

## **Supplementary Information**

**Energy efficiency and biological interactions define the core microbiome of deep oligotrophic groundwater**

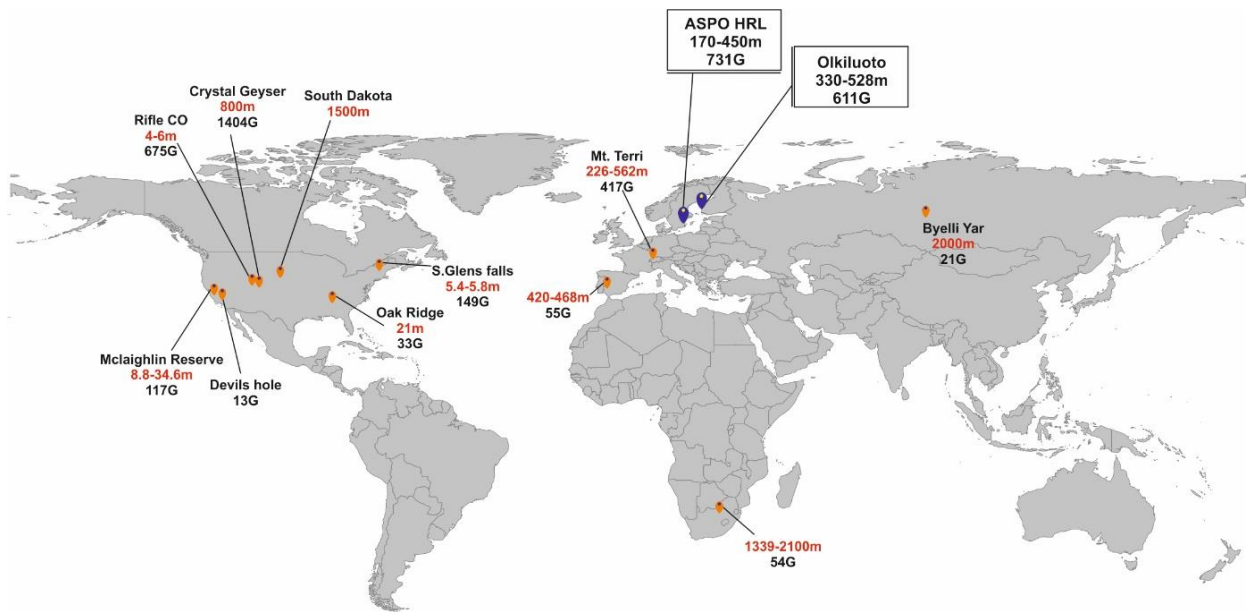
Mehrshad, et al.

**Supplementary Figures and legends**

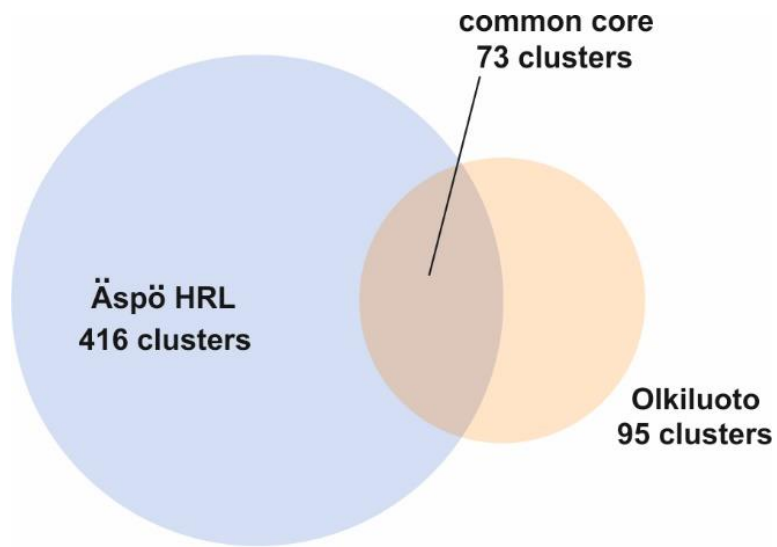
**Supplementary Background**

**Supplementary Methods**

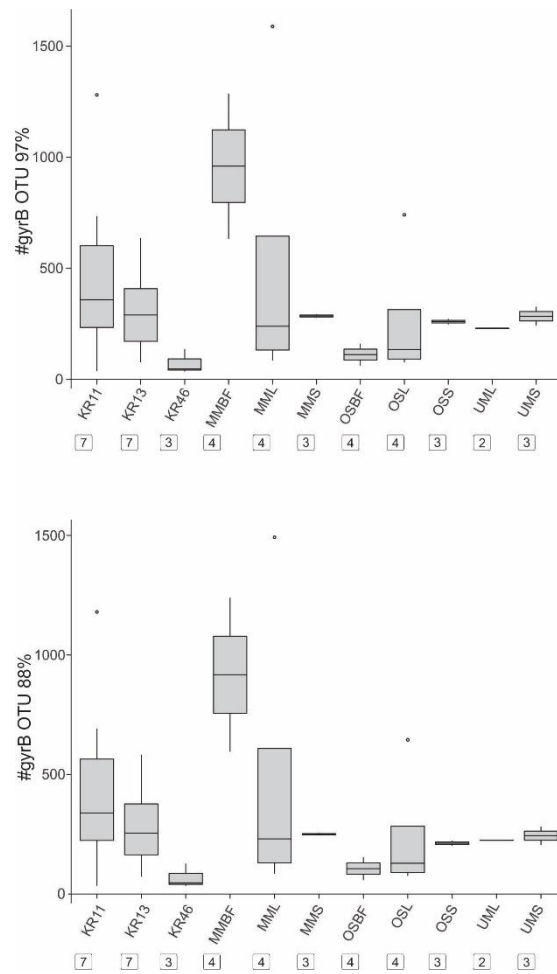
## Supplementary Figures and legends



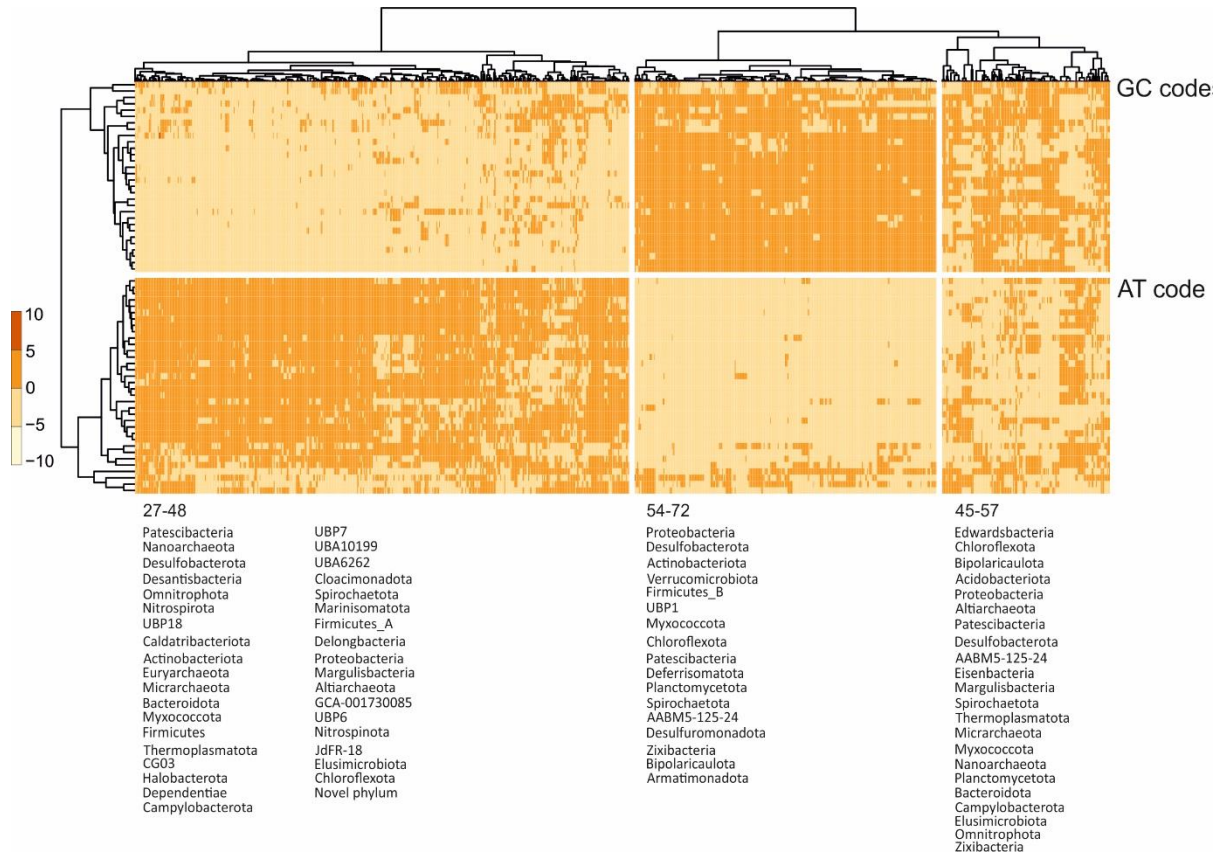
**Supplementary Fig. S1-** Geographical distribution of publicly available metagenomic datasets sequenced from oligotrophic groundwater samples (landfill groundwater, oil-influenced, and shale samples are not included in this representation). The depth ranges (as mentioned in the corresponding publication) of the samples (in red) and the amount of publicly available sequenced data (presented as Gb) are shown for each location.



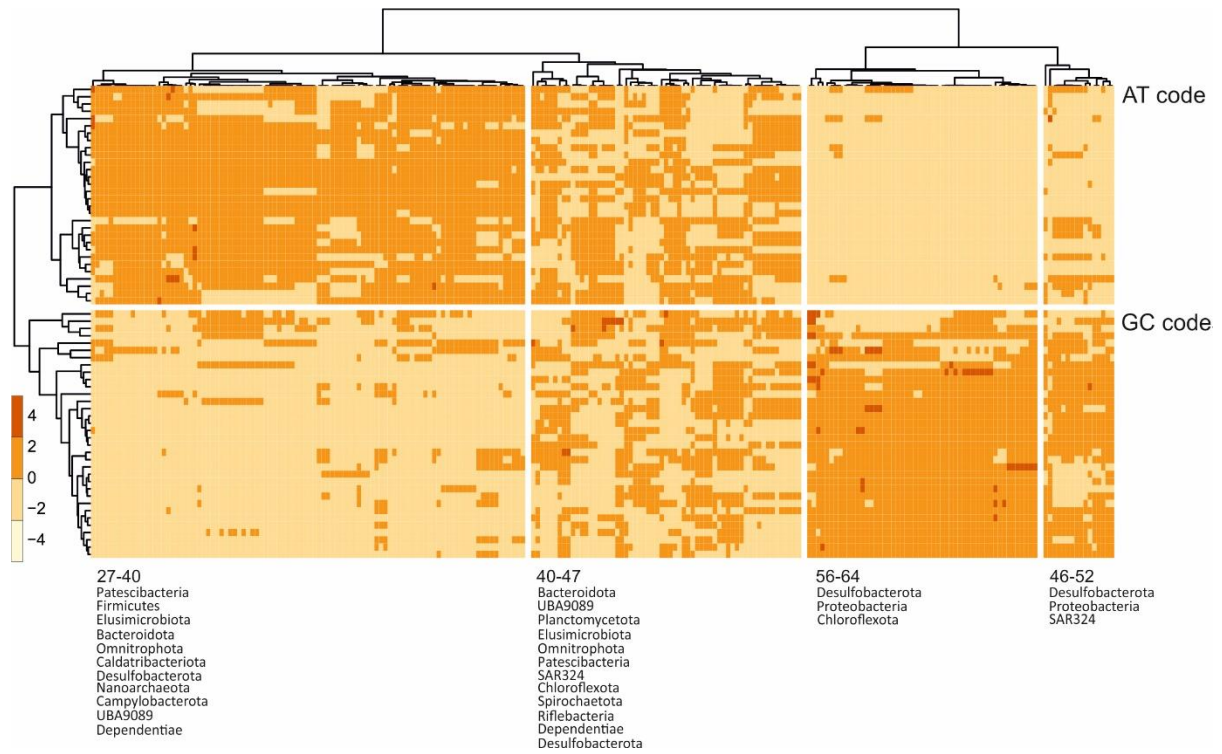
**Supplementary Fig. S2-** Overlap of the microbiome of the Fennoscandian Shield deep groundwater flowing in two disconnected sites. Site name and number of clusters are mentioned in the diagram.



**Supplementary Fig. S3-** Species richness in metagenomes generated from each deep groundwater site based on gene *gyrB* (centerline, median; hinge limits, 25 and 75% quartiles; whiskers, 1.5x interquartile range; points, outliers). Numbers next to the sample names represent the number of metagenomes generated and analysed for this figure.

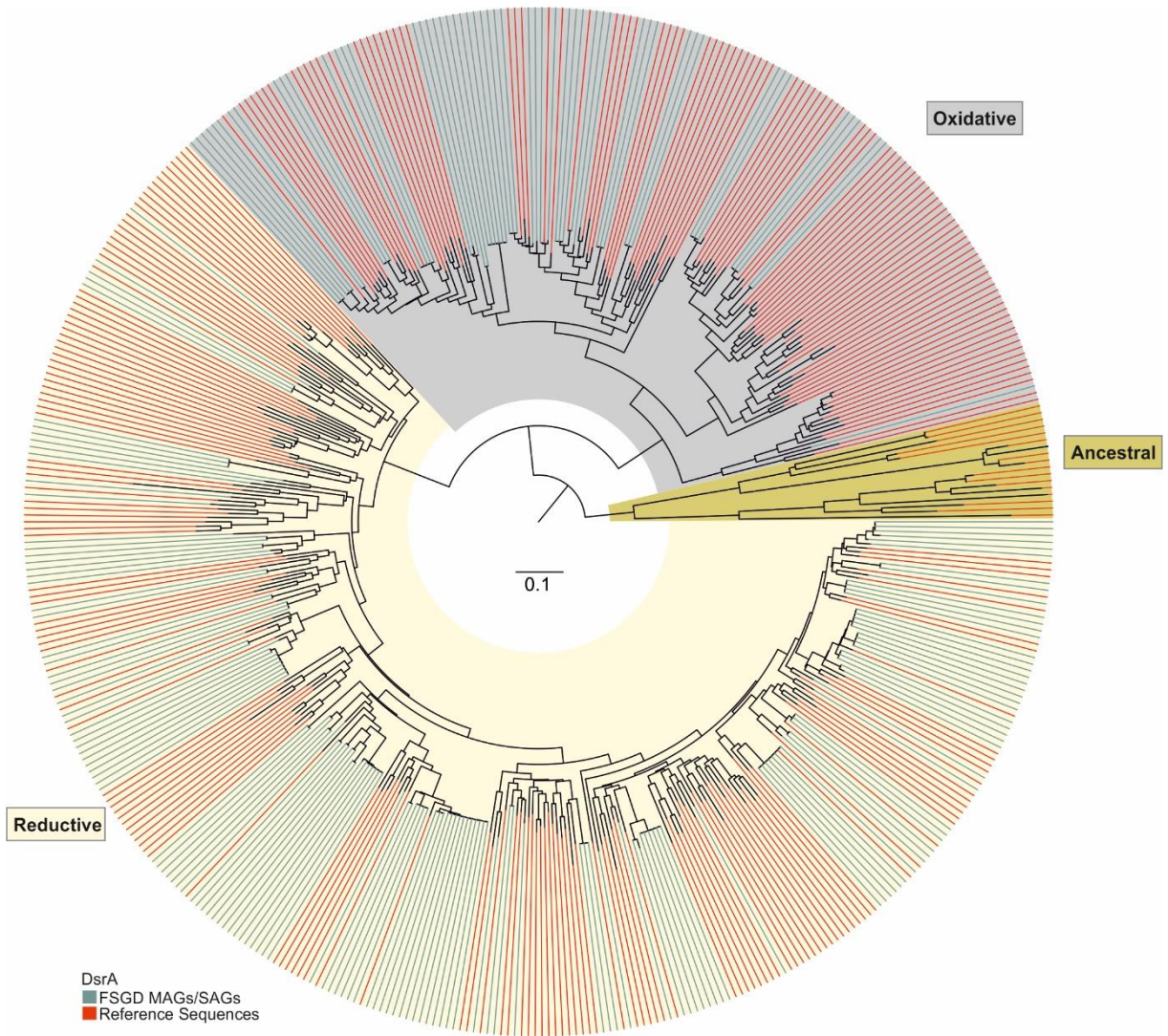


**Supplementary Fig. S4-** Representation of the frequency (the expected number of codons, given the input sequences, per 1000 bases) of utilization of synonymous codons across MAGs and SAGs of the FSGD.

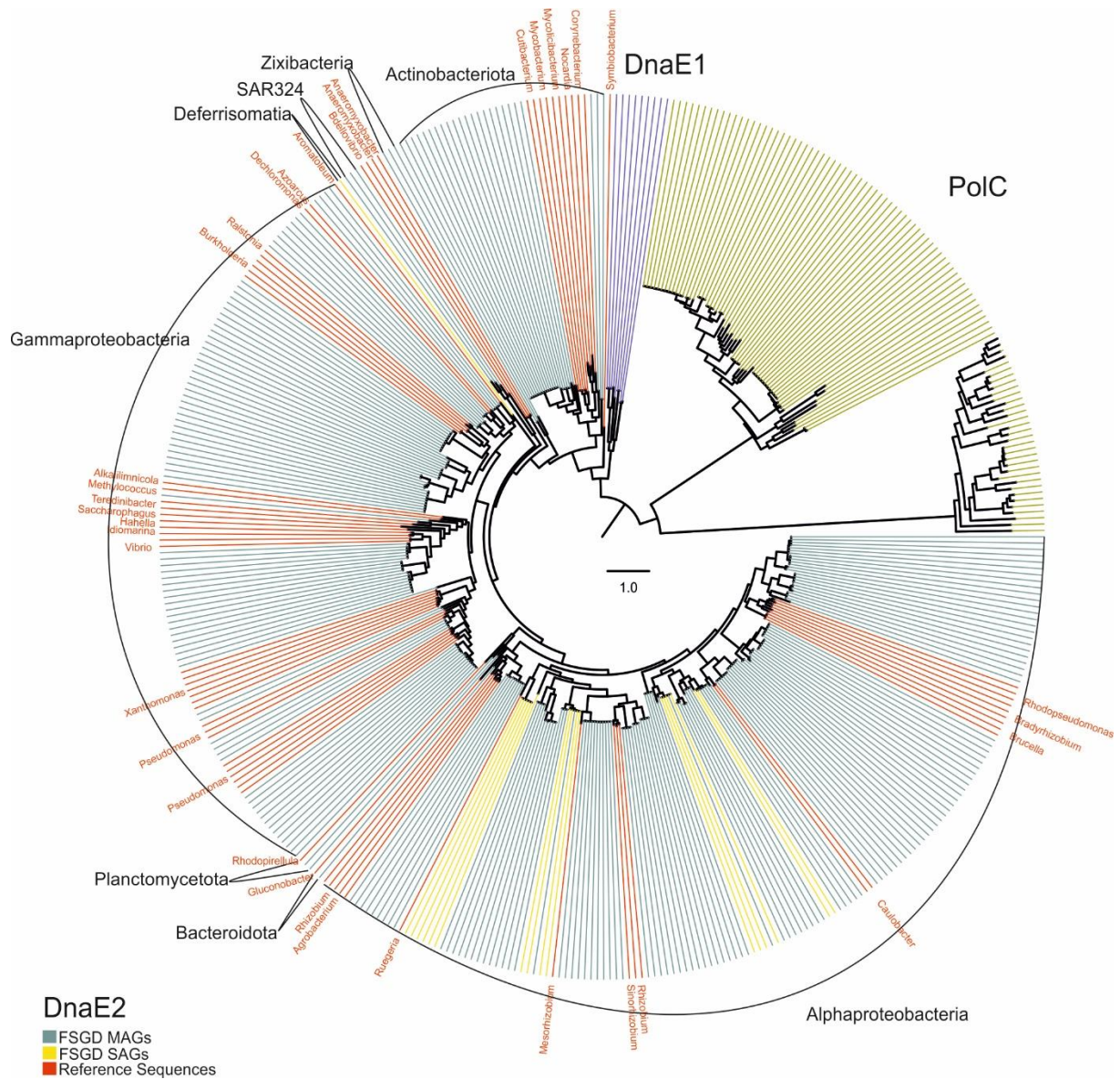


**Supplementary Fig. S5-** Representation of the frequency (the expected number of codons, given the input sequences, per 1000 bases) of utilization of synonymous codons across highly expressed (TPM >10000 arbitrary threshold) MAGs and SAGs of the FSGD.





**Supplementary Fig. S6-** Reconstructed phylogeny of the DsrA protein sequences recovered from FSGD reconstructed MAGs and SAGs (dark blue) together with reference DsrA sequences (in red). Different types of DsrA protein according to their function are mentioned in the figure. The alignment and the phylogenetic tree are publicly available via figshare at DOI:10.6084/m9.figshare.14166638 and DOI:10.6084/m9.figshare.14166650.v1 respectively.



**Supplementary Fig. S7-** Reconstructed phylogeny of the C-family polymerases. The tree is arbitrarily rooted by the PolC branch. Taxonomy of the leaves are written for each group. The reference genomes are the reviewed sequences retrieved from uniprot for each polymerase type (red). The alignment used to reconstruct this phylogeny and the unrooted tree are available via figshare at DOI:10.6084/m9.figshare.12170310.v1 and DOI:10.6084/m9.figshare.13298513 respectively.



## 1 **Background**

2           During the last decades, it has been demonstrated that microorganisms inhabit the sub-  
3 surface biome down to several kilometers below the surface<sup>1</sup>. Although microbial life in the  
4 deep crust biosphere has been gaining interest<sup>2</sup>, due to its difficulty to access it is still one of  
5 the least understood environments on earth<sup>3</sup>. Consequently, many novel taxa are being  
6 identified in the deep biosphere<sup>4,5</sup>.

7           The continental deep biosphere is estimated to contain  $2 \text{ to } 6 \times 10^{29}$  cells<sup>6,7</sup>, containing  
8 members from all domains of life along with viruses<sup>8-14</sup>. Even though the available energy flux  
9 is very low compared to that of the photosynthetically fixed carbon on the surface<sup>15</sup>,  
10 microorganisms are widely spread in oceanic crust fluids<sup>16-18</sup>, marine sediments<sup>19,20</sup>,  
11 terrestrial rocks<sup>21,22</sup>, and granitic groundwaters<sup>23-25</sup>. In addition, many of these deep  
12 biosphere microbes are both alive and active<sup>10,26-30</sup>.

13           In continental groundwaters, microbial activity is strongly positively correlated to the  
14 proximity of the photosynthesis-fueled surface<sup>7</sup> and thus, water-bearing deep fracture  
15 systems are extremely oligotrophic<sup>31</sup>. A recent study showed that there is a steady flow of  
16 surface organisms to the deep subsurface with a selection event resulting in some taxa  
17 adapting to the new conditions while others perish<sup>24</sup>. In addition, Lopez-Fernandez et al.<sup>28</sup>  
18 showed the presence of viable taxa in a deep continental crystalline rock and that any non-  
19 viable cells are rapidly degraded and recycled into new biomass<sup>28</sup>. The deep biosphere is  
20 suggested to be adapted to the low energy conditions by e.g. small cell size and streamlined  
21 genomes<sup>32</sup>. Microbial populations with the potential to initiate biofilm formation were also  
22 identified in deep terrestrial subsurface waters<sup>11</sup>. This close proximity likely promotes  
23 syntrophy and cycling of nutrients between populations<sup>33</sup>.

## 24 **Site lithologies**

25           The two sampling locations, Äspö HRL and Olkiluoto Island drillholes, are situated on  
26 opposite sides of the Baltic Sea on the Swedish (Lat N 57° 26' 4'' Lon E 16° 39' 36'') and Finnish  
27 (Lat N 61° 14' 31'', Lon E 21° 29' 23'') coasts. The crystalline bedrock of Sweden and Finland is  
28 part of the Precambrian Fennoscandian Shield that is predominantly made up of granite and  
29 quartz monzodiorite of quartz and aluminosilicate minerals including mica and feldspar.

30           The bedrock in the Äspö HRL region consists of overall well preserved (locally low-grade  
31 metamorphism and discrete foliation occur) Palaeoproterozoic ~1.8 Ga granitoids of the  
32 Transscandinavian Igneous Belt<sup>34</sup>. At the Äspö HRL site, these rocks have a composition  
33 ranging from diorite to granite<sup>35</sup>. The rocks are also cut by a number of deformation zones and  
34 open fractures, which frequently have surfaces covered by high and/or low temperature 1-15  
35 mm thick secondary-mineral precipitates including calcite, chlorite, pyrite, clay minerals,  
36 epidote, adularia, and hematite<sup>36</sup>. Hence, the fluids flowing in the fractures are in direct  
37 contact with both Precambrian granitoids and a variety of secondary minerals of variable age.

38           Olkiluoto is located in the southern Satakunta region in south-western Finland. The ~1.8  
39 Ga Palaeoproterozoic bedrock of the southern Satakunta region is composed of supracrustal,  
40 metasedimentary and metavolcanics rocks deformed and metamorphosed during the  
41 Svecofennian orogeny<sup>37</sup>. The main rock types at Olkiluoto are mica and veined gneisses,  
42 migmatite granite, grey gneisses and diabase. Minor veins and dykes are quartz feldspar  
43 gneisses and amphibolites. Mica and veined gneisses contain calcite, sulphides and clay  
44 mineral fracture fillings<sup>38</sup>.

45

46

## 47 **Supplementary Methods**

48 In this study, the extremely oligotrophic deep groundwaters of two sites excavated in the  
49 Fennoscandian Shield have been sampled and extensively analyzed via a “multi-omics”  
50 approach.

51

## 52 **Site description**

53 The Äspö Hard Rock Laboratory (HRL) is located in the south-east of Sweden (Lat N 57°  
54 26' 4" Lon E 16° 39' 36") and is excavated in the Proterozoic crystalline bedrock of the  
55 Fennoscandian Shield extending to a depth of 460 m below sea level (mbsl) with 3600 m of  
56 total tunnel length<sup>42</sup>. The tunnel provides access to investigate the microbial life in the deep  
57 Fennoscandian Shield groundwater.

58 Building of the Äspö Hard Rock Laboratory was initiated in 1986 and the site  
59 circumvents many problems normally associated with sources of contamination in  
60 groundwater. While flow of groundwater towards the tunnel results in mixing of some  
61 groundwaters, it also flushes away any anthropogenic contamination introduced into the  
62 deep biosphere. In addition, boreholes drilled far into the bedrock means that oxygen does  
63 not penetrate to the sampled fracture waters and enable waters to be sampled under *in situ*  
64 conditions. Finally, the risk of influencing the microbial community due to materials used to  
65 close the borehole sections<sup>43</sup> is minimized by flushing three section volumes prior to sampling.

66 The island of Olkiluoto on the south-west coast of Finland will host a deep geological  
67 repository for the final disposal of spent nuclear fuel (Lat N 61° 14' 31", Lon E 21° 29' 23").  
68 Groundwater is accessed via deep drillholes at Olkiluoto fitted with multipackers that allow  
69 isolation of fracture fluids. Fractures were pumped for 2–3 months prior to microbiological

70 sampling. During this time, chemical parameters (including dissolved O<sub>2</sub> and oxidation-  
71 reduction potential) were monitored to ensure that the water was representative of the  
72 isolated fracture. Deep drillholes and excavation of the underground tunnels at Olkiluoto can  
73 form transient drawdown of groundwater in connected fractures due to changes in hydraulic  
74 head (pressure). This can also cause mixing of different groundwater-types at Olkiluoto, but  
75 deep groundwaters are highly reducing and there is no evidence of oxygen penetration<sup>44</sup>.  
76 Sterility and anoxic conditions during sampling of groundwaters at Olkiluoto, Finland, were  
77 ensured by filtering directly through 0.2 µm filters as described in<sup>33</sup>.

### 78 **Groundwater samples**

79 We collected groundwater samples from five different boreholes: SA1229A-1 (171.3  
80 mbsl), KA3105A-4 (415.2 mbsl), KA2198A (294.1 mbsl), KA3385A-1 (448.4 mbsl), and  
81 KF0069A01 (454.8 mbsl) with varying geochemical conditions as described below. A total of  
82 27 metagenomes were generated from the respective Äspö HRL boreholes that are listed in  
83 Supplementary Data 1. The biofilm metagenomes were formed on rock and glass surfaces  
84 after 33 days in flow cells attached to the KA2198A and KF0069A01 boreholes (total  $n = 8$ ) as  
85 described in Wu et al<sup>11</sup>. Further metagenomes were generated from planktonic cells captured  
86 on 0.22 µm filters from boreholes SA1229A, KA3105-4, and KA3385A (total  $n = 6$ ; termed 'large  
87 cells') and cells that passed through the 0.22 µm filters from same three boreholes ( $n = 9$ ;  
88 termed 'small cells') as previously described<sup>32</sup>. Finally, planktonic cells captured on 0.1 µm  
89 filters from boreholes SA1229A and KA3385A that were extracted in this study using the  
90 MOBIO PowerWater DNA Isolation Kit or phenol-chloroform ( $n = 4$ ).  
91 These groundwaters carried iron as Fe<sup>2+</sup>, contain dissolved sulfide (HS<sup>-</sup>), had temporally stable  
92 chemistry and δ<sup>18</sup>O, and neutral pH<sup>32</sup>. However, differential chemical composition and δ<sup>18</sup>O  
93 values enable us to characterize the origin and age of these groundwaters<sup>45</sup>. Groundwaters

94 SA1229A-1, KA3105A-4, and KA2198A have stable chloride concentrations and  $\delta^{18}\text{O}$  values,  
95 similar to those corresponding to Baltic Sea water<sup>46</sup>. Consequently, these groundwaters have  
96 a marine signature and are most likely composed of Baltic Sea water mixed with minor  
97 proportions of meteoric water and/or older more saline water residing in the bedrock  
98 fractures<sup>47</sup>. The precise infiltration age of this groundwater was unknown, but estimated to  
99 be <20 years or even less<sup>11</sup>. They were termed 'modern marine' ('MM') waters, concretely  
100 'MM-171.3' for SA1229A-1, 'MM-294.1' for KA2198A, and 'MM-415.2' for KA3105A-4  
101 groundwaters. Borehole KA3385A-1 contained water with chloride concentrations and  $\delta^{18}\text{O}$   
102 values in between the groundwaters with saline signature with a very long residence time<sup>48</sup>  
103 and marine groundwaters. Therefore, this groundwater was classified as thoroughly mixed  
104 ('TM-448.4') and is composed of unknown proportions of two or more water types such that  
105 the age of this groundwater cannot be assessed<sup>25</sup>. Borehole KF0069A01 had a typical signature  
106 of low dissolved organic carbon and other anions, high chloride and sulfate concentrations  
107 derived from mineral weathering, an age of millions of years<sup>11</sup>, and was defined as 'old saline'  
108 ('OS-454.8'). The MM-171.3, MM-415.2, and TM-448.4 boreholes were sampled from  
109 November to December 2016 while MM-294.1 and OS-454.8 were both sampled between  
110 May to June 2013 (Supplementary Data 1).

111 **Olkiluoto, Finland.** Groundwater was collected from three drillholes that access fracture  
112 fluids at different depths; OL-KR11 (366.7-383.5 mbsl), OL-KR13 (330.5-337.9 mbsl), and OL-  
113 KR46 (528.7-531.5 mbsl). Multiple samples were collected during 2016 (OL-KR11  $n = 7$ , OL-  
114 KR13  $n = 7$  and OL-KR46  $n = 3$ ). Geochemical parameters of the groundwater were monitored  
115 throughout the sampling period and they are available in Supplementary Data 1. At Olkiluoto,  
116 the groundwater chemistry is stratified with depth. Salinity increases with depth and brackish  
117 sulfate-rich groundwater is found up to ~400 m depth, beyond which sulfate-free saline



118 groundwater dominates (Posiva 2013;  
119 [http://www.posiva.fi/en/databank/posiva\\_reports/olkiluoto\\_site\\_description\\_2011.1871.xh](http://www.posiva.fi/en/databank/posiva_reports/olkiluoto_site_description_2011.1871.xhtml#XnkY5C2ZNTY)  
120 [tml#XnkY5C2ZNTY](http://www.posiva.fi/en/databank/posiva_reports/olkiluoto_site_description_2011.1871.xhtml#XnkY5C2ZNTY)). OL-KR11 and OL-KR13 drillholes both access brackish groundwater types  
121 (residence time 2,500-8,500 years). Drillhole OL-KR46 accesses a deeper saline groundwater  
122 with a residence time >10,000 years<sup>44</sup>.

123

#### 124 **Biomass collection, DNA extraction, and metagenome sequencing**

125 **Äspö HRL.** Planktonic cells were collected after flushing five borehole section volumes on  
126 sterile polyvinylidene fluoride (PVDF), hydrophilic, 0.1 µm, 47 mm Durapore membrane filters  
127 (Merck Millipore) under *in situ* conditions by connecting a High-Pressure Stainless Steel Filter  
128 Holder (Millipore) with a downstream needle valve and pressure gauge directly to the  
129 borehole. After filtering an appropriate volume of groundwater, each filter was rolled and  
130 placed in a sterile cryogenic tube (Thermo Scientific) and immediately frozen in liquid  
131 nitrogen. Samples were frozen at the sampling site to allow transport to the laboratory  
132 without any changes in the microbial community. Tubes were stored at -80 °C until further  
133 processing. DNA of samples MM-171.3-PC and TM-448.4-PC were extracted using the  
134 phenol/chloroform/isoamyl alcohol (24:24:1) method using Phase Lock tubes (Eppendorf).  
135 Firstly, 840 µL of TE buffer, pH 8 plus 94 µL of lysozyme (100 mg/mL) were added to each filter  
136 before incubation at 37 °C for 30 min. Then, 60 µL of 10 % sodium dodecyl sulfate (SDS) and 6  
137 µL of proteinase K (20 mg/ml) were added, mixed, and incubated at 50 °C for 20 min.  
138 Afterwards, an equal volume of phenol/chloroform/isoamyl alcohol was added to the cell  
139 lysate, mixed by inverting, and transferred to a Phase Lock Gel tube before centrifugation at  
140 1500 × g for 10 min. Then, another equal volume of phenol/chloroform was added and mixed  
141 before centrifuging at 1500 × g for 10 min. The nucleic acid was precipitated by adding an

142 equal volume of ice-cold isopropanol and 0.1 volume of 3M sodium acetate, pH 5.2 and  
143 incubating at -20 °C for 60 min. After precipitation, the nucleic acids were centrifuged at 16000  
144 × g and 4 °C for 20 min. The supernatant was discarded, and the pellet was rinsed with 500 µL  
145 of cold 80 % ethanol. Finally, the pellet was dried at 55 °C on a heat block and re-suspended  
146 in 50 µL of TE buffer prior incubation overnight at 4°C. The next day the DNA was incubated  
147 at 70 °C for 10 min to help dissolve the last of the DNA. DNA samples termed MM-171.3-PW,  
148 MM-415.2-PW, and TM-448.4-PW were extracted using the MO BIO PowerWater DNA  
149 isolation kit, following the manufacturer's instructions except that the final DNA re-suspension  
150 was performed using 60 µL of eluent<sup>32</sup>. The quality and quantity of the extracted DNA by both  
151 methods were analyzed with a Thermo Scientific Nanodrop 2000 and Qubit 2.0 Fluorometer  
152 (Life Technologies), respectively. Extracted DNA was stored at -20 °C. Twenty-seven  
153 metagenomic datasets were generated from the samples collected from the Äspö HRL.  
154 Detailed statistics of the generated metagenomes and the respective sequencing platform are  
155 shown in Supplementary Data 1. DNA was extracted from the MM-294.1 and OS-454.8  
156 samples as explained in the reference<sup>11</sup>.

157 **Olkiluoto.** To collect biomass for DNA analysis, approximately 10 L of groundwater was  
158 pumped directly into a chilled sterile Nalgene filtration unit fitted with a 0.22 µm pore size  
159 Isopore polycarbonate membrane (Millipore) and connected to a vacuum pump. After  
160 filtration, the membrane filters were rolled and stored in 1.5 mL sterile screwcap tubes. Filters  
161 collected for DNA extraction were preserved in 750 mL LifeGuard Soil Preservation Solution  
162 (MoBio, Carlsbad, CA, United States) and transferred to the laboratory on dry ice. Filters were  
163 stored at -20 °C until further processing. DNA content was extracted using a phenol-  
164 chloroform protocol<sup>29</sup> with the following modifications. Firstly, filter pieces were subject to  
165 bead-beating (2 × 15 s) prior to incubation in lysozyme for 2 h at 37 °C and secondly, lysate

166 was incubated in Proteinase K (200 mg/mL final concentration) for 2 h. Extracted DNA was  
167 measured using the Qubit dsDNA High Sensitivity Assay kit (Thermo Fisher Scientific, Inc.). A  
168 total of 17 metagenomic datasets were generated from Olkiluoto, the statistics are shown in  
169 Supplementary Data 1.

170

### 171 **RNA extraction and metatranscriptome sequencing**

172 The groundwaters were sampled from the Äspö HRL under *in situ* conditions using two  
173 different sampling methods. Firstly, by connecting a sampling device with an in-built fixation  
174 system as described in Lopez-Fernandez et al.<sup>10</sup> from June 2015 to March 2016  
175 (Supplementary Data 1). Secondly, by connecting a high-pressure stainless steel filter holder  
176 (Merck Millipore, USA) with a downstream needle valve and pressure gauge as described in  
177 Lopez-Fernandez et al.<sup>25</sup> from September 2015 to January 2016 (Supplementary Data 1). In  
178 both cases, planktonic cells were collected on sterile hydrophilic polyvinylidene fluoride  
179 (PVDF) membranes with 0.1 µm poresize (47 mm Durapore, Merck Millipore, USA) under *in*  
180 *situ* conditions. The cell collection using both sampling methods, RNA extraction, and cDNA  
181 generation was performed as previously described<sup>10</sup>. Detailed statistics of the nine sequenced  
182 metatranscriptomes are shown in the Supplementary Data 1.

183

### 184 **Single cell collection and amplification**

185 The metagenomics samples were augmented with 564 single-cell amplified genomes (SAGs).  
186 The SAGs originate from MM-171.3 (borehole SA1229A-1,  $n=118$ ), MM-415.2 (borehole  
187 KA3105A-4,  $n=15$ ), and TM-448.4 (borehole KA3385A-1,  $n=148$ ) borehole samples from the  
188 Äspö HRL along with OL-KR11 ( $n=138$ ), OL-KR13 ( $n=117$ ), and OL-KR46 ( $n=28$ ) borehole

189 samples from the Olkiluoto. SAGs were amplified, sequenced, and assembled by the Joint  
190 Genome Institute (JGI), USA.

191 SAGs were de-replicated separately from the MAGs. SAGs that are containing the exact  
192 match of 16S rRNA and those with average nucleotide identity higher than 95% were  
193 combined in order to retrieve a higher number of good quality SAGs. These combined SAGs  
194 are referred to as several-SAG (s-SAG). A total of 22 SAGs were also sequenced from different  
195 water types of Äspö HRL at the SciLifeLab, Sweden as a pilot study. The SAGs were sequenced  
196 using the Illumina platform and assembled using MEGAHIT<sup>49</sup>.

197 SAGs were clustered separately from the MAGs using fastANI (v. 1.1) with 95% average  
198 nucleotide identity and 70% minimum overlap. Then, SAGs belonging to the same cluster were  
199 analyzed using the 'merge' command in checkm (v. 1.0.7) to find sets of within-cluster SAGs  
200 that could potentially be merged in order to increase the completeness. Initially, this resulted  
201 in nine pairs of SAGs where the estimated combined genome completeness would increase.  
202 GC-profiles were calculated for these 18 SAGs using 1 kbp sliding windows and similar profiles  
203 for each pair were validated by manual inspection. Next, redundancies within each pair were  
204 investigated by aligning contigs from SAGs with nucmer (v. 3.23) with default settings. Aligned  
205 regions were only kept on the longer contigs, by clipping the corresponding stretch from the  
206 shorter contigs. If clipping resulted in a contigs <300 bp, that contig was removed completely.  
207 In addition, if more than 25% of a contig aligned to another contig, the shorter contig was  
208 removed completely. After removing redundant regions, contigs from the SAG pairs were  
209 combined. All s-SAGs were again checked for completeness and contamination using Checkm  
210 and those with >5% contamination was kept as original SAGs. Detailed information regarding  
211 the SAGs is shown in Supplementary Data 1.

212

### 213 **Species richness using gene *gyrB***

214 To evaluate the species richness of each ecosystem captured by the metagenomes analyzed  
215 in this study, we extracted the *gyrB* genes from all metagenomic assemblies, evaluated the  
216 annotation by checking the conserved domains of the gene, and then clustered them at the  
217 97% and 88% identity threshold defined for this gene to reconstruct *gyrB* mOTUs<sup>39,40</sup> using  
218 CD-hit<sup>41</sup> (Supplementary Fig. S3).

219

### 220 **Fennoscandian Shield genomic database (FSGD)**

221 The generated “multi-omics” data were used to construct a comprehensive genomic and  
222 metatranscriptomic database of different water-types of the extremely oligotrophic deep  
223 groundwater.

224 The sequenced metagenomes were quality checked and trimmed using Trimmomatic<sup>50</sup>  
225 (v. 0.36) with settings to trim the Illumina TruSeq adapter (‘TruSeq3-PE-2.fa:2:30:15  
226 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:31’). For the six samples sequenced on  
227 the MiSeq platform, reads were first cropped to 125 bp by trimming the right end of the reads  
228 prior to adapter and quality trimming as above. This was done to recover more paired-end  
229 reads from these samples that all had lower quality bases at the end of reverse reads. Each  
230 dataset was assembled separately as well as co-assemblies on those datasets originating from  
231 the same water type in each sampling site using MEGAHIT<sup>49</sup> (v. 1.1) with settings (--k-min 21 -  
232 -k-max 141 --k-step 12 --min-count 2). Contigs  $\geq 2$ kb in each assembly were automatically  
233 binned using metabat2<sup>51</sup> with default setting. CheckM<sup>52</sup> was used to estimate the genome  
234 completeness of MAGs, SAGs, and s-SAGs. Those with completeness  $\geq 50\%$  and contamination  
235  $\leq 5\%$  were considered for down-stream analysis. In addition, SAGs with  $< 50\%$  completeness



236 were considered for down-stream analysis if they clustered with another SAG in the SAG-  
237 specific fastANI step (see above) and if they had a genome size  $\geq 500$  kbp.

238 **Genome de-replication.** MAGs and SAGs were clustered using fastANI<sup>53</sup> (v. 1.1) at  $\geq 95\%$   
239 identity and  $\geq 70\%$  coverage threshold. Those genomes in a single cluster are considered as  
240 representatives of a single population.

241 **Genome taxonomy and phylogeny.** Taxonomic affiliation of the FSGD MAGs/SAGs was  
242 assigned using GTDB-tk (v. 0.2.2) with reference to the release 86 database<sup>54</sup>. The alignments  
243 generated by the GTDB-tk for bacteria and archaea were curated and used for phylogeny  
244 reconstruction using FastTree<sup>55</sup> (v. 2.1.10) with parameters '-wag -gamma'.

245 **Gene annotation and functional analysis.** Prodigal<sup>56</sup> (v. 2.6.2) was run in metagenomic  
246 mode ('-p meta') for predicting protein-coding genes in the assembled contigs. This was  
247 followed by functional annotation of the predicted proteins using eggno-mapper<sup>57</sup> (v. 2.2.1)  
248 with the eggno-g\_5.0 database, and pfam\_scan.pl (v. 1.6) with the 31.0 release of the PFAM  
249 database. Reconstructed MAGs and SAGs were initially annotated using Prokka<sup>58</sup> (v. 1.12)  
250 followed by further annotation with eggno-mapper and pfam\_scan.pl using the same  
251 databases as for the metagenomic assemblies. Enzyme EC numbers, and KEGG orthologs,  
252 pathways and modules were assigned from the eggno-mapper output. All annotations of key  
253 genes were manually inspected for their conserved domains and their annotations were  
254 further evaluated using NLM's Conserved Domain Database (CDD) search<sup>59</sup> and phylogeny.

255 **Genome presence/absence patterns.** Metagenomics reads were mapped against all  
256 MAGs/SAGs that passed the criteria for downstream analysis using bowtie2<sup>60</sup> (v. 2.3.3.1) with  
257 parameters '--very-sensitive --no-unal'. This was followed by removal of duplicates using  
258 MarkDuplicates from the picard suite (v. 2.18.6). Only contigs that  $\geq 50\%$  of their length was  
259 covered by the mapped reads were considered for further analysis. Mapped reads were

260 counted using featureCounts<sup>61</sup> (v. 1.6.1) with settings '-M -B' to count multi-mapping and only  
261 count read-pairs with both ends aligned. The raw counts were normalized as transcripts per  
262 million (TPM) in order to calculate MAGs/SAGs abundance in each metagenome. Based on the  
263 calculated average TPM per contig; the MAGs/SAGs were considered detected in the  
264 metagenome if they show value  $\geq 1$  and not detected if the value is  $< 1$ . These strict mapping  
265 thresholds identify closely related isolates to the reconstructed MAGs/SAGs.

266 **Computation of Isoelectric point and codon usage frequency.** The isoelectric point  
267 calculation for the protein sequences as well as amino acid features were calculated using  
268 pepstat software in the EMBOSS package (v. 6.6.0)<sup>62</sup>. The codon usage frequency of the coding  
269 regions was calculated using software cusp in the EMBOSS package (v. 6.6.0)<sup>62</sup>.

270 **DnaE2 phylogeny.** The phylogeny of the C-family polymerases was reconstructed by using  
271 the reference genomes are the reviewed sequences retrieved from uniprot for each  
272 polymerase type. The annotation of the protein coding sequences with evaluated annotation  
273 as DnaE2 was verified by using this phylogeny. Sequences were aligned using Kalign<sup>63</sup> (2.04)  
274 and FastTree (v. 2.1.10) was used for creating the maximum-likelihood tree (JTT +CAT model,  
275 gamma approximation).

276 **Dissimilatory sulfur metabolism.** The phylogeny of DsrA as a key gene in the dissimilatory  
277 sulfur metabolism was generated by using reference *dsrA* genes (both oxidative and reductive  
278 types) together with the genes annotated as *dsrA* in the reconstructed MAGs/SAGs of our  
279 study ( $\geq 200$  amino acids length). Sequences were aligned using Muscle<sup>64</sup> in MEGA7<sup>65</sup> and  
280 evolutionary relationships were visualized by constructing a maximum-likelihood  
281 phylogenetic tree (JTT +CAT model). All residues were used, and the tree was bootstrapped  
282 with 100 replicates. MAGs and SAGs containing *dsrA* gene were further inspected for the  
283 presence of *aprAB* (K00394 and K00395), *sat* (K00958), *dsrAB* (K11180 and K11181), *dsrC*

284 (K11179), *dsrD* (PF08679), and *dsrEFH* (K07235, K07236, and K07237). These genes were  
285 searched for using eggno-mapper<sup>57</sup> (v. 2.2.1 with the eggno-5.0 database) and  
286 pfam\_scan.pl (v. 1.6 with the 31.0 release of the PFAM database) and annotations were  
287 manually validated for each gene. The contribution of each MAG/SAG to the sulfur cycle was  
288 inferred according to the pattern of presence/absence of these genes as suggested by  
289 Anantharaman et. al<sup>66</sup>. For the final inference members of each cluster were considered  
290 together.

291 **Metatranscriptome analysis.** The sequenced metatranscriptomes were quality checked  
292 and trimmed using Trimmomatic (v. 0.36)<sup>50</sup>. The rRNA reads were filtered out using cmsearch  
293 (v. 1.1.3)<sup>67</sup>. The remaining reads were mapped against the FSGD MAGs/SAGs. The expressed  
294 genetic content of each MAG/SAG were extracted at the threshold of 100 TPM and the pattern  
295 of gene expression and the expressed content of each MAG/SAG was analyzed.

296

## 297 **References**

- 298 1. Magnabosco, C. *et al.* The biomass and biodiversity of the continental subsurface. *Nat.*  
299 *Geosci.* **11**, 707–717 (2018).
- 300 2. Colman, D. R., Poudel, S., Stamps, B. W., Boyd, E. S. & Spear, J. R. The deep, hot  
301 biosphere: Twenty-five years of retrospection. *Proc. Natl. Acad. Sci.* **114**, 6895–6903  
302 (2017).
- 303 3. Cario, A., Oliver, G. C. & Rogers, K. L. Exploring the Deep Marine Biosphere: Challenges,  
304 Innovations, and Opportunities. *Front. Earth Sci.* **7**, 1–9 (2019).
- 305 4. Sackett, J. D. *et al.* Four draft Single-Cell genome sequences of novel , nearly identical  
306 Kiritimatiellaeta strains isolated from the continental deep subsurface. *Microbiol*  
307 *Resour Announc* e01249-18 (2019).

- 308 5. Anantharaman, K. *et al.* Thousands of microbial genomes shed light on interconnected  
309 biogeochemical processes in an aquifer system. *Nat. Commun.* **7**, 1–11 (2016).
- 310 6. Flemming, H.-C. & Wuertz, S. Bacteria and archaea on Earth and their abundance in  
311 biofilms. *Nat. Rev. Microbiol.* **17**, 247–260 (2019).
- 312 7. Magnabosco, C. *et al.* A metagenomic window into carbon metabolism at 3 km depth  
313 in Precambrian continental crust. *ISME J.* **10**, 730–741 (2016).
- 314 8. Al-Shayeb, B. *et al.* Clades of huge phages from across Earth’s ecosystems. *Nature*  
315 (2020) doi:10.1038/s41586-020-2007-4.
- 316 9. Schwank, K. *et al.* An archaeal symbiont-host association from the deep terrestrial  
317 subsurface. *ISME J.* **13**, 2135–2139 (2019).
- 318 10. Lopez-Fernandez, M. *et al.* Metatranscriptomes reveal that all three domains of life are  
319 active but are dominated by Bacteria in the Fennoscandian crystalline granitic  
320 continental deep biosphere. *MBio* **9**, 1–15 (2018).
- 321 11. Wu, X. *et al.* Potential for hydrogen-oxidizing chemolithoautotrophic and diazotrophic  
322 populations to initiate biofilm formation in oligotrophic, deep terrestrial subsurface  
323 waters. *Microbiome* **5**, 1–13 (2017).
- 324 12. Daly, R. A. *et al.* Microbial metabolisms in a 2.5-km-deep ecosystem created by  
325 hydraulic fracturing in shales. *Nat. Microbiol.* **1**, 16146 (2016).
- 326 13. Engelhardt, T., Kallmeyer, J., Cypionka, H. & Engelen, B. High virus-to-cell ratios indicate  
327 ongoing production of viruses in deep subsurface sediments. *ISME J.* **8**, 1503–1509  
328 (2014).
- 329 14. Borgonie, G. *et al.* Eukaryotic opportunists dominate the deep-subsurface biosphere in  
330 South Africa. *Nat. Commun.* **6**, (2015).
- 331 15. Jørgensen, B. B. Shrinking majority of the deep biosphere. *Proc. Natl. Acad. Sci. U. S. A.*

- 332           **109**, 15976–15977 (2012).
- 333 16. Li, J. *et al.* Recycling and metabolic flexibility dictate life in the lower oceanic crust.  
334           *Nature* **579**, 250–255 (2020).
- 335 17. Heard, A. W. *et al.* South African crustal fracture fluids preserve paleometeoric water  
336           signatures for up to tens of millions of years. *Chem. Geol.* **493**, 379–395 (2018).
- 337 18. Robador, A. *et al.* Nanocalorimetric characterization of microbial activity in deep  
338           subsurface oceanic crustal fluids. *Front. Microbiol.* **7**, (2016).
- 339 19. Bradley, J. A., Amend, J. P. & LaRowe, D. E. Survival of the fewest: Microbial dormancy  
340           and maintenance in marine sediments through deep time. *Geobiology* **17**, 43–59  
341           (2019).
- 342 20. Wasmund, K., Mußmann, M. & Loy, A. The life sulfuric: microbial ecology of sulfur  
343           cycling in marine sediments. *Environ. Microbiol. Rep.* **9**, 323–344 (2017).
- 344 21. Momper, L. *et al.* Major phylum-level differences between porefluid and host rock  
345           bacterial communities in the terrestrial deep subsurface. *Environ. Microbiol. Rep.* **9**,  
346           501–511 (2017).
- 347 22. Suko, T. *et al.* Geomicrobiological properties of Tertiary sedimentary rocks from the  
348           deep terrestrial subsurface. *Phys. Chem. Earth* **58–60**, 28–33 (2013).
- 349 23. Arbour, T. J., Gilbert, B. & Banfield, J. F. Diverse Microorganisms in Sediment and  
350           Groundwater Are Implicated in Extracellular Redox Processes Based on Genomic  
351           Analysis of Bioanode Communities. *Front. Microbiol.* **11**, (2020).
- 352 24. Borgonie, G. *et al.* New ecosystems in the deep subsurface follow the flow of water  
353           driven by geological activity. *Sci. Rep.* **9**, 1–16 (2019).
- 354 25. Lopez-Fernandez, M., Åström, M., Bertilsson, S. & Dopson, M. Depth and Dissolved  
355           Organic Carbon Shape Microbial Communities in Surface Influenced but Not Ancient



- 356 Saline Terrestrial Aquifers. *Front. Microbiol.* **9**, 1–16 (2018).
- 357 26. Cai, L. *et al.* Active and diverse viruses persist in the deep sub-seafloor sediments over  
358 thousands of years. *ISME J.* **13**, 1857–1864 (2019).
- 359 27. Lopez-Fernandez, M., Broman, E., Simone, D., Bertilsson, S. & Dopson, M. Statistical  
360 analysis of community RNA transcripts between organic carbon and geogas-fed  
361 continental deep biosphere groundwaters. *MBio* **10**, e01470-19 (2019).
- 362 28. Lopez-Fernandez, M. *et al.* Investigation of viable taxa in the deep terrestrial biosphere  
363 suggests high rates of nutrient recycling. *FEMS Microbiol. Ecol.* **94**, 1–9 (2018).
- 364 29. Bagnoud, A. *et al.* Reconstructing a hydrogen-driven microbial metabolic network in  
365 Opalinus Clay rock. *Nat. Commun.* **7**, 1–10 (2016).
- 366 30. Lau, M. C. Y. *et al.* An oligotrophic deep-subsurface community dependent on syntrophy  
367 is dominated by sulfur-driven autotrophic denitrifiers. *Proc. Natl. Acad. Sci. U. S. A.* **113**,  
368 E7927–E7936 (2016).
- 369 31. Lever, M. A. *et al.* Life under extreme energy limitation: A synthesis of laboratory- and  
370 field-based investigations. *FEMS Microbiol. Rev.* **39**, 688–728 (2015).
- 371 32. Wu, X. *et al.* Microbial metagenomes from three aquifers in the Fennoscandian shield  
372 terrestrial deep biosphere reveal metabolic partitioning among populations. *ISME J.* **10**,  
373 1192–1203 (2016).
- 374 33. Bell, E. *et al.* Active sulfur cycling in the terrestrial deep subsurface. *ISME J.* **14**, 1260–  
375 1272 (2020).
- 376 34. Drake, H., Åström, M. E., Tullborg, E. L., Whitehouse, M. & Fallick, A. E. Variability of  
377 sulphur isotope ratios in pyrite and dissolved sulphate in granitoid fractures down to  
378 1km depth - Evidence for widespread activity of sulphur reducing bacteria. *Geochim.*  
379 *Cosmochim. Acta* **102**, 143–161 (2013).

- 380 35. Alakangas, L. J., Mathurin, F. A. & Åström, M. E. Diverse fractionation patterns of Rare  
381 Earth Elements in deep fracture groundwater in the Baltic Shield – Progress from  
382 utilisation of Diffusive Gradients in Thin-films (DGT) at the Äspö Hard Rock Laboratory.  
383 *Geochim. Cosmochim. Acta* **269**, 15–38 (2020).
- 384 36. Drake, H., Tullborg, E. L., Hogmalm, K. J. & Åström, M. E. Trace metal distribution and  
385 isotope variations in low-temperature calcite and groundwater in granitoid fractures  
386 down to 1km depth. *Geochim. Cosmochim. Acta* **84**, 217–238 (2012).
- 387 37. Kärki, A. & Paulamäki, S. *Petrology of Olkiluoto. Posiva Report 2006-02. Posiva Oy,*  
388 *Olkiluoto* vol. 2  
389 [http://www.posiva.fi/en/databank/posiva\\_reports/petrology\\_of\\_olkiluoto.1871.xhtm](http://www.posiva.fi/en/databank/posiva_reports/petrology_of_olkiluoto.1871.xhtm)  
390 [l?xm\\_freetext=petrology&xm\\_col\\_type=4&cd\\_order=col\\_report\\_number&cd\\_offset=](http://www.posiva.fi/en/databank/posiva_reports/petrology_of_olkiluoto.1871.xhtml?xm_freetext=petrology&xm_col_type=4&cd_order=col_report_number&cd_offset=0#.X44JgS-Q3OQ)  
391 [0#.X44JgS-Q3OQ](http://www.posiva.fi/en/databank/posiva_reports/petrology_of_olkiluoto.1871.xhtml?xm_freetext=petrology&xm_col_type=4&cd_order=col_report_number&cd_offset=0#.X44JgS-Q3OQ) (2006).
- 392 38. Pitkanen, P., Partamies, S. & Luukkonen, A. *Hydrogeochemical Interpretation of*  
393 *Baseline Groundwater Conditions at the Olkiluoto Site. Posiva Report 2003-07, Posiva*  
394 *Oy, Olkiluoto.*  
395 [http://www.posiva.fi/en/databank/posiva\\_reports/hydrochemical\\_interpretation\\_of\\_](http://www.posiva.fi/en/databank/posiva_reports/hydrochemical_interpretation_of)  
396 [baseline\\_groundwater\\_conditions\\_at\\_the\\_olkiluoto\\_site.1871.xhtml?xm\\_freetext=Hy](http://www.posiva.fi/en/databank/posiva_reports/hydrochemical_interpretation_of)  
397 [drogeochemical+Interpretation+of+Baseline+Groundwater+Conditions+at+the+Olkilu](http://www.posiva.fi/en/databank/posiva_reports/hydrochemical_interpretation_of)  
398 [oto+Site&xm\\_col\\_ty](http://www.posiva.fi/en/databank/posiva_reports/hydrochemical_interpretation_of) (2003).
- 399 39. Poirier, S. *et al.* Deciphering intra-species bacterial diversity of meat and seafood  
400 spoilage microbiota using gyrB amplicon sequencing: A comparative analysis with 16S  
401 rDNA V3-V4 amplicon sequencing. *PLoS One* **13**, 1–26 (2018).
- 402 40. Caro-Quintero, A. & Ochman, H. Assessing the unseen bacterial diversity in microbial  
403 communities. *Genome Biol. Evol.* **7**, 3416–3425 (2015).

- 404 41. Li, W. & Godzik, A. Cd-hit: A fast program for clustering and comparing large sets of  
405 protein or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).
- 406 42. Hallbeck, L. & Pedersen, K. Characterization of microbial processes in deep aquifers of  
407 the Fennoscandian Shield. *Appl. Geochemistry* **23**, 1796–1819 (2008).
- 408 43. Drake, H. *et al.* Extreme fractionation and micro-scale variation of sulphur isotopes  
409 during bacterial sulphate reduction in deep groundwater systems. *Geochim.*  
410 *Cosmochim. Acta* **161**, 1–18 (2015).
- 411 44. Posiva Oy. *Olkiluoto Site Description 2011*. vol. 31 (2011).
- 412 45. Laaksoharju, M., Gascoyne, M. & Gurban, I. Understanding groundwater chemistry  
413 using mixing models. *Appl. Geochemistry* **23**, 1921–1940 (2008).
- 414 46. Mathurin, F. A., Åström, M. E., Laaksoharju, M., Kalinowski, B. E. & Tullborg, E. L. Effect  
415 of tunnel excavation on source and mixing of groundwater in a coastal granitoidic  
416 fracture network. *Environ. Sci. Technol.* **46**, 12779–12786 (2012).
- 417 47. Laaksoharju, M. *et al.* Hydrogeochemical evaluation and modelling performed within  
418 the Swedish site investigation programme. *Appl. Geochemistry* **23**, 1761–1795 (2008).
- 419 48. Louvat, D., Luc, J. & Franc, J. È spo Origin and residence time of salinity in the A  
420 groundwater system. *Appl. Geochemistry* **14**, 917–925 (1999).
- 421 49. Li, D., Liu, C.-M. M., Luo, R., Sadakane, K. & Lam, T.-W. W. MEGAHIT: An ultra-fast single-  
422 node solution for large and complex metagenomics assembly via succinct de Bruijn  
423 graph. *Bioinformatics* **31**, 1674–1676 (2014).
- 424 50. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic : a flexible trimmer for Illumina  
425 sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
- 426 51. Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for accurately  
427 reconstructing single genomes from complex microbial communities. *PeerJ* **3**, e1165

- 428 (2015).
- 429 52. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:  
430 assessing the quality of microbial genomes recovered from isolates, single cells, and  
431 metagenomes. *Genome Res.* **25**, 1043–55 (2015).
- 432 53. Jain, C., Rodriguez-r, L. M. & Aluru, S. High throughput ANI analysis of 90K prokaryotic  
433 genomes reveals clear species boundaries. *Nat. Commun.* **9**, 5114 (2018).
- 434 54. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny  
435 substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996 (2018).
- 436 55. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2--approximately maximum-likelihood  
437 trees for large alignments. *PLoS One* **5**, e9490 (2010).
- 438 56. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site  
439 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 440 57. Huerta-Cepas, J. *et al.* Fast genome-wide functional annotation through orthology  
441 assignment by eggNOG-mapper. *Mol. Biol. Evol.* **34**, 2115–2122 (2017).
- 442 58. Seemann, T. Prokka: Rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–  
443 2069 (2014).
- 444 59. Lu, S. *et al.* CDD/SPARCLE: The conserved domain database in 2020. *Nucleic Acids Res.*  
445 **48**, D265–D268 (2020).
- 446 60. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*  
447 **9**, 357–359 (2012).
- 448 61. Liao, Y., Smyth, G. K. & Shi, W. FeatureCounts: An efficient general purpose program for  
449 assigning sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
- 450 62. Rice, P., Longden, I. & Bleasby, A. EMBOSS: The European Molecular Biology Open  
451 Software Suite. *Trends Genet.* **16**, 276–277 (2000).

- 452 63. Lassmann, T. & Sonnhammer, E. L. L. Kalign--an accurate and fast multiple sequence  
453 alignment algorithm. *BMC Bioinformatics* **6**, 298 (2005).
- 454 64. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high  
455 throughput. *Nucleic Acid Res.* **32**, 1792–1797 (2004).
- 456 65. Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular Evolutionary Genetics Analysis  
457 Version 7.0 for Bigger Datasets. *Mol. Biol. Evol.* **33**, 1870–1874 (2016).
- 458 66. Anantharaman, K. *et al.* Expanded diversity of microbial groups that shape the  
459 dissimilatory sulfur cycle. *ISME J.* **12**, 1715–1728 (2018).
- 460 67. Nawrocki, E. P. & Eddy, S. R. Infernal 1.1: 100-fold faster RNA homology searches.  
461 *Bioinformatics* **29**, 2933–2935 (2013).
- 462