13):i142-i150.

375  14. Chen IA, Chu K, Palaniappan K, Ratner A, Huang J, Huntemann M, et al. The IMG/M data

376  management and analysis system v.6.0: new tools and advanced capabilities. Nucleic Acids Res.

377  2020;49(D1):D751–D763.

378  15. Mukherjee S, Stamatis D, Bertsch J, Ovchinnikova G, Sundaramurthi JC, Lee J, et al.

379  Genomes OnLine Database (GOLD) v.8: overview and updates. Nucleic Acids Res.

380  2020;49(D1):D723-D733.

381  16. Chaumeil PA, Mussig AJ, Hugenholtz P, Parks DH. GTDB-Tk: a toolkit to classify genomes

382  with the Genome Taxonomy Database. Bioinformatics. 2019;36(6):1925–7.

383  17. Eren AM, Esen OC, Quince C, Vineis JH, Morrison HG, Sogin ML, et al. Anvi'o: an

384  advanced analysis and visualization platform for 'omics data. PeerJ. 2015;3:1-29.

385    18. Garber AI, Nealson KH, Okamoto A, McAllister SM, Chan CS, Barco RA, et al. FeGenie: a

386    comprehensive tool for the identification of iron genes and iron gene neighborhoods in genome

387    and metagenome assemblies. Front Microbiol. 2020;11:1-23.

388    19. Fish JA, Chai B, Wang Q, Sun Y, Brown CT, Tiedje JM, et al. FunGene: the functional gene

389    pipeline and repository. Front Microbiol. 2013;4:1-14.

390    20. Sondergaard D, Pedersen CN, Greening C. HydDB: a web tool for hydrogenase classification

391    and analysis. Sci Rep. 2016;6:1-8.

392    21. Muller AL, Kjeldsen KU, Rattei T, Pester M, Loy A. Phylogenetic and environmental

393    diversity of DsrAB-type dissimilatory (bi)sulfite reductases. ISME J. 2015;9(5):1152-65.

394    22. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.

395    Nucleic Acids Res. 2004;32(5):1792-7.

396    23. Anantharaman K, Hausmann B, Jungbluth SP, Kantor RS, Lavy A, Warren LA, et al.

397    Expanded diversity of microbial groups that shape the dissimilatory sulfur cycle. ISME J.

398    2018;12(7):1715-28.

399    24. Stepanauskas R, Fergusson EA, Brown J, Poulton NJ, Tupper B, Labonte JM, et al.

400    Improved genome recovery and integrated cell-size analyses of individual uncultured microbial

401    cells and viral particles. Nat Commun. 2017;8(1):1-10.

402    25. Kuske CR, Barns SM, Grow CC, Merrill L, Dunbar J. Environmental survey for four

403    pathogenic bacteria and closely related species using phylogenetic and functional genes. J

404    Forensic Sci. 2006;51(3):548-58.

405    26. Miller CS, Handley KM, Wrighton KC, Frischkorn KR, Thomas BC, Banfield JF. Short-read

406    assembly of full-length 16S amplicons reveals bacterial diversity in subsurface sediments. PLoS

407    One. 2013;8(2):1-11.

408    27. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal

409    RNA gene database project: improved data processing and web-based tools. Nucleic Acids Res.

410    2013;41:D590-596.

411    28. Schmieder R, Edwards R. Fast identification and removal of sequence contamination from

412    genomic and metagenomic datasets. PLoS One. 2011;6(3):1-11.

413    29. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a

414    new genome assembly algorithm and its applications to single-cell sequencing. J Comput Biol.

415    2012;19(5):455-77.

416    30. Jain C, Rodriguez RL, Phillippy AM, Konstantinidis KT, Aluru S. High throughput ANI

417    analysis of 90K prokaryotic genomes reveals clear species boundaries. Nat Commun.

418    2018;9(1):1-8.

419    31. Kopylova E, Noe L, Touzet H. SortMeRNA: fast and accurate filtering of ribosomal RNAs

420    in metatranscriptomic data. Bioinformatics. 2012;28(24):3211-7.

421    32. Sheik CS, Reese BK, Twing KI, Sylvan JB, Grim SL, Schrenk MO, et al. Identification and

422    removal of contaminant sequences from ribosomal gene databases: lessons from the Census of

423    Deep Life. Front Microbiol. 2018;9:1-14.

424    33. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods.

425    2012;9(4):357-9.

426    34. Anders S, Pyl PT, Huber W. HTSeq--a Python framework to work with high-throughput

427    sequencing data. Bioinformatics. 2015;31(2):166-9.

428    35. Jones RM, Goordial JM, Orcutt BN. Low energy subsurface environments as extraterrestrial

429    analogs. Front Microbiol. 2018;9:1-18.

430    36. Alneberg J, Karlsson CMG, Divne AM, Bergin C, Homa F, Lindh MV, et al. Genomes from

431    uncultivated prokaryotes: a comparison of metagenome-assembled and single-amplified

432    genomes. Microbiome. 2018;6(1):1-14.

433    37. Nigro LM, Hyde AS, MacGregor BJ, Teske A. Phylogeography, salinity adaptations and

434    metabolic potential of the candidate division KB1 bacteria based on a partial single cell genome.

435    Front Microbiol. 2016;7:1-9.