1  **Title: The economical lifestyle of CPR bacteria in groundwater allows little**

2  **preference for environmental drivers**

3  **Authors:** Narendrakumar M. Chaudhari[1,2], Will A. Overholt[1], Perla Abigail Figueroa-Gonzalez[3],

4  Martin Taubert[1], Till L. V. Bornemann[3], Alexander J. Probst[3], Martin Hölzer[4,5†], Manja Marz[4,5,6]

5  and Kirsten Küsel[1,2*]

6  **Institutions:**

7  [1]Aquatic Geomicrobiology, Institute of Biodiversity, Friedrich Schiller University, Jena,

8  Germany.

9  [2]German Center for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Leipzig,

10  Germany.

11  [3]Department for Chemistry, Environmental Microbiology and Biotechnology, Group for Aquatic

12  Microbial Ecology (GAME), University Duisburg-Essen, Essen, Germany.

13  [4]RNA Bioinformatics and High Throughput Analysis, Friedrich Schiller University, Jena,

14  Germany.

15  [5]European Virus Bioinformatics Center, Friedrich Schiller University, Jena, Germany.

16  [6]FLI Leibniz Institute for Age Research, Jena, Germany.

17  [†]Present address: Methodology and Research Infrastructure, MF1 Bioinformatics, Robert Koch

18  Institute, Berlin, Germany.

19  **Emails:**

20  Narendrakumar M. Chaudhari: narendrakumar.chaudhari@uni-jena.de

21  Will A. Overholt: will.overholt@uni-jena.de

22  Perla Abigail Figueroa-Gonzalez: abigail.figueroa-gonzalez@uni-due.de

23  Martin Taubert: martin.taubert@uni-jena.de

24  Till L. V. Bornemann: till.bornemann@uni-due.de

25  Alexander J. Probst: alexander.probst@uni-due.de

26  Martin Hölzer: hoelzer.martin@gmail.com

27  Manja Marz: manja@uni-jena.de

28  [*]Kirsten Küsel: kirsten.kuesel@uni-jena.de

29  [*]Corresponding author.

## Abstract

31   The highly diverse *Cand*. Patescibacteria are predicted to have minimal biosynthetic and

32   metabolic pathways, which hinders understanding of how their populations differentiate to

33   environmental drivers or host organisms. Their metabolic traits to cope with oxidative stress are

34   largely unknown. Here, we utilized genome-resolved metagenomics to investigate the adaptive

35   genome repertoire of Patescibacteria in oxic and anoxic groundwaters, and to infer putative host

36   ranges.

37   Within six groundwater wells, *Cand*. Patescibacteria was the most dominant (up to 79%) super-

38   phylum across 32 metagenomes obtained from sequential 0.2 and 0.1 μm filtration. Of the

39   reconstructed 1275 metagenome-assembled genomes (MAGs), 291 high-quality MAGs were

40   classified as *Cand*. Patescibacteria. *Cand*. Paceibacteria and *Cand*. Microgenomates were

41   enriched exclusively in the 0.1 μm fractions, whereas candidate division ABY1 and *Cand*.

42   Gracilibacteria were enriched in the 0.2 μm fractions. Patescibacteria enriched in the smaller 0.1

43   μm filter fractions had 22% smaller genomes, 13.4% lower replication measures, higher fraction

44   of rod-shape determining proteins, and genomic features suggesting type IV pili mediated cell-

45   cell attachments. Near-surface wells harbored Patescibacteria with higher replication rates than

46   anoxic downstream wells characterized by longer water residence time. Except prevalence of

47   superoxide dismutase genes in Patescibacteria MAGs enriched in oxic groundwaters (83%), no

48   major metabolic or phylogenetic differences were observed based on oxygen concentrations. The

49   most abundant Patescibacteria MAG in oxic groundwater encoded a nitrate transporter, nitrite

50   reductase, and F-type ATPase, suggesting an alternative energy conservation mechanism.

51   Patescibacteria consistently co-occurred with one another or with members of phyla

52    Nanoarchaeota, Bacteroidota, Nitrospirota, and Omnitrophota. However, only 8% of MAGs

53    showed highly significant one-to-one association, mostly with Omnitrophota. Genes coding for

54    motility and transport functions in certain Patescibacteria were highly similar to genes from other

55    phyla (Omnitrophota, Proteobacteria and Nanoarchaeota).

56    Other than genes to cope with oxidative stress, we found little genomic evidence for niche

57    adaptation of Patescibacteria to oxic or anoxic groundwaters. Given that we could detect specific

58    host preference only for a few MAGs, we propose that the majority of Patescibacteria can attach

59    to multiple hosts just long enough to loot or exchange supplies with an economic lifestyle of

60    little preference for geochemical conditions.

61    **Keywords**

62    Candidate Phyla Radiation (CPR), *Cand.* Patescibacteria, Economic lifestyle, Metagenomics,

63    Microbial ecology.

64    # Introduction

65    Metagenomic sequencing of diverse environments has enabled the recovery of genomic

66    information from a vast majority of uncultivated microbial dark matter, significantly expanding

67    the tree of life. *Cand.* Patescibacteria is a superphylum also known as Candidate Phyla Radiation

68    (CPR) that constitutes a major portion of this expanded tree of life [1]. Patescibacteria, initially

69    recovered from groundwater and aquatic sediments [2,3], are now shown to inhabit a broad range

70    of surface and subsurface habitats, such as marine water, freshwater, freshwater beach sands [4]

71    hydrothermal vents [5], cold-water geyser [6,7], plant rhizosphere [8], alpine permafrost [9],

72    permafrost thaw ponds [10], and many more habitats [11] including the human oral cavity [12–

3

73    14]. Nevertheless, they dominate the groundwater, where they comprise 20-70% of the total

74    microbial community [15–18] along with thermokarst lakes [19] and hypersaline soda lake

75    sediments [20].

76    Patescibacteria have small genomes characterized by predicted minimal biosynthetic and

77    metabolic pathways, and are reported to have an anaerobic, fermentative lifestyle [21,22]. These

78    traits may be responsible for their high abundance in nutrient-limited groundwater habitats,

79    which are mainly anoxic. Interestingly, oxic surface soils are a major source of CPR bacteria

80    inhabiting modern groundwater (stored within last 50 years) [23], as these organisms are easily

81    mobilized into soil seepage water [17,24], but their metabolic traits to cope with oxidative stress

82    are largely unknown. Divergent trends in the preference for several hydrochemical parameters or

83    specific host preferences seem to result in the differentiation of CPR bacteria in groundwater

84    [17]. Similarly, little species-level overlap of metagenome-assembled genomes (MAGs) across

85    varying groundwater sites suggests that CPR communities differ based on specific environmental

86    factors including host populations [18].

87    Most Patescibacteria cells are estimated to have ultra-small diameters ranging from 0.1 μm to 0.3

88    μm [11,15,21] with few exceptions like Saccharimonadia (candidate division TM7) that may be

89    as large as 0.7 μm in diameter [25]. Small cell sizes of Patescibacteria accompanied by reduced

90    genomes [3,21,22] suggest host-associated lifestyles. Indeed, specific studies on Patescibacteria

91    isolates along with co-culture and microscopic analyses provided evidence of their symbiotic

92    associations with other organisms e.g. with *Paramecium bursaria*, a ciliated protist in freshwater

93    [26], or with Actinobacteria (*Actinomyces odontolyticus*, *Propionibacterium propionicus*,

94    *Schaalia meyeri*) in the human oral cavity [12,27–29]. Similarly, CPR bacteria attach as

4

95    episymbionts to putative bacterial hosts through pilin-like appendages in pristine groundwater

96    [18].

97    In contrast, single cell genomic and biophysical observations from 46 globally distributed

98    groundwater sites did not support the prevailing view that Patescibacteria are dominated by

99    symbionts [11]. The authors suggest that their unusual genomic features and prevalent

100   auxotrophies may be the result of ancestral, primitive energy metabolism that relies on

101   fermentation. Additionally, genome streamlining in free-living prokaryotes in the open ocean is a

102   known mechanism to reduce functional redundancy and conserve energy [30]. Minimizing

103   energy expenditure and nutrient demands has constituted a selective advantage for

104   *Prochlorococcus* in surface waters where nutrients are scarce at the expense of versatility and

105   competitiveness in changing conditions [31], and the same could be true for CPR bacteria

106   dominating oligotrophic subsurface waters. Thus, there is the need to disentangle which lineages

107   of CPR bacteria are host-dependent and which are free-living, and how much variation in terms

108   of lifestyle, metabolism and gene content exists between those which show a preference for

109   certain geochemical conditions.

110   In this study, we took advantage of a well-studied modern groundwater system within the

111   Hainich Critical Zone Exploratory (CZE) located in Thuringia, Germany [32], dominated by

112   CPR bacteria, that exhibits large environmental gradients from oxic to anoxic conditions

113   accompanied by different well-specific microbiomes [33]. Using 291 manually curated MAGs

114   we aimed to identify the adaptive genomic repertoire of CPR bacteria. Sequential filtration was

115   performed to gather clues about possible physical association of ultra-small Patescibacteria with

116   larger sized host ranges. We also inferred putative hosts for Patescibacteria based on the co-

117    occurrence patterns with other microorganisms within the transect, especially based on

118    abundances of all the MAGs enriched in the 0.2 μm filter fractions.


# Results

## Patescibacteria represent more than 50% of all prokaryotes in Hainich groundwater

121    *Cand*. Patescibacteria dominated the groundwater community representing on average more than

122    $50 \pm 18\%$ (range 23-79%) prokaryotes across 32 metagenomes obtained from groundwater of six

123    wells that was sequentially filtered through 0.2 μm and 0.1 μm filters, based on the proportion of

124    the quality-controlled metagenomic reads mapped to the 16S rRNA database (SILVA SSU

125    rRNA Ref NR99) [34]. Three major classes within the phylum were detected: *Cand*.

126    Parcubacteria/Paceibacteria ($36.2 \pm 17.1\%$, range 13-65.7%), *Cand*. Microgenomatia ($7.2 \pm$

127    $3.1\%$, range 2-12%) and candidate division ABY1 ($3.2 \pm 1.4\%$, range 1.1-5.1%). Patescibacteria

128    were found to be highly abundant in both filter fractions. Their relative abundances were

129    significantly higher (two-proportions z-test, p-value 1.16e-05) in the 0.1 μm filter fractions (67.6

130    $\pm 9.1\%$, range 54.1-78.5%) than in the 0.2 μm filter fractions ($35.5 \pm 8.9\%$, range 23.1-51.1%),

131    (Figure 1).

132    Within the detected Patescibacteria, site specific and filter size preferences were observed

133    (Figure 2). The shallowest well at the top of the hillslope, H14, showed a relatively higher

134    percentage of Saccharimondales compared to other wells. *Candidatus* Staskwiczbacteria showed

135    preference for wells H14 and H43 (characterized by hypoxic/ anoxic environments with low

136    nitrate), and *Candidatus* Wolfebacteria, UBA9983, and *Candidatus* Liptonbacteria for well H52

137    (characterized by anoxic environment and longest water residence time). *Candidatus*

138 Magasanikbacteria and UBA9983 showed preference for 0.2 μm filter fractions of all the wells,

139 whereas *Candidatus* Woesebacteria was enriched in all the 0.1 μm filter fractions.

**Dominance of Patescibacteria in Hainich groundwater communities enabled recovery of**

**hundreds of high quality MAGs**

142 Metagenomic assembly and binning of all individual groundwater samples (n = 32) yielded a

143 total of 1275 non-redundant manually refined MAGs from various bacterial and archaeal species.

144 Among these MAGs, 584 MAGs were classified as *Cand*. Patescibacteria by GTDB-Tk and 291

145 of them were classified as CPR with high confidence score by a random forest classifier within

146 Anvi'o v6.1 [35,36], trained with a set of CPR specific single copy genes extracted from

147 previously published CPR genomes [15,37] (Additional file 1). Most of these 291 MAGs

148 belonged to the classes: *Cand*. Paceibacteria (163 MAGs) followed by candidate division ABY1

149 (49 MAGs), and *Cand*. Microgenomatia (46 MAGs) (Figure 3, A). The details about all the

150 Patescibacteria MAGs are provided in Additional file 2. The phylogenetic tree constructed from

151 the multiple alignment of 68 core protein sequences confirmed the taxonomic placement of

152 Patescibacteria MAGs (Figure 3, B).

**Differences in the genome sizes of Patescibacteria based on cell size enrichment**

154 We identified 110 Patescibacteria MAGs enriched in the 0.1 μm filter fractions based on their

155 average *rpoB* gene-count-normalized coverage (See Methods) being 5-fold higher than in the 0.2

156 μm filter fractions. Of these, 82 MAGs were further classified as *Cand*. Paceibacteria, and 23 as

157 *Cand*. Microgenomatia. Both classes were absent in the MAGs enriched in 0.2 μm filter

158 fractions. Similarly, 33 Patescibacteria MAGs were enriched 5-fold more in the 0.2 μm filter

159 fractions, with 22 of those belonging to the candidate division ABY1, and 5 to *Cand*.

7

160    Gracilibacteria. Again, none of the genomes classified in these two classes were enriched in the

161    0.1 µm filter fractions.

162    The average genome size of all Patescibacteria MAGs enriched in the 0.1 µm filter fractions

163    (688.7 ± 139.4 kb) was significantly smaller (Dunn's test, p = 1.02e-06) than that of the

164    Patescibacteria MAGs enriched in the 0.2 µm filter fractions (883.1 ± 204.3 kb), (Figure 4, A).

165    There was no significant difference in the genome completeness and contamination values

166    between the two groups.

167    When we analyzed the gene compositions of the two sets of Patescibacteria genomes, the genes

168    encoding type-IV pilus assembly proteins (PilC, PilM, PilO) were significantly overrepresented

169    (two-proportions z-test, p = 1.4e-04) in Patescibacteria enriched in the 0.1 µm filter fractions

170    (~88% of these genomes) as compared to those from the 0.2 µm filter fractions (~64% of these

171    genomes). Similarly, genes encoding cell division proteins FtsW and FtsI were present in 93%

172    and 36% of the Patescibacteria MAGs enriched in 0.1 µm filter fractions, respectively. In

173    comparison, the same genes were present in only 70% and 3% MAGs enriched in the 0.2 µm

174    filter fractions (two-proportions z-test, p = 6.2e-04 and 4.7e-04). The gene encoding for the rod-

175    shape determining protein (MreB) was also more likely to be found in Patescibacteria MAGs

176    enriched in the 0.1 µm filter fraction (95% in the 0.1 µm-enriched vs 75% in the 0.2 µm-

177    enriched, two-proportions z-test, p = 1.8e-03). Additionally, genes involved in colanic acid

178    biosynthesis (*wcaH* and *wcaF*) were uniquely present in ~10% of the Patescibacteria enriched in

179    the 0.1 µm filter fractions.

180    Conversely, the L-lactate dehydrogenase gene was detected in 12% of the MAGs enriched in the

181    0.2 µm filter fractions and was entirely absent in the 0.1 µm-enriched MAGs. A similar pattern

182    was found for the tryptophan synthase genes, *trpA* and *trpB*, which were detected in 15% and

183    18% of the MAGs enriched in the 0.2 μm filter fractions, but absent in Patescibacteria MAGs

184    enriched in the 0.1 μm filter fractions.

**Growth dynamics of Patescibacteria using *in situ* measure of replication**

186    Patescibacteria MAGs had comparatively higher estimated growth measures (GRiD values) in

187    the near surface wells of the groundwater transect (wells H14 and H32), in comparison to the

188    downstream wells (Figure 5, A). Specifically, these Patescibacteria showed significantly higher

189    GRiD values at well H14 as compared to the downstream wells H41 and H43, and significantly

190    higher GRiD values at well H32 as compared to all other wells present downstream. Notably, the

191    wells with highest mean GRiD values for Patescibacteria were also the wells with lowest number

192    of Patescibacteria MAGs. (Additional file 4, Figure S1).

193    The GRiD values were significantly higher (Welch Two Sample t-test, $p = 8.73e-07$) in

194    Patescibacteria MAGs enriched in 0.2 μm filter fractions ($1.40 \pm 0.27$) as compared to

195    Patescibacteria MAGs enriched in the 0.1 μm filter fractions ($1.25 \pm 0.029$). When we compared

196    the GRiD values of individual classes of Patescibacteria between 0.1 and 0.2 μm filter fractions,

197    only MAGs from class Paceibacteria showed significantly higher GRiD values in the 0.2 μm

198    filter fractions (Welch Two Sample t-test, $p = 6.14e-03$, Figure 5, B).

**Limited metabolic and biosynthetic capabilities in Patescibacteria**

200    Metabolic reconstructions based on KEGG modules revealed that the metabolic repertoire of the

201    analyzed Patescibacteria genomes did not show a clear separation by their taxonomy (Figure 6)

202    nor followed a particular pattern in oxic and anoxic wells (Additional file 5, Figure S2). All

9

203 Patescibacteria MAGs lacked central energy metabolism and biosynthetic pathways for most

204 amino acids and vitamins. The tri-carboxylic acid (TCA) cycle was missing in 81.8% of the

205 Patescibacteria MAGs and was incomplete for the remaining 18.2% of the MAGs. Glycolysis

206 was incomplete in all MAGs, pentose phosphate pathway (PPP) was incomplete in 92% of the

207 MAGs, and reductive PPP was absent in 97% of the MAGs. Biosynthesis pathways for most of

208 the amino acids (except serine, glycine and sometimes asparagine) and vitamins (except

209 cobalamin and thiamin) were missing in most of the Patescibacteria MAGs. In addition, electron

210 transport chain complexes (I-IV) were not identified, with exception of gene encoding for the F-

211 Type ATPase (from ETC complex V) in 59.7% of the Patescibacteria.

212 However, Patescibacteria possessed some notable genes, namely those coding for copper

213 transporter (*copA*) and cobalt transporter (*corA*) that are usually found in pathogenic bacteria

214 [38,39]. Also, carbohydrate active enzymes (CAZy) responsible for degradation of starch (11%

215 MAGs), polyphenolics (25% MAGs) and chitin (11% MAGs) were observed. At least 13% of

216 the MAGs had more than one type of CAZy. Patescibacteria also encoded genes for small chain

217 fatty acids (SCFA) and alcohol conversion functions e.g. D-lactate dehydrogenase (25% MAGs),

218 L-lactate dehydrogenase (4% MAGs), and conversion of pyruvate to Acetyl-CoA (K00174, 14%

219 MAGs). Acetate kinase was found in only 6% of the Patescibacteria MAGs. A mutually

220 exclusive presence of D- and L-lactate dehydrogenases was observed.

221 **Genomic signs of adaptive response of Patescibacteria to oxic and anoxic conditions**

222 We classified 134 Patescibacteria MAGs as 5-fold enriched in oxic wells (H32, H41 and H51)

223 and 64 Patescibacteria MAGs as 5-fold enriched in anoxic wells (H14, H43 and H52). No

224 taxonomic preference for oxic or anoxic conditions was observed. Patescibacteria MAGs

225      enriched in oxic sites showed some unique features with respect to their ability to resist oxidative

226      stress. We found that superoxide dismutase genes (at least one of the *sodA*, *sodB*, *sodC*, *sodF*,

227      *sodM*, *sodN* or *chrC* genes) were encoded by significantly higher proportion (82.8%) of the

228      Patescibacteria MAGs enriched in oxic wells than in anoxic wells (65.6%) (two-proportions z-

229      test, p = 8.8e-03), but there was no evidence for other stress regulator genes (*oxyR*, *soxR*, *soxS*,

230      *rpoS*). There were no relevant metabolic pathways or genes specific to the 64 Patescibacteria

231      MAGs enriched in anoxic wells (Additional file 5, Figure S2).

232      Correlation of the genomic coverages (relative abundances) of the Patescibacteria MAGs

233      enriched in oxic wells with the dissolved oxygen concentrations revealed highly significant

234      positive correlations for 28 MAGs (Additional file 6). Most of these MAGs belonged to class

235      *Cand*. Paceibacteria (family UBA1539/*Yonathbacteraceae*) and genus GWC2-37-13 from order

236      UBA1406/*Roizmanbacterales*. Most of these MAGs (82%) carried superoxide dismutase gene

237      (K04564) essential for protection against free superoxide radicals in oxic environments.

238      We chose the most abundant, high quality Patescibacteria MAGs from oxic well H41 (H41-

239      bin288, 0.1 μm filter fraction, relative abundance = 0.75% ± 0.15) and anoxic well H52 (H52-

240      bin095, 0.1 μm filter fraction, relative abundance = 2.28% ± 0.37) as model organisms to

241      illustrate the commonalities and divergences in their genomes (Figure 7). We also included the

242      second most abundant Patescibacteria MAG from the same oxic well H41 (H41-bin049, 0.1 μm

243      filter fraction, relative abundance = 0.41% ± 0.02) from the same taxonomic family as the anoxic

244      representative. This was done to rule out the genomic differences due to the relatively distant

245      evolutionary history of the first pair (H41-bin288 and H52-bin095). The representative MAGs

246      H41-bin288 and H41-bin049 from the oxic well H41 showed positive correlations with oxygen

247   (R = 0.88, p = 2.0e-02 and R = 0.75, n.s., respectively), while the representative MAG from

248   anoxic well (H52-bin095) showed a negative correlation (R = −0.43, n.s.).

249   Features specific to both representative genomes from oxic well H41 were genes coding for F-

250   type H$^+$-transporting ATPase (subunit a, b, c, α, β and γ), NitT/TauT family transporter (involved

251   in transport of inorganic ions like nitrate, sulfonate, and bicarbonate), and nitrite reductase (*nirK*

252   involved in conversion of nitrite to nitric oxide). On the other hand, genes related to sugar

253   sensing and multiple sugar transport systems (ABC.MS.S), and lactate dehydrogenase

254   (fermentation) were specific to the anoxic representative. Common genes or functions were

255   found for all three representative genomes, e.g. genes encoding type IV pilus assembly proteins

256   (PilB, PilC, PilM, and PilO) as well as competence proteins (ComEC, ComFC), useful for DNA

257   uptake from exogenous sources, superoxide dismutase (SOD2) for protection against superoxide

258   radicals, transporters of metal ions like zinc, copper, calcium, nickel. We also identified genes

259   encoding for rod-shape determining proteins, like RodA with additionally related genes encoding

260   for proteins like MreB and MreC in the anoxic representative.

261   **Co-occurrence patterns of Patescibacteria with other microbial species**

262   A co-occurrence network generated using metagenomic abundances of MAGs revealed that

263   many species of Patescibacteria were consistently co-occurring with one another, as well as with

264   species of other bacteria and archaea (Figure 8). The average normalized genome coverages for

265   all the studied MAGs across both filter fractions of all the wells are provided in Additional file 7.

266   The most common one-to-one associations were observed with MAGs from the phyla

267   Nanoarchaeota (mostly order Pacearchaeales), Bacteroidota, MBNT15, and Bdellovibrionota. A

268   small isolated cluster within the network showed indirect but close associations of

12

269    Patescibacteria with multiple members of the phylum Nitrospirota (genus RGB.16.64.22), and

270    phylum Omnitrophota (Figure 8).


271    Under the assumption that Patescibacteria were physically associated with larger host cells, we

272    simplified our co-occurrence network to further refine the associations in the 0.2 µm filter

273    fractions (using the 5-fold coverage cut-off as compared to 0.1 µm filter fractions). This follow-

274    up co-occurrence network showed one-to-one associations of MAGs of the phylum

275    Omnitrophota (class koll11) with MAGs from Patescibacteria (each one from the classes

276    Paceibacteria, Microgenomatia, and candidate division ABY1). One of the MAGs from class

277    Paceibacteria showed association with a Proteobacteria MAG (order Rickettsiales), while a

278    MAG from candidate division ABY1 showed direct connections with two Bacteroidota MAGs.

279    Another MAG from class Gracilibacteria showed direct connections with 5 Nitrospirota MAGs

280    from the same genus UBA1546 (Figure 8). The sequence coverages of these highlighted genome

281    pairs or clusters across the metagenomes are compared in Additional file 8, Figure S3 and

282    Additional file 9, Figure S4. Two Actinobacteria MAGs belonging to the species

283    *Aurantimicrobium* sp003194085 also showed associations with Patescibacteria. The first

284    *Aurantimicrobium* interacted with a Patescibacteria (*Cand*. Paceibacteria) MAG, and the second

285    with multiple Patescibacteria (2 *Cand*. Paceibacteria, 2 *Cand*. Gracilibacteria and 3 candidate

286    division ABY1) MAGs.


287    When we searched for sequence similarity of all gene open reading frames (ORFs) from all

288    Patescibacteria MAGs to ORFs from all other bacterial and archaeal MAGs in the present study

289    using blastn [40], we found various ORFs from other taxa highly similar to Patescibacteria ORFs

290    (95% sequence identity covering 85% length of the query and hit sequences). The most ORFs

291    that matched were between members of genus UBA10092 of Patescibacteria (class

292    Paceibacteria) and two members of the family UBA12090 of Omnitrophota (34 and 16 ORFs,

293    respectively). They included genes encoding for twitching motility protein PilT (K02669), P-

294    type Cu+ transporter (K17686) and lipopolysaccharide export system permease protein

295    (K11720). Between members of genus UBA11707 of Patescibacteria (class ABY1) and genus

296    UBA1573 of Proteobacteria (family Micavibrionaceae), 14 such ORFs, including gene encoding

297    for ABC-2 type transport system ATP-binding protein (K01990), were observed. Thirteen such

298    ORFs, including gene for ABC-2 type transport system permease protein (K01992), were

299    observed between members of the family Zambryskibacteraceae of class Paceibacteria and genus

300    ASMP01 of Nanoarchaeota.

301    To have an idea about the temporal co-occurrence patterns of other groundwater microbes with

302    Patescibacteria, we additionally utilized time-series data based on 16S rRNA gene amplicon

303    sequencing from the same groundwater transect from three wells (H41, H43 and H52) measured

304    over more than six years [41]. We observed that Patescibacteria co-occurred mostly with

305    members of phyla Proteobacteria (mostly order Burkholderiales) and Nitrospirota (order

306    Thermodesulfovibrionia), in the well H41; Verrucomicrobiota, in the well H43 and

307    Planctomycetota (mostly genus *Brocadia*) in the well H52. Similarly, a Patescibacteria MAG

308    was identified to co-occur with multiple Thermodesulfovibrionia MAGs belonging to the

309    phylum Nitrospirota in this study.

## Discussion

311    Our comprehensive metagenomic analyses revealed that modern pristine groundwater of the

312    Hainich CZE is clearly dominated by *Cand*. Patescibacteria with an average relative abundance

313    of 50% across all wells and a maximum of 79% in the 0.1 μm filter fraction. Compared to other

14

314   groundwater communities dominated by CPR bacteria ranging from 2-28% [16], 3-40% [18], 10-

315   28% [7] and 36-65% [15], the exceptionally high abundance of CPR bacteria discovered in this

316   study is distributed over distinct geochemical zones spanning oxic and anoxic conditions [17,33].

317   Although the spatial distribution patterns of the different *Cand*. Patescibacteria taxa (Figure 2)

318   were less pronounced than those observed in other bacteria in groundwater of the Hainich CZE

319   [33,41], and despite their streamlined genomes, we could highlight certain environmental

320   preferences of the *Cand*. Patescibacteria. Access to 587 manually curated MAGs of *Cand*.

321   Patescibacteria, assigned to different filter fractions, allowed us to shed some light on genomic

322   characteristics linked to their cell size and a putative free living or host attached lifestyle.

323   Patescibacteria have been described mostly in anoxic or hypoxic environments [42,43]. Our data

324   show no major metabolic or taxonomic differences in Patescibacteria enriched in oxic and anoxic

325   groundwater wells. Significantly higher proportion of superoxide dismutase genes in

326   Patescibacteria MAGs enriched in oxic groundwater wells compared to those in anoxic wells is

327   an example of spatial differentiation that might be due to an environmental selection mechanism,

328   as these enriched species have an advantage to withstand the presence of oxygen radicals when

329   exposed to high $O_2$ concentrations. More than 80% of the Patescibacteria MAGs enriched in oxic

330   wells could potentially resist superoxide radicals, and more than 20% showed a positive

331   correlation to oxygen concentrations, in particular those belonging to class *Cand*. Paceibacteria

332   (family UBA1539/*Yonathbacteraceae*) and to order UBA1406/*Roizmanbacterales*. But even

333   closely related Patescibacteria species showed different preferences for oxygen concentrations in

334   terms of metabolic pathways (Figure 7).

335   The permanently high $O_2$ concentration in well H32 (2.23 ± 0.56 mg/L) and especially in well

336   H41 (4.83 ± 1.7 mg/L) [41,44], did not lead to enrichment of groundwater Patescibacteria MAGs

15

337    with genetic traits of energy harvesting mechanisms through aerobic respiration. Exposure to

338    oxygen is not exceptional for *Cand.* Patescibacteria, as oxic soils are the main source for their

339    vertical translocation into shallow groundwater [17,24]. *Cand.* Patescibacteria represent only

340    0.55% of the total bacterial soil community in the preferential forest surface-recharge area of the

341    Hainich CZE (Herrmann *et al.* 2021, unpublished observations). Despite this low abundance,

342    these ultra-small organisms are readily mobilized from soil, especially during winter months

343    when ionic strength of the seepage is very low (Herrmann *et al*. 2021, unpublished observations),

344    and as such constitute the largest fraction of taxa shared between seepage and shallow

345    groundwater [17].

346    The most abundant Patescibacteria MAG from oxic well H41 (H41-bin288) had genes that

347    encode for nitrite transport and its subsequent reduction into nitric oxide involving

348    ferricytochrome c. Also, this genome possessed a gene for F-Type ATPase to generate energy by

349    ATP formation and it did not encode genes for fermentation (L- or D-lactate dehydrogenase).

350    This collectively suggests the possibility of an alternative anaerobic respiration mechanism in

351    this particular genome. Despite the low *in situ* concentrations of nitrite, it might be alternatively

352    provided by the nitrification process. This relates to the fact that well H41 is characterized as a

353    nitrification hotspot with measured rates of $0.48 \pm 0.09$ and $0.64 \pm 0.39$ nmol $NO_x$ liter$^{-1}$ h$^{-1}$ [45]

354    and to the high relative abundances of *Nitrospira* on the metagenome level and *Thaumarchaeota*

355    on the metatranscriptome level [46]. Presence of genes coding for multiple subunits of F-Type

356    ($H^+$ transporting) ATPase in this genome confirms the existence of supplementary ATP synthesis

357    machinery, which are commonly observed in aerobic bacteria [47]. Similarly, notable features

358    specific to both representative genomes from oxic well H41 included genes involved in the

359    transport of inorganic ions like nitrate, sulfonate, and bicarbonate.

360   The almost complete absence of the aerobic respiration machinery i.e. the electron transport

361   chain complexes, terminal oxidases / electron acceptors, and gene products associated with the

362   TCA cycle, along with widespread presence of L- or D-lactate dehydrogenases confirms the

363   previously postulated fermentative lifestyles of Patescibacteria [11,15,48] in members of the

364   three lineages OD1 (Parcubacteria), OP11 (Microgenomates), and BD1-5 (Gracilibacteria).

365   Parcubacteria were proposed to produce acetate, ethanol, lactate, and hydrogen as fermentation

366   products based on metagenomic and proteomic analysis [3,15,48]. Presence of L- or D-lactate

367   dehydrogenase genes in one third of the Patescibacteria MAGs indicates specificity for

368   fermentation substrates. In one tenth of the MAGs enriched in 0.1 µm filter fractions, specificity

369   for L-lactate could be observed based on the exclusive presence of L-lactate dehydrogenase

370   genes. Presence of multiple carbohydrate active enzymes (CAZy) in many Patescibacteria

371   suggests their potential for degradation of multiple complex compounds like starch, chitin, and

372   polyphenolics.

373   The spatial differentiation of *Cand.* Patescibacteria could also be indirectly caused by the

374   preference of a putative host organism for certain environmental conditions. The oxic, nitrate-

375   rich (15.71 mg/L) groundwater of well H41 was dominated by Nitrospirota MAGs, and 5 of

376   them co-occurred with a single Patescibacteria MAG (H52-bin081_1, *Cand.* Gracilibacteria) and

377   had similar abundance patterns (Additional file 9, Figure S4). As some Nitrospirota MAGs (n =

378   51) were enriched exclusively in oxic wells, their preference might have determined the

379   distribution pattern of putative CPR episymbionts. Nitrospirota species were also found to be

380   consistently co-occurring with Patescibacteria in some of the studied wells based on OTU

381   abundances from 16S rRNA gene amplicon sequencing data collected over 6.5 years [41] as well

17

382    as MAG abundances from this study across the groundwater transect. At the minimum, these

383    observations suggest common niche preferences between some members of these two phyla.

384    To elucidate other possible associations of Patescibacteria with other prokaryotes, we utilized

385    above mentioned time-series data that revealed consistent co-occurrence of Patescibacteria OTUs

386    with OTUs from Proteobacteria, Verrucomicrobiota, and Planctomycetota in addition to OTUs

387    from Nitrospirota [41]. When we looked into the genomic characteristics of all Patescibacteria

388    and all other MAGs, we found various ORFs from other taxa highly similar with Patescibacteria,

389    between members of (i) class Paceibacteria and family Omnitrophota, (ii) class ABY1 and

390    family Micavibrionaceae, and (iii) family Zambryskibacteraceae of class Paceibacteria and genus

391    ASMP01 of Nanoarchaeota, suggesting probable acquisition of motility and transport functions

392    from other bacteria or archaea.

393    Network analysis based on abundances of all MAGs of both filter fractions revealed that the

394    members of the phyla Bacteroidota, MBNT15, and Bdellovibrionota along with members of

395    phyla Nitrospirota and Omnitrophota had direct specific connections with some Patescibacteria.

396    Furthermore, we restricted the network analysis only to MAGs enriched on the 0.2 μm filter

397    fractions (57 Patescibacteria and 423 other MAGs) in order to identify Patescibacteria that would

398    be potentially attached to other larger host cells. This narrowed-down analysis showed

399    interactions of Patescibacteria with few specific MAGs of the phyla Bacteroidota, Nitrospirota,

400    Omnitrophota, and Actinobacteria. Our co-occurrence analysis did not reveal direct connections

401    of Actinobacteria MAGs with any of the Saccharibacteria, although Actinobacteria are reported

402    as host for Saccharibacteria (TM7) in human oral cavity [12,27,29]. However, direct network

403    connections of *Aurantimicrobium* species, members of the phylum Actinobacteria with multiple

404    other Patescibacteria MAGs from classes Paceibacteria, Gracilibacteria, and candidate division

405    ABY1 hint towards possible host-symbiont relationships in these particular pairs.

406    Direct one-to-one connections with members of other phyla were found in only 5 out of 57

407    (8.77%) Patescibacteria MAGs enriched in 0.2 µm filter fractions, suggesting that the majority of

408    groundwater Patescibacteria of the Hainich CZE is not specifically associated with one single

409    host, but associations with multiple hosts cannot be ruled out. The attachments between cells are

410    often fragile and may be partly or completely disrupted during filtration and sample processing

411    steps, and hence are difficult to track using sequential filtration. An even lower percentage of

412    associations (<1.5%) based on potentially co-sorted SAGs containing DNA from heterogeneous

413    sources was reported from Beam *et al*. 2020 [11].

414    On average, Patescibacteria enriched in 0.1 µm filter fractions had 22% smaller genome size

415    than those enriched in 0.2 µm filter fractions, and it has been previously shown that smaller cell

416    size is linked to genome reduction [49,50]. This genome size difference might be due to

417    differences in average cell sizes of *Cand*. Paceibacteria and *Cand*. Microgenomatia that were

418    preferentially enriched within 0.1 µm filter fractions; and candidate division ABY1, and *Cand*.

419    Gracilibacteria that were preferentially enriched within the 0.2 µm filter fractions. Smaller

420    genomes in tiny CPRs might be the result of genome streamlining leading to lack of complex

421    energy metabolism and biosynthetic capabilities which makes them rely on other cells through

422    cell-cell attachment.

423    We found Type IV pilus assembly proteins in a higher proportion of Patescibacteria enriched in

424    0.1 µm filter fractions. These proteins are responsible for formation of pilin-like appendages that

425    are involved in a variety of functions like adherence to host cells, locomotion, DNA uptake as

426    well as protein secretion in bacteria [51], which would support physical association with other

427    microbes. Type IV pili (T4P) are essential for virulence of some Gram-negative pathogenic

428    bacteria [52] and also found in Gram-positive bacteria with a different pilus assembly

429    mechanism involving a sortase [53]. Pili like appendages were microscopically shown to form

430    surface attachment of CPR bacteria with other (host) large cells [18]. The symbiotic association

431    of TM7i (*Cand*. Saccharibacteria) with its host *Leucobacter aridocollis* J1, mediated by T4P was

432    identified in a co-culture experiment [54]. As pilus mediated attachments are often fragile, small

433    Patescibacteria cells passing through the 0.2 µm filters do not necessarily indicate lack of cell-

434    cell attachment with larger bacterial cells. Many of these ultra-small Patescibacteria appear to

435    have a rod-shaped morphology, as genes encoding the rod shape-determining protein (MreB)

436    were found in a higher proportion of MAGs enriched in 0.1 µm filter fractions. The recent

437    reconstruction of the last bacterial common ancestor (LBCA) genome of CPR lineage suggests a

438    rod-shaped morphology [55]. However, most of the reported morphologies for the

439    Patescibacteria are cocci [12,18,21]. Although we cannot rule out that some of the larger rod-

440    shaped Patescibacteria could still pass through the 0.2 µm filter pores, this would not explain the

441    enrichment in the 0.1 µm filter fractions. More direct microscopic visualization is needed to

442    verify the morphology of these ultra-small Patescibacteria.

443    We found higher growth rates of Patescibacteria in near-surface wells (H14, H32) of the

444    groundwater transect than in the ones more downstream. Growth of CPR bacteria is stimulated

445    after attachment to host-cells [18]. As cell-cell aggregations might be more prone to dispersal

446    limitations in a dense rock matrix, surface-near wells could have higher probabilities of host

447    interactions. But our co-occurrence analysis did not reveal direct connections of CPR MAGs

448    with higher growth rates with other MAGs.

449     Groundwater of the very shallow well H14, located uphill of the transect, shows a fast response

450     to weather events [56], and is characterized by both the highest bacterial diversity and the

451     presence of well-known surface heterotrophs; whereas core groundwater species dominated

452     groundwater microbiomes in the downstream direction [33]. This well, along with the other near-

453     surface well (H32) showed the lowest relative abundances of Patescibacteria and of

454     Patescibacteria MAGs, although those that were detected had higher expected replication rates

455     on average. A possible explanation for this pattern is that surface exported members were

456     replicating within the soil before being flushed into the groundwater. Other, more successful

457     groundwater CPR groups may have slower growth and replication rates within the transect due to

458     much lower microbial cell densities and less available organic carbon. Indeed, some taxa such as

459     those belonging to *Cand.* Saccharimonadia, which had among the highest growth rates, did not

460     flourish within other wells of the groundwater transect. We hypothesize that they might be more

461     adapted to soil habitats, which was also observed in previous studies [17].

462     The predominance of particular CPR species in oxic (H41) and anoxic (H52) wells appears to be

463     the result of environmental preference or exploitation of other organisms for cellular

464     requirements in the nutrient deficient groundwater. Some potential hosts supporting an

465     episymbiotic lifestyle could be identified. The environmental preference of some of these hosts,

466     e.g. Nitrospirota for oxygen and nitrogen in well H41, would explain the predominance of their

467     potential Patescibacteria episymbiont in H41, with an estimated episymbiont-to-host ratio of

468     3.6:1 based on coverages of Patescibacteria and Nitrospirota MAGs in total coverage of all

469     binned genomes. But the vast majority of the ultra-small Patescibacteria in the groundwater

470     appears to be free-living, self-sufficient with their minimal genomes [11,42], adapted to

471     oligotrophic conditions with low growth rates, and equipped with genes to cope with oxidative

472    stress only if needed. We found evidence that the majority has the capability to attach to other

473    cells, which appears to also include other Patescibacteria, and this attachment might be not very

474    specific or for longer time periods, just long enough to loot or exchange supplies.

## Conclusions

476    The Candidate Phyla Radiation represent the largest phylogenetic diversity within the bacterial

477    domain, which has not been reflected in the metabolic versatility of genomic representatives

478    studied to date. Here we leveraged a well characterized aquifer transect, that is dominated by

479    members of the CPR and spans large biogeochemical gradients, to explicitly explore genomic

480    adaptations to environmental conditions. The most significant and surprising result was the high

481    level of similarity in predicted metabolic functions and expected lifestyles that spanned large

482    redox gradients from fully oxic to completely anoxic groundwater, both within the larger CPR

483    clade as well as at finer phylogenetic resolutions. One noteworthy exception was a differential

484    abundance in superoxide dismutase, a potentially useful indicator of oxygen exposure in CPR

485    genomes recovered from other environments or already deposited to sequence databases. Due to

486    a suspected dependence on other bacterial hosts, we searched among >1200 constructed MAGs

487    and a larger amplicon dataset for potential partners, finding that only 8% of CPR MAGs

488    exhibited significant one-to-one relationships. Therefore, we propose that most members of the

489    CPR form non-specific associations, attaching to multiple hosts to supplement their energetic

490    demands within oligotrophic groundwaters.

## Methods and Materials

492    **Groundwater sampling, DNA extraction and sequencing**

22

493    Samples were collected from a groundwater transect system spanning through a ~6 km long zone

494    including forest, pasture and agricultural land within the Hainich Critical Zone Exploratory

495    (CZE) located in Thuringia, Germany. The Hainich CZE was established and extensively studied

496    by Collaborative Research Center AquaDiva [32]. The groundwater was collected from 6 wells

497    (H14, H41, H43, H51, H52) in January 2019 and (H32) in November 2018 spanning various

498    zones of the transect. For each well, on average $61.3 \pm 35.4$ liters of groundwater was filtered

499    through 0.2 µm filters (Omnipore Hydrophilic PTFE membrane, Merck Chemicals GmbH)

500    followed by 0.1 µm filters in triplicates (except for well H32 where there were only two

501    replicates out of which one from the November 2018 sampling campaign was used as biological

502    replicate). All the 32 filter fractions were immediately frozen and stored under -80°C. The DNA

503    was extracted using a phenol/chloroform protocol, the libraries generated with an NEBNext

504    Ultra FS DNA preparation kit, and sequenced on an Illumina NextSeq 500 system with paired-

505    end library ($2 \times 150$ bp).

506    On an average $9.8 \pm 1.15$ Gb of raw DNA sequence data were obtained from each of the 32 filter

507    fractions. Of which, $86.12 \pm 0.57$ % of the reads were of very high quality (at least quality score

508    Q40). Subsequent quality control steps like adapter trimming, PhiX detection and removal using

509    BBDuk (bbtools version 37.09, written by Brian Bushnell, last modified March 30, 2017) further

510    improved the quality of the reads. These high-quality reads were then used for metagenomic

511    assembly and followed by genome binning steps.

512    **Metagenomic assembly, genome binning and refinement**

513    The quality controlled reads of each individual filter fraction replicate were assembled and

514    scaffolded using metaSPAdes v3.13 [57]. Scaffolds larger than 1 kb were used for downstream

515    analyses. Genome binning was carried out using three binning algorithms - Abawaca v1.07 [15],

516    ESOM [58,59] and Maxbin2 v2.2.4 [60]. The values 3000 and 5000 bp as well as 5000 and

517    10000 bp were used as *-min* and *-max* parameters to calculate 4-mer frequencies for Abawaca

518    and ESOM (the script esomWrapper.pl, https://github.com/tetramerFreqs/Binning), and both the

519    40 and 107 marker gene sets were utilized in Maxbin2. DASTool v1.1 [14] was used to

520    determine the best bins among these approaches. Bins were further refined manually inside the

521    Anvi'o workflow v6.1 [35,36]. The quality of the refined bins (completeness and

522    contamination/redundancy) was also calculated based on domain-level single-copy core genes

523    within Anvi'o. Genomes from each assembly were de-replicated using dRep v2.6.2 [61] at 99%

524    ANI to remove strain level redundancy across sites, resulting into 1275 representative MAGs.

525    Genome coverages were calculated within Anvi'o, and were normalized using number of RNA

526    polymerase B (*rpoB*) genes identified within the metagenomic reads.

527    **Taxonomic assignments, gene annotations and pathway predictions**

528    Overall community composition of each metagenome was determined using phyloFlash v3.4

529    [62] based on proportions of reads mapped to SILVA SSU rRNA Ref NR99 database, Release

530    138 [34]. Taxonomic classification of individual MAGs was performed by GTDB-Tk v0.3.2 [63]

531    using GTDB Release 89 as reference database. Out of the 1275 genomes GTDB-Tk classified

532    587 genomes as *Cand*. Patescibacteria at phylum level. We used *anvi-script-gen-CPR-classifier*

533    script from Anvi'o v6.1 [35,36] which uses supervised machine learning model (random forest

534    classifier) to train the program and *anvi-script-predict-CPR-genomes* for predicting the

535    probability of the MAGs to confirm the CPR genomes. The training is based on the profile of

536    previously published 139 single copy core genes from hundreds of CPR genomes from Brown *et*

537    *al.* [15] and Campbell *et al.* [37] as input. This model confirmed 291 out of 587 genomes as CPR

538    with a high confidence score (75% or more). While the model was inconclusive in case of the

24

539    remaining 174 genomes based on low confidence score (less than 75%) and the remaining 122

540    genomes were discarded due to their completion levels below 50%.

541    The gene annotations, coding sequences, respective protein sequences, coverage calculations and

542    other mapping statistics for all the genomes were exported by *anvi-summarize* program from

543    within the Anvi'o workflow. The annotations were also carried out using Prodigal v2.6.3 [64].

544    Distilled and Refined Annotation of Metabolism (DRAM) [65] was used to generate pathway /

545    metabolism summaries. At least one proper (other than hypothetical, uncharacterized or gene

546    with unknown function) annotation from KEGG [66], MEROPs [67], Pfam [68] or dbCAN [69]

547    was considered. This generated a single tab delimited annotation file listing the best hits from all

548    these databases as well as summaries focused on most important pathways and functions. The

549    pathway coverages (completeness) of central metabolism pathways were calculated based on

550    KEGG modules definitions (https://www.genome.jp/kegg/module.html).

551    **Phylogenetic analysis**

552    Single copy core bacterial genes were detected in all the 1087 bacterial MAGs using hmm

553    profile (default 'Bacteria_71' hmm profile in Anvi'o v6.1), their protein sequences were

554    extracted and aligned using MUSCLE [70] from within the Anvi'o [35,36]. A phylogenetic tree

555    based on multiple sequence alignment of the 68 core proteins present in all bacterial MAGs

556    (1087) was constructed using Approximate Maximum Likelihood in FastTree v2.1.11 SSE3,

557    OpenMP [71] with 1000 bootstrap replications. The subset of the tree was used for arranging the

558    metabolic pathways of 291 selected Patescibacteria MAGs in Figure 6.

559    **In-situ measurement of replication**

560    The forward sequencing reads from all the metagenomes were mapped to the MAGs to calculate

561    the sequence coverage of individual contigs. These coverage profiles were utilized to calculate

562    Growth Rate InDex (GRiD) [72] which is directly proportional to the growth rates of the cells in

563    a given environment. GRiD measures the difference in genome copies closer to the origin of

564    replication compared to the terminus caused by ongoing replication forks. The coverage cut-off

565    of 0.7 was used to remove extremely low coverage contigs.

566    **Statistical analyses**

567    The difference in the mean genome sizes of the MAGs enriched in different filter fractions were

568    compared using Kruskal-Wallis rank sum test followed by pairwise Dunn's test in R [73]. The

569    proportions of gene annotations (KEGG) in the MAGs enriched in different filter fractions or

570    oxic and anoxic wells were compared with two-proportions z-test with Yates' continuity

571    correction in R. The p-values were adjusted for multiple testing using 'fdr' correction unless

572    otherwise mentioned.

573    **Co-occurrence network analysis**

574    We used normalized average genome coverages of all the 1275 MAGs across all the

575    metagenomes as the approximation of abundance profiles of species from respective

576    metagenomes. This abundance matrix was used to calculate proportionality of the coverage

577    profiles in R package propR v4.2.6 [74]. A $\rho$ cutoff of 0.95 was used for network creation to

578    highlight only the most relevant co-occurrences. The network was generated using the R package

579    igraph v1.2.6 [75] and exported to Cytoscape v3.8.2 [76] for visualization using R package Rcy3

580    v2.8.1 [77].

581    **Search for ORF similarity**

582    We carried out blastn [40] search on all the annotated ORFs for Patescibacteria MAGs as a query

583    against all the ORFs of all the MAGs other than Patescibacteria. We filtered the results based on

584    95% sequence identity over 95% query and hit ORF length with e-value cut off of 1.0e-5. We

585    chose only one hit in case of more than one hits for the same query sequence.

## Declarations

### Availability of data and material

588    Data used for this study were deposited into the European Nucleotide Archive (ENA). The raw

589    metagenomic sequencing reads were deposited under ENA project accession PRJEB36505,

590    assemblies for individual samples were deposited under ENA project accession PRJEB36523.

### Competing interests

592    The authors declare that they have no competing interests.

### Funding

27

603 EVE, a joint effort of both the Helmholtz Centre for Environmental Research - UFZ

604 (http://www.ufz.de/) and the German Centre for Integrative Biodiversity Research (iDiv) Halle-

605 Jena-Leipzig (http://www.idiv-biodiversity.de/).

606 **Authors' contributions**

607 NMC, KK, WAO, MT, AJP designed this study. WAO, KK, AJP, TLVB, and MT planned,

608 designed, and conducted the metagenomic sampling approach. MM, MH helped during

609 metagenomic sequencing. NMC, WAO, TLVB, and AJP performed the metagenomic analysis.

610 NMC manually curated and performed comparative genome analysis of the MAGs. PAFG

611 conducted the metabolic reconstruction analysis of representative MAGs. NMC, KK, WAO

612 wrote the manuscript with the help of all authors.

613 **Acknowledgements**

# References

622 1. Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, et al. A new view of the tree
623 of life. Nature Microbiology. 2016;1:16048.

624 2. Elshahed MS, Najar FZ, Aycock M, Qu C, Roe BA, Krumholz LR. Metagenomic analysis of the
625 microbial community at Zodletone Spring (Oklahoma): insights into the genome of a member of the

28

626    novel candidate division OD1. Applied and environmental microbiology. Am Soc Microbiol;
627    2005;71:7598–602.

628    3. Wrighton KC, Thomas BC, Sharon I, Miller CS, Castelle CJ, VerBerkmoes NC, et al. Fermentation,
629    hydrogen, and sulfur metabolism in multiple uncultivated bacterial phyla. Science. American Association
630    for the Advancement of Science; 2012;337:1661–5.

631    4. Mohiuddin MM, Salama Y, Schellhorn HE, Golding GB. Shotgun metagenomic sequencing reveals
632    freshwater beach sands as reservoir of bacterial pathogens. Water Research. 2017;115:360–9.

633    5. Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, et al. Insights into the
634    phylogeny and coding potential of microbial dark matter. Nature. 2013;499:431–7.

635    6. Probst AJ, Castelle CJ, Singh A, Brown CT, Anantharaman K, Sharon I, et al. Genomic resolution of a
636    cold subsurface aquifer community provides metabolic insights for novel microbes adapted to high $CO(2)$
637    concentrations. Environ Microbiol. England; 2017;19:459–74.

638    7. Probst AJ, Ladd B, Jarett JK, Geller-McGrath DE, Sieber CM, Emerson JB, et al. Differential depth
639    distribution of microbial function and putative symbionts through sediment-hosted aquifers in the deep
640    terrestrial subsurface. Nature microbiology. Nature Publishing Group; 2018;3:328–36.

641    8. Correa-Galeote D, Bedmar EJ, Fernández-González AJ, Fernández-López M, Arone GJ. Bacterial
642    Communities in the Rhizosphere of Amilaceous Maize (Zea mays L.) as Assessed by Pyrosequencing.
643    Frontiers in Plant Science. 2016;7:1016.

644    9. Frey B, Rime T, Phillips M, Stierli B, Hajdas I, Widmer F, et al. Microbial diversity in European alpine
645    permafrost and active layers. FEMS Microbiology Ecology [Internet]. 2016;92. Available from:
646    https://doi.org/10.1093/femsec/fiw018

647    10. Wurzbacher C, Nilsson RH, Rautio M, Peura S. Poorly known microbial taxa dominate the
648    microbiome of permafrost thaw ponds. The ISME journal. Nature Publishing Group; 2017;11:1938–41.

649    11. Beam JP, Becraft ED, Brown JM, Schulz F, Jarett JK, Bezuidt O, et al. Ancestral absence of electron
650    transport chains in Patescibacteria and DPANN. Frontiers in microbiology. Frontiers; 2020;11:1848.

651    12. He X, McLean JS, Edlund A, Yooseph S, Hall AP, Liu S-Y, et al. Cultivation of a human-associated
652    TM7 phylotype reveals a reduced genome and epibiotic parasitic lifestyle. Proc Natl Acad Sci U S A.
653    2015;112:244–9.

654    13. Baker JL, Morton JT, Dinis M, Alvarez R, Tran NC, Knight R, et al. Deep metagenomics examines
655    the oral microbiome during dental caries, revealing novel taxa and co-occurrences with host molecules.
656    Genome research. Cold Spring Harbor Lab; 2021;31:64–74.

657    14. Shaiber A, Willis AD, Delmont TO, Roux S, Chen L-X, Schmid AC, et al. Functional and genetic
658    markers of niche partitioning among enigmatic members of the human oral microbiome. Genome
659    biology. BioMed Central; 2020;21:1–35.

660    15. Brown CT, Hug LA, Thomas BC, Sharon I, Castelle CJ, Singh A, et al. Unusual biology across a
661    group comprising more than 15% of domain Bacteria. Nature. Nature Publishing Group; 2015;523:208–
662    11.

663    16. Danczak R, Johnston M, Kenah C, Slattery M, Wrighton KC, Wilkins M. Members of the Candidate
664    Phyla Radiation are functionally differentiated by carbon-and nitrogen-cycling capabilities. Microbiome.
665    Springer; 2017;5:1–14.

666    17. Herrmann M, Wegner C-E, Taubert M, Geesink P, Lehmann K, Yan L, et al. Predominance of Cand.
667    Patescibacteria in groundwater is caused by their preferential mobilization from soils and flourishing
668    under oligotrophic conditions. Frontiers in microbiology. Frontiers; 2019;10:1407.

669    18. He C, Keren R, Whittaker ML, Farag IF, Doudna JA, Cate JH, et al. Genome-resolved metagenomics
670    reveals site-specific diversity of episymbiotic CPR bacteria and DPANN archaea in groundwater
671    ecosystems. Nature microbiology. Nature Publishing Group; 2021;6:354–65.

672    19. Vigneron A, Cruaud P, Langlois V, Lovejoy C, Culley AI, Vincent WF. Ultra-small and abundant:
673    Candidate phyla radiation bacteria are potential catalysts of carbon transformation in a thermokarst lake
674    ecosystem. Limnology and Oceanography Letters. Wiley Online Library; 2020;5:212–20.

675    20. Vavourakis CD, Andrei A-S, Mehrshad M, Ghai R, Sorokin DY, Muyzer G. A metagenomics
676    roadmap to the uncultured genome diversity in hypersaline soda lake sediments. Microbiome. Springer;
677    2018;6:1–18.

678    21. Luef B, Frischkorn KR, Wrighton KC, Holman H-YN, Birarda G, Thomas BC, et al. Diverse
679    uncultivated ultra-small bacterial cells in groundwater. Nature communications. Nature Publishing Group;
680    2015;6:1–8.

681    22. Castelle CJ, Banfield JF. Major new microbial groups expand diversity and alter our understanding of
682    the tree of life. Cell. Elsevier; 2018;172:1181–97.

683    23. Gleeson T, Befus KM, Jasechko S, Luijendijk E, Cardenas MB. The global volume and distribution of
684    modern groundwater. Nature Geoscience. 2016;9:161–7.

685    24. Krüger M, Potthast K, Michalzik B, Tischer A, Küsel K, Deckner FF, et al. Drought and rewetting
686    events enhance nitrate leaching and seepage-mediated translocation of microbes from beech forest soils.
687    Soil Biology and Biochemistry. Elsevier; 2021;154:108153.

688    25. Albertsen M, Hugenholtz P, Skarshewski A, Nielsen KL, Tyson GW, Nielsen PH. Genome sequences
689    of rare, uncultured bacteria obtained by differential coverage binning of multiple metagenomes. Nature
690    biotechnology. Nature Publishing Group; 2013;31:533–8.

691    26. Gong J, Qing Y, Guo X, Warren A. "Candidatus Sonnebornia yantaiensis", a member of candidate
692    division OD1, as intracellular bacteria of the ciliated protist Paramecium bursaria (Ciliophora,
693    Oligohymenophorea). Systematic and applied microbiology. Elsevier; 2014;37:35–41.

694    27. Bor B, Collins A, Murugkar P, Balasubramanian S, To T, Hendrickson E, et al. Insights obtained by
695    culturing Saccharibacteria with their bacterial hosts. Journal of dental research. SAGE Publications Sage
696    CA: Los Angeles, CA; 2020;99:685–94.

697    28. Cross KL, Campbell JH, Balachandran M, Campbell AG, Cooper SJ, Griffen A, et al. Targeted
698    isolation and cultivation of uncultivated bacteria by reverse genomics. Nature biotechnology. Nature
699    Publishing Group; 2019;37:1314–21.

700    29. Murugkar PP, Collins AJ, Chen T, Dewhirst FE. Isolation and cultivation of candidate phyla radiation
701    Saccharibacteria (TM7) bacteria in coculture with bacterial hosts. Journal of oral microbiology. Taylor &
702    Francis; 2020;12:1814666.

703    30. Giovannoni SJ, Thrash JC, Temperton B. Implications of streamlining theory for microbial ecology.
704    The ISME journal. Nature Publishing Group; 2014;8:1553–65.

705    31. Dufresne A, Garczarek L, Partensky F. Accelerated evolution associated with genome reduction in a
706    free-living prokaryote. Genome biology. Springer; 2005;6:1–10.

707    32. Küsel K, Totsche KU, Trumbore SE, Lehmann R, Steinhäuser C, Herrmann M. How deep can surface
708    signals be traced in the critical zone? Merging biodiversity with biogeochemistry research in a central
709    German Muschelkalk landscape. Frontiers in Earth Science. Frontiers; 2016;4:32.

710    33. Yan L, Herrmann M, Kampe B, Lehmann R, Totsche KU, Küsel K. Environmental selection shapes
711    the formation of near-surface groundwater microbiomes. Water research. Elsevier; 2020;170:115341.

712    34. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal RNA gene
713    database project: improved data processing and web-based tools. Nucleic acids research. Oxford
714    University Press; 2012;41:D590–6.

715    35. Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, et al. Anvi'o: an advanced
716    analysis and visualization platform for 'omics data. PeerJ. PeerJ Inc.; 2015;3:e1319.

717    36. Eren AM, Kiefl E, Shaiber A, Veseli I, Miller SE, Schechter MS, et al. Community-led, integrated,
718    reproducible multi-omics with anvi'o. Nature microbiology. Nature Publishing Group; 2021;6:3–6.

719    37. Campbell JH, O'Donoghue P, Campbell AG, Schwientek P, Sczyrba A, Woyke T, et al. UGA is an
720    additional glycine codon in uncultured SR1 bacteria from the human microbiota. Proceedings of the
721    National Academy of Sciences. National Acad Sciences; 2013;110:5540–5.

722    38. Kersey CM, Agyemang PA, Dumenyo CK. CorA, the magnesium/nickel/cobalt transporter, affects
723    virulence and extracellular enzyme production in the soft rot pathogen Pectobacterium carotovorum.
724    Molecular plant pathology. Wiley Online Library; 2012;13:58–71.

725    39. Porcheron G, Garénaux A, Proulx J, Sabri M, Dozois CM. Iron, copper, zinc, and manganese
726    transport and regulation in pathogenic Enterobacteria: correlations between strains, site of infection and
727    the relative importance of the different metal transport systems for virulence. Frontiers in cellular and
728    infection microbiology. Frontiers; 2013;3:90.

729    40. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. Journal of
730    molecular biology. Elsevier; 1990;215:403–10.

731    41. Yan L, Hermans SM, Totsche KU, Lehmann R, Herrmann M, Küsel K. Groundwater bacterial
732    communities evolve over time in response to recharge. Water Research. Elsevier; 2021;117290.

733    42. Tian R, Ning D, He Z, Zhang P, Spencer SJ, Gao S, et al. Small and mighty: adaptation of
734    superphylum Patescibacteria to groundwater environment drives their genome simplicity. Microbiome.
735    BioMed Central; 2020;8:1–15.

736    43. Castelle CJ, Brown CT, Thomas BC, Williams KH, Banfield JF. Unusual respiratory capacity and
737    nitrogen metabolism in a Parcubacterium (OD1) of the Candidate Phyla Radiation. Scientific reports.
738    Nature Publishing Group; 2017;7:1–12.

739    44. Kohlhepp B, Lehmann R, Seeber P, Küsel K, Trumbore SE, Totsche KU. Aquifer configuration and
740    geostructural links control the groundwater quality in thin-bedded carbonate–siliciclastic alternations of
741    the Hainich CZE, central Germany. Hydrology and Earth System Sciences. Copernicus GmbH;
742    2017;21:6091–116.

743    45. Opitz S, Küsel K, Spott O, Totsche KU, Herrmann M. Oxygen availability and distance to surface
744    environments determine community composition and abundance of ammonia-oxidizing prokaroytes in
745    two superimposed pristine limestone aquifers in the Hainich region, Germany. FEMS Microbiology
746    Ecology. 2014;90:39–53.

747    46. Wegner C-E, Gaspar M, Geesink P, Herrmann M, Marz M, Küsel K, et al. Biogeochemical Regimes
748    in Shallow Aquifers Reflect the Metabolic Coupling of the Elements Nitrogen, Sulfur, and Carbon.
749    Applied and Environmental Microbiology. 2019;85:e02346-18.

750    47. Ozawa K, Meikari T, Motohashi K, Yoshida M, Akutsu H. Evidence for the presence of an F-type
751    ATP synthase involved in sulfate respiration in Desulfovibrio vulgaris. J Bacteriol. American Society for
752    Microbiology; 2000;182:2200–6.

753    48. Nelson WC, Stegen JC. The reduced genomes of Parcubacteria (OD1) contain signatures of a
754    symbiotic lifestyle. Frontiers in microbiology. Frontiers; 2015;6:713.

755    49. Levin PA, Angert ER. Small but mighty: cell size and bacteria. Cold Spring Harbor Perspectives in
756    Biology. Cold Spring Harbor Lab; 2015;7:a019216.

757    50. Kempes CP, Wang L, Amend JP, Doyle J, Hoehler T. Evolutionary tradeoffs in cellular composition
758    across diverse bacteria. The ISME journal. Nature Publishing Group; 2016;10:2145–57.

759    51. Melville S, Craig L. Type IV pili in Gram-positive bacteria. Microbiology and molecular biology
760    reviews. Am Soc Microbiol; 2013;77:323–41.

761    52. Craig L, Volkmann N, Arvai AS, Pique ME, Yeager M, Egelman EH, et al. Type IV pilus structure by
762    cryo-electron microscopy and crystallography: implications for pilus assembly and functions. Molecular
763    cell. Elsevier; 2006;23:651–62.

764    53. Mandlik A, Swierczynski A, Das A, Ton-That H. Pili in Gram-positive bacteria: assembly,
765    involvement in colonization and biofilm development. Trends in microbiology. Elsevier; 2008;16:33–40.

766    54. Xie B, Wang J, Nie Y, Chen D, Hu B, Wu X, et al. EpicPCR-Directed Cultivation of a Candidatus
767    Saccharibacteria Symbiont Reveals a Type IV Pili-dependent Epibiotic Lifestyle. bioRxiv [Internet]. Cold
768    Spring Harbor Laboratory; 2021; Available from:
769    https://www.biorxiv.org/content/early/2021/07/08/2021.07.08.451036

770    55. Coleman GA, Davín AA, Mahendrarajah TA, Szánthó LL, Spang A, Hugenholtz P, et al. A rooted
771    phylogeny resolves early bacterial evolution. Science. American Association for the Advancement of
772    Science; 2021;372.

773  56. Lehmann K, Lehmann R, Totsche KU. Event-driven dynamics of the total mobile inventory in
774  undisturbed soil account for significant fluxes of particulate organic carbon. Science of The Total
775  Environment. Elsevier; 2021;756:143774.

776  57. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. metaSPAdes: a new versatile metagenomic
777  assembler. Genome research. Cold Spring Harbor Lab; 2017;27:824–34.

778  58. Ultsch A, Mörchen F. ESOM-Maps: tools for clustering, visualization, and classification with
779  Emergent SOM. DATA BIONICS RESEARCH GROUP, UNIVERSITY OF MARBURG; 2005.

780  59. Dick GJ, Andersson AF, Baker BJ, Simmons SL, Thomas BC, Yelton AP, et al. Community-wide
781  analysis of microbial genome sequence signatures. Genome biology. Springer; 2009;10:1–16.

782  60. Wu Y-W, Simmons BA, Singer SW. MaxBin 2.0: an automated binning algorithm to recover
783  genomes from multiple metagenomic datasets. Bioinformatics. Oxford University Press; 2016;32:605–7.

784  61. Olm MR, Brown CT, Brooks B, Banfield JF. dRep: a tool for fast and accurate genomic comparisons
785  that enables improved genome recovery from metagenomes through de-replication. The ISME Journal.
786  2017;11:2864–8.

787  62. Gruber-Vodicka HR, Seah BK, Pruesse E. phyloFlash: Rapid small-subunit rRNA profiling and
788  targeted assembly from metagenomes. Msystems. Am Soc Microbiol; 2020;5:e00920-20.

789  63. Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH. GTDB-Tk: a toolkit to classify genomes with
790  the Genome Taxonomy Database. Oxford University Press; 2020.

791  64. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic gene
792  recognition and translation initiation site identification. BMC bioinformatics. BioMed Central;
793  2010;11:1–11.

794  65. Shaffer M, Borton MA, McGivern BB, Zayed AA, La Rosa SL, Solden LM, et al. DRAM for
795  distilling microbial metabolism to automate the curation of microbiome function. Nucleic acids research.
796  Oxford University Press; 2020;48:8883–900.

797  66. Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K. KEGG: new perspectives on genomes,
798  pathways, diseases and drugs. Nucleic acids research. Oxford University Press; 2017;45:D353–61.

799  67. Rawlings ND, Barrett AJ, Bateman A. MEROPS: the peptidase database. Nucleic acids research.
800  Oxford University Press; 2010;38:D227–33.

801  68. El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, et al. The Pfam protein families
802  database in 2019. Nucleic acids research. Oxford University Press; 2019;47:D427–32.

803  69. Zhang H, Yohe T, Huang L, Entwistle S, Wu P, Yang Z, et al. dbCAN2: a meta server for automated
804  carbohydrate-active enzyme annotation. Nucleic acids research. Oxford University Press; 2018;46:W95–
805  101.

806  70. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. Nucleic
807  acids research. Oxford University Press; 2004;32:1792–7.

808  71. Price MN, Dehal PS, Arkin AP. FastTree 2–approximately maximum-likelihood trees for large
809  alignments. PloS one. Public Library of Science San Francisco, USA; 2010;5:e9490.

810  72. Emiola A, Oh J. High throughput in situ metagenomic measurement of bacterial replication at ultra-
811  low sequencing coverage. Nature communications. Nature Publishing Group; 2018;9:1–8.

812  73. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna, Austria:
813  R Foundation for Statistical Computing; 2020. Available from: https://www.R-project.org/

814  74. Quinn TP, Richardson MF, Lovell D, Crowley TM. propr: an R-package for identifying
815  proportionally abundant features using compositional data analysis. Scientific reports. Nature Publishing
816  Group; 2017;7:1–9.

817  75. Csardi G, Nepusz T. The igraph software package for complex network research. InterJournal.
818  2006;Complex Systems:1695.

819  76. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: a software
820  environment for integrated models of biomolecular interaction networks. Genome research. Cold Spring
821  Harbor Lab; 2003;13:2498–504.

822  77. Gustavsen JA, Pai S, Isserlin R, Demchak B, Pico AR. RCy3: Network biology using Cytoscape from
823  within R. F1000Research. Faculty of 1000 Ltd; 2019;8.

824  78. Kassambara A. rstatix: Pipe-Friendly Framework for Basic Statistical Tests [Internet]. 2020.
825  Available from: https://CRAN.R-project.org/package=rstatix

826  ## Figure Legends

827  **Figure 1: Community composition of the groundwater samples based on metagenomic reads**

828  **mapped against the SILVA (SSU rRNA Ref NR99) database.** Each column represents a metagenomic

829  sample replicate for specified filter fractions from respective wells of the limestone-mudstone strata that

830  host the multi-story upper aquifer assemblage (HTU; wells H14, H32, H43, H52) and the karstified main

831  aquifer (HTL; wells H41, H51).

832  **Figure 2: Community composition showing taxonomic preferences of Patescibacteria in wells and**

833  **filter fractions across the Hainich transect.** The cross section of the studied groundwater transect (from

834  Kohlhepp *et al.*, 2017 [44], modified) shows the karstified main aquifer [HTL; (wells studied: H41, H51)]

835  that is characterized by higher surface-connection to preferential recharge areas and the hanging thin-

836  bedded alternating limestone-mudstone strata that host the multi-story upper aquifer assemblage (HTU;

837  wells studied H14, H32, H43, H52). Height above mean sea level (amsl), in meters, is shown along the y-

838     axis and length of hillslope is shown in meters along the x-axis. The colored pie charts show percentages

839     of taxa within Patescibacteria at order level. The underlined taxon, *Parcubacteria;other* (all Parcubacteria

840     other than the mentioned Parcubacteria orders merged together) was most abundant among

841     Patescibacteria in all the filter fractions of all the wells. The grey pie charts show the relative percentage

842     of Patescibacteria in the total community. The table includes levels of various hydrochemical parameters

843     of the studied wells, including the dissolved oxygen, measured during July 2014 - April 2017 [33].


844     **Figure 3: Phylogenetic placement of Patescibacteria MAGs after binning and refinement. A.**

845     Genome completeness distribution of the MAGs classified as Patescibacteria by GTDB-Tk alone (174,

846     orange-colored bars), and by both GTDB-Tk and Anvi'o (291, teal-colored bars). **B.** Phylogenetic tree

847     based on 68 core proteins from all bacterial MAGs (1087) using Maximum Likelihood in FastTree2 with

848     1000 bootstrap replications. Bacterial taxa other than Patescibacteria were collapsed together and only

849     Patescibacteria are colored as per their taxonomic assignments from GTDB-Tk. The bootstrap values of 0.9

850     and above are indicated by filled circles. The phylogenetic tree is supplied as Additional file 3.


851     **Figure 4: Distribution of genome sizes of Patescibacteria MAGs enriched in 0.1 μm and 0.2 μm**

852     **filter fractions. A**. For all 291 high-quality Patescibacteria MAGs, the ratio of average normalized

853     genome coverage in 0.1 μm filter fractions to 0.2 μm filter fractions from metagenomes was used to form

854     three groups: '0.1 μm filter' - MAGs where this ratio was at least 5, '0.2 μm filter' - MAGs where this

855     ratio was ⅕ or less, and 'None' - MAGs other than first two groups. The mean genome sizes were

856     significantly different (Kruskal-Wallis rank sum test, p = 2.24e-06). Pairwise Dunn's test showed the

857     genome sizes were significantly different between '0.1 μm filter' and '0.2 μm filter' (fdr adjusted p =

858     1.02e-06), and between '0.2 μm filter' and 'None' (fdr adjusted p = 9.08e-05). **B.** The scatter plot shows

859     the distribution of $\log_2$ filter enrichment factors (the ratio of average normalized genome coverage in 0.2

860     μm filter fractions to 0.1 μm filter fractions from metagenomes) of Patescibacteria MAGs, as the function

861     of their genome sizes. The dashed lines indicate the cut-off value of 5 and ⅕ for filter enrichment factors

862     on the y-axis.

863   **Figure 5: The estimated growth rate index (GRiD) distribution of Patescibacteria MAGs across the**

864   **metagenomes. A.** Well-wise GRiD distribution of all Patescibacteria. **B.** GRiD distribution of classes of

865   Patescibacteria in 0.1 μm and 0.2 μm filter fractions. The statistical significance was calculated by using

866   the t_test function with FDR correction in R package *rstatix* [78].


867   **Figure 6: Metabolic and functional repertoire of the high quality Patescibacteria MAGs.** The

868   heatmap shows completeness of pathways and presence/absence of the functions in 291 high-quality

869   Patescibacteria genomes annotated within DRAM [65], arranged according to their phylogenetic

870   placement. Clade background colors within the phylogenetic tree represent respective taxonomic classes

871   of Patescibacteria. Colored triangles next to each genome represent their enrichment in 0.1 μm filter

872   fractions (green), 0.2 μm filter fractions (red), anoxic wells (blue) and oxic wells (orange), respectively.

873   Electron transport chain complexes I-IV, sulfur metabolism functions, and photosynthesis related genes

874   were absent from almost all the MAGs. A similar heatmap arranged as per the 5-fold enrichment of the

875   MAGs in oxic and anoxic wells is provided as Additional file 5, Figure S2.


876   **Figure 7: Cell schematic representing the functional repertoire of most abundant model**

877   **Patescibacteria from oxic and anoxic groundwater wells.** The common and genome specific gene

878   features are shown for the three representative genomes based on KEGG pathways. The pie diagrams next

879   to each reaction or function state the presence of respective enzymes or proteins in the three model

880   organisms as per the color key (oxic representatives in green and blue, and the anoxic representative in

881   pink), while absence is indicated by white color.


882   **Figure 8: Co-occurrence network among the MAGs recovered from the studied groundwater wells.**

883   The proportionality network was constructed using normalized average coverages of the MAGs enriched

884   (by 5-fold coverage difference) in 0.2 μm filter fractions as compared to 0.1 μm filter fractions to retain

885   Patescibacteria possibly attached to other microbial hosts. The filled oval regions highlight the direct one-

886   to-one associations of Patescibacteria MAGs paired with Omnitrophota MAGs. The zoomed-in cluster

887   shows direct associations of Patescibcteria MAGs (filled red circles) with multiple Nitrospirota (filled

888    blue circles) and Bacteroidota (filled cyan circles) MAGs highlighted with black outlines and arrows,

889    while grey outlines and arrows indicate indirect associations. For construction of the proportionality

890    network, ρ (rho) value cut-off of 0.95 was used.

## Additional files

892    **Additional file 1:** Single copy genes from publically available CPR genomes used to predict CPR MAGs

893    in this study. The file was taken from Anvi'o codebase (https://github.com/merenlab/anvio).

894    **Additional file 2:** Genome statistics and taxonomic assignments of Patescibacteria MAGs in this study.

895    **Additional file 3:** The Newick tree file for phylogenetic tree shown in Figure 3, B.

896    **Additional file 4, Figure S1:** Correlation of average normalized genome coverages of Patescibacteria

897    MAGs from respective wells with respective GRiD values.

898    **Additional file 5, Figure S2:** Metabolic and functional repertoire of high quality Patescibacteria MAGs.

899    The heatmap shows completeness of pathways and presence/absence of functions in 291 high-quality

900    Patescibacteria genomes annotated with DRAM, arranged according to their enrichment in oxic and

901    anoxic wells based on 5-fold coverage criterion.

902    **Additional file 6:** Correlations of average normalized genome coverages of Patescibacteria MAGs

903    enriched in oxic wells with dissolved oxygen and nitrate concentration.
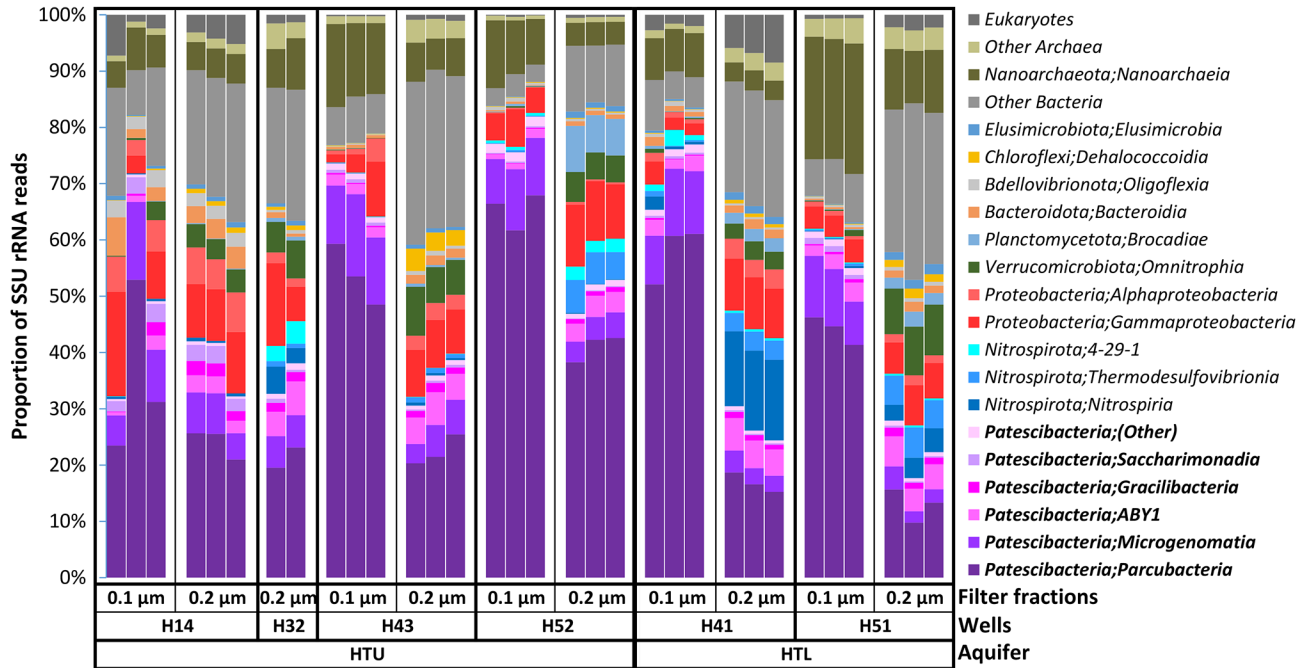
904    **Additional file 7:** Genomic coverages of 1275 microbial MAGs in all studied metagenomes.
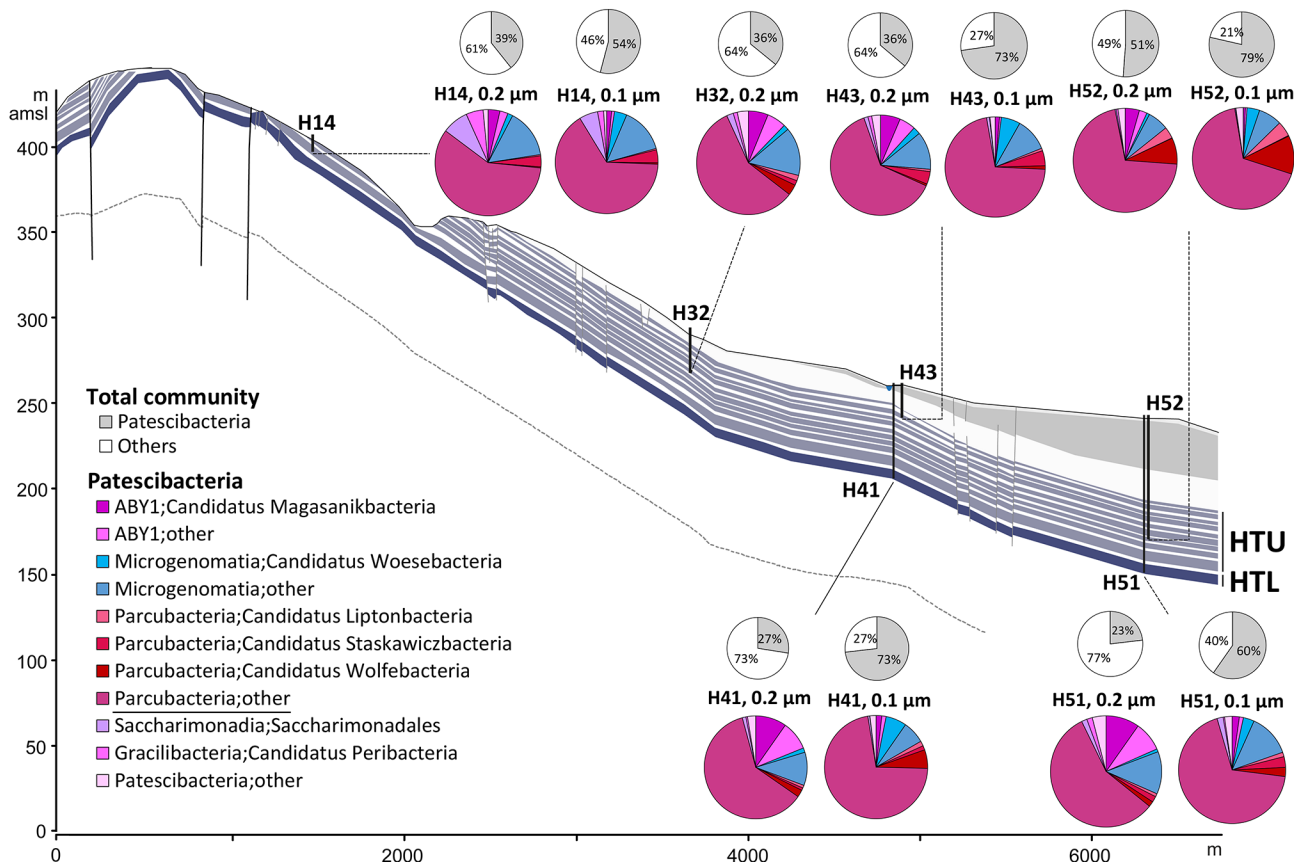
905    **Additional file 8, Figure S3:** Coverage distribution of selected MAGs from the network in Figure 8.

906    Only the direct one-to-one pairs of Patescibacteria with other MAGs are plotted.

907    **Additional file 9, Figure S4:** Coverage distribution of selected MAGs from the highlighted cluster in

908    network in Figure 8. Only the direct connections of Patescibacteria with other MAGs are plotted.
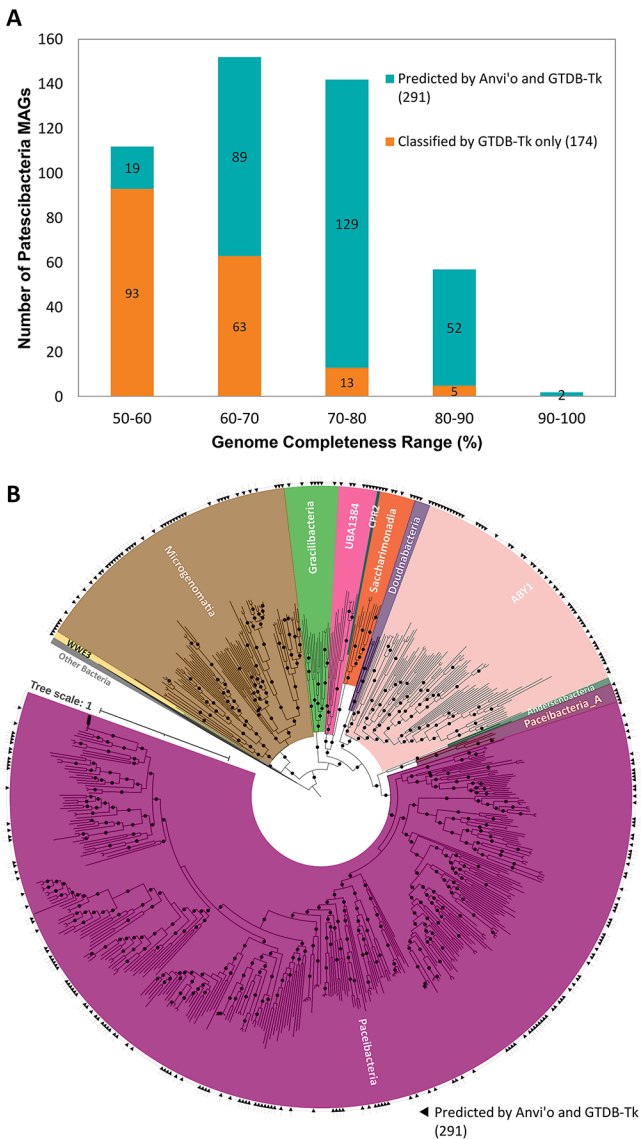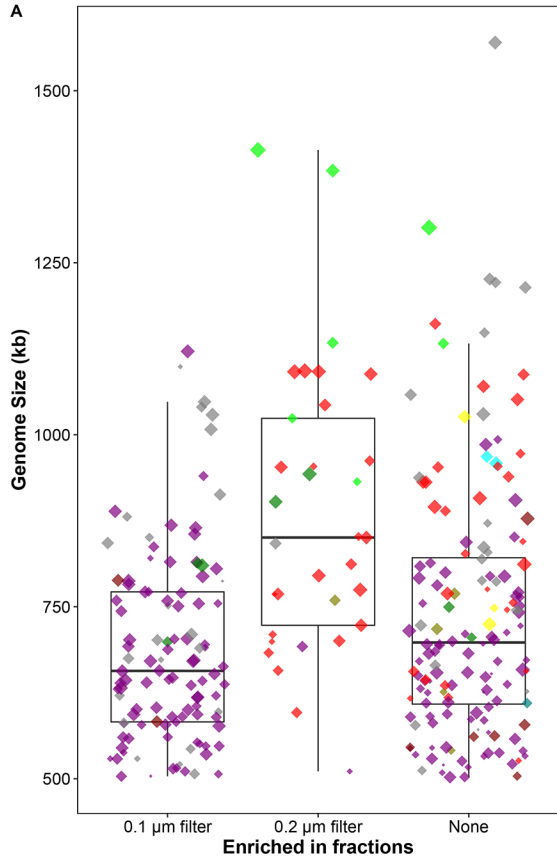
909

37

**Figure 1**

**Figure 2**



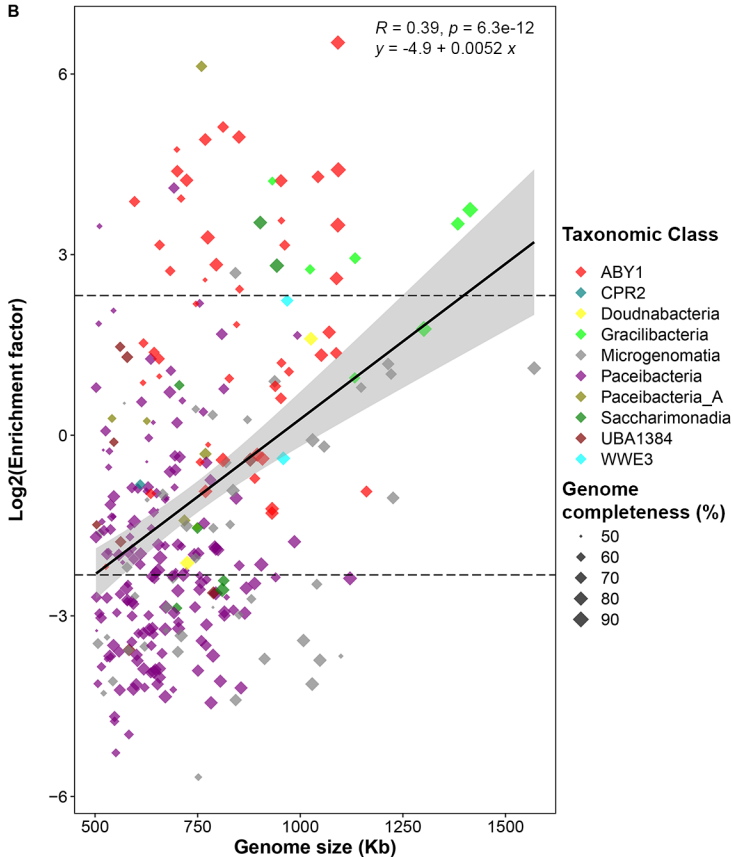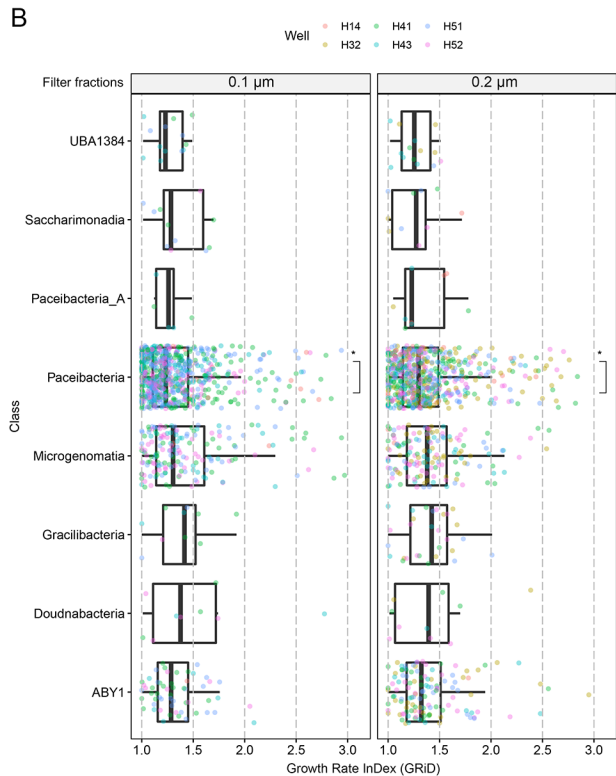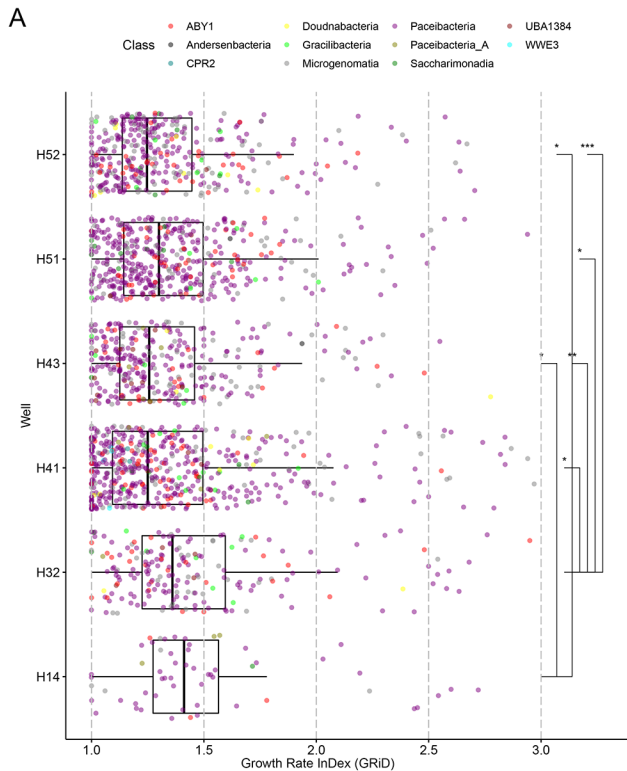| | | pH | Dissolved Oxygen (mg/L) | Ammonium (mg/L) | Nitrate (mg/L) | Sulphate (mg/L) |
|---|---|---|---|---|---|---|
| HTU | H14 | 6.98 ± 0.09 (6.8 - 7.2) | 0.61 ± 0.58 (0.1 - 2.54) | 0.01 ± 0.02 (0 - 0.06) | 1.29 ± 0.21 (0.77 - 1.52) | 26.72 ± 2.16 (23.88 - 30.2) |
| | H32 | 7.31 ± 0.07 (7.2 - 7.5) | 2.23 ± 0.56 (1.31 - 3.41) | 0.01 ± 0.02 (0 - 0.11) | 28.51 ± 8.22 (12.57 - 40.58) | 73.12 ± 5.25 (63.18 - 91.64) |
| | H43 | 7.14 ± 0.07 (7 - 7.3) | 0 | 0.09 ± 0.06 (0 - 0.27) | 1.55 ± 3.92 (0.01 - 11.99) | 38.52 ± 1.94 (35.15 - 47.06) |
| | H52 | 7.31 ± 0.06 (7.1 - 7.4) | 0 | 0.41 ± 0.1 (0.13 - 0.58) | 5.35 ± 4.37 (0.07 - 16.32) | 88.66 ± 8.22 (72.81 - 102.95) |
| HTL | H41 | 7.25 ± 0.17 (7.1 - 8.1) | 4.83 ± 1.7 (1.77 - 8.04) | 0.12 ± 0.1 (0 - 0.33) | 10.16 ± 4.41 (2.51 - 23.33) | 91.62 ± 20.76 (59.44 - 140.48) |
| | H51 | 7.15 ± 0.09 (6.9 - 7.3) | 2.73 ± 0.31 (2.21 - 3.29) | 0.04 ± 0.12 (0 - 0.68) | 8.12 ± 3.27 (4.87 - 21.05) | 289.47 ± 19.95 (253.96 - 337.19) |

# Figure 3

**Figure 4**



Legend (Figure B):

Taxonomic Class
- ABY1
- CPR2
- Doudnabacteria
- Gracilibacteria
- Microgenomatia
- Paceibacteria
- Paceibacteria_A
- Saccharimonadia
- UBA1384
- WWE3

Genome completeness (%)
- 50
- 60
- 70
- 80
- 90

Panel B annotation: $R = 0.39$, $p = 6.3\text{e-}12$; $y = -4.9 + 0.0052\,x$

# Figure 5

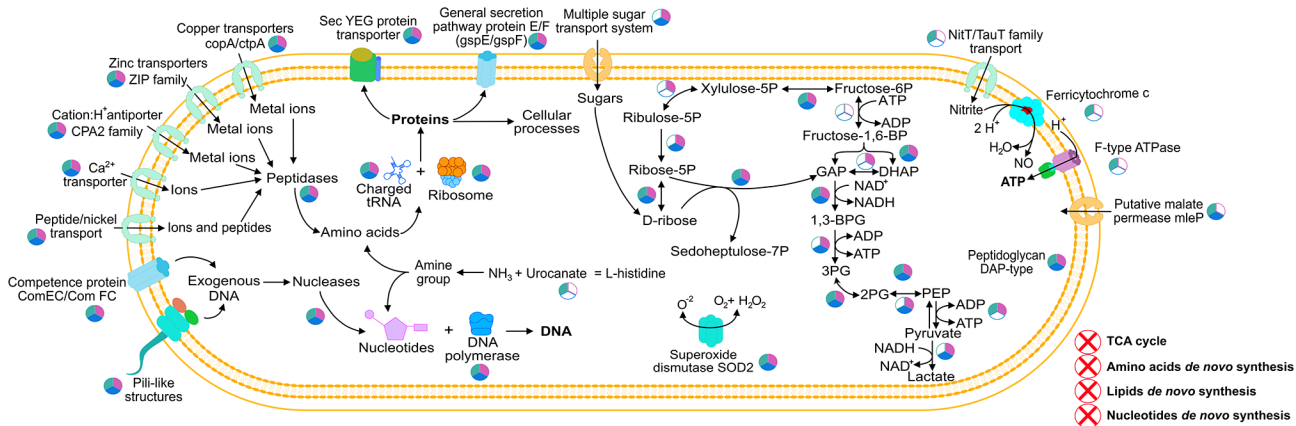# Figure 6

# Figure 7



| | Well | Filter fraction | Average normalized genome coverage | Length (bp) | # of contigs | N50 | GC content | % completion | % redundancy | Order | Family | Genus |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| H41-bin049 | H41 | 0.1μm | 482.45 | 582,449 | 29 | 36,229 | 38.65 | 56.33 | 0.00 | UBA6257 | UBA9933 | - |
| H52-bin095 | H52 | 0.1μm | 1764.92 | 502,359 | 21 | 45,086 | 42.87 | 74.65 | 0.00 | UBA6257 | UBA9933 | - |
| H41-bin288 | H41 | 0.1μm | 1080.23 | 620,744 | 45 | 18,731 | 48.57 | 74.65 | 0.00 | UBA9983_A | UBA2163 | C7867-001 |

**Figure 8**



Legend:
- Patescibacteria
- Nitrospirota
- Bacteroidota
- Chloroflexota
- Omnitrophota
- Planctomycetota
- Other Bacteria
- Archaea