

1 **Evolution in action: Habitat-transition leads to genome-streamlining**  
2 **in Methylophilaceae (Betaproteobacteriales)**

3

4 Running title: Habitat-transition leads to genome-streamlining

5

6 Michaela M. Salcher<sup>1,2\*</sup>, Daniel Schaeffle<sup>2,3</sup>, Melissa Kaspar<sup>2</sup>, Stefan M. Neuenschwander<sup>2</sup>,  
7 Rohit Ghai<sup>1</sup>

8

9 <sup>1</sup>Department of Aquatic Microbial Ecology, Institute of Hydrobiology, Biology Centre CAS, Na  
10 Sádkách 7, 37005 České Budějovice, Czech Republic

11 <sup>2</sup>Limnological Station, Institute of Plant and Microbial Biology, University of Zurich,  
12 Seestrasse 187, 8802 Kilchberg, Switzerland

13 <sup>3</sup>Institute of Medical Microbiology, University of Zurich, Gloriastrasse 28/30, 8006 Zurich,  
14 Switzerland

15

16 \*corresponding author:

17 Michaela M. Salcher

18 Department of Aquatic Microbial Ecology, Institute of Hydrobiology, Biology Centre CAS  
19 Na Sádkách 7, 37005 České Budějovice, Czech Republic

20 Phone: +420 387 775 836

21 Email: [michaelasalcher@gmail.com](mailto:michaelasalcher@gmail.com)

22

23 The authors declare that they have no competing interests.

## 24 **Abstract**

25 The most abundant aquatic microbes are small in cell and genome size. Genome-  
26 streamlining theory predicts gene loss caused by evolutionary selection driven by  
27 environmental factors, favouring superior competitors for limiting resources. However,  
28 evolutionary histories of such abundant, genome-streamlined microbes remain largely  
29 unknown. Here we reconstruct the series of steps in the evolution of some of the most  
30 abundant genome-streamlined microbes in freshwaters ('*Ca. Methylophilus*') and oceans  
31 (marine lineage OM43). A broad genomic spectrum is visible in the family Methylophilaceae  
32 (Betaproteobacteriales), from sediment microbes with medium-sized genomes (2-3 Mbp  
33 genome size), an occasionally blooming pelagic intermediate (1.7 Mbp), and the most  
34 reduced pelagic forms (1.3 Mbp). We show that a habitat transition from freshwater sediment  
35 to the relatively oligotrophic pelagial was accompanied by progressive gene loss and  
36 adaptive gains. Gene loss has mainly affected functions not necessarily required or  
37 advantageous in the pelagial or are encoded by redundant pathways. Likewise, we identified  
38 genes providing adaptations to oligotrophic conditions that have been transmitted  
39 horizontally from pelagic freshwater microbes. Remarkably, the secondary transition from the  
40 pelagial of lakes to the oceans required only slight modifications, i.e., adaptations to higher  
41 salinity, gained via horizontal gene transfer from indigenous microbes. Our study provides  
42 first genomic evidence of genome-reduction taking place during habitat transitions. In this  
43 regard, the family Methylophilaceae is an exceptional model for tracing the evolutionary  
44 history of genome-streamlining as such a collection of evolutionarily related microbes from  
45 different habitats is practically unknown for other similarly abundant microbes (e.g., '*Ca.*  
46 *Pelagibacterales*', '*Ca. Nanopelagicales*').

## 47 **Keywords**

48 Genome-streamlining, Genome reduction, Horizontal gene transfer, Evolutionary selection,  
49 Habitat transition, Comparative Genomics, Bacteria

## 50 Introduction

51 Marine and freshwater pelagic habitats are numerically dominated by very small  
52 microbes (cell volumes  $<0.1 \mu\text{m}^3$ ) that seem to be perfectly adapted to nutrient-poor  
53 (oligotrophic) conditions by successfully competing for dissolved organic matter and nutrients  
54 at low nM concentrations due to higher surface-to-volume ratios and superior transport  
55 systems [1]. Small-sized cells also enjoy other benefits such as reduced replication costs and  
56 mortality rates by size selective protistan predators [2]. The genomes of such oligotrophs are  
57 characterized by being very small (streamlined,  $<1.5 \text{ Mbp}$ ) with highly conserved core  
58 genomes and few pseudogenes, compacted intergenic spacers, reduced numbers of  
59 paralogs, and a low genomic GC content [3, 4]. While genetic drift has been proposed as the  
60 evolutionary mechanism behind the reduced genomes of symbionts, parasites and  
61 commensals, selection driven by environmental factors has been suggested as the primary  
62 driving force in the case of free-living oligotrophs [3]. The most abundant organisms on earth,  
63 bacteria of the marine SAR11 lineage ('Ca. Pelagibacter ubique', Alphaproteobacteria) serve  
64 as models for genome streamlining in the oceans [5] and their freshwater sister lineage LD12  
65 is also known to be of similarly small size [6, 7]. Other examples of aquatic microbes with  
66 small cells and reduced genomes can be found among Actinobacteria (marine 'Ca.  
67 Actinomarina minuta' [8], freshwater 'Ca. Nanopelagicales' [9, 10], freshwater luna1 lineage  
68 [11, 12]), Thaumarchaeota (marine 'Ca. Nitrosopelagicus brevis')[13], and  
69 Betaproteobacteriales (freshwater 'Ca. Methylopumilus planktonicus' [14], marine OM43  
70 lineage [15, 16]).

71 The latter are methylotrophs that are specialized in using reduced one-carbon ( $\text{C}_1$ )  
72 compounds like methanol, methylamine and formaldehyde as sole energy and carbon  
73 sources by means of a modular system of different pathways for their oxidation,  
74 demethylation and assimilation [17]. The family Methylophilaceae (Betaproteobacteriales) is  
75 among the most important methylotrophs playing a key role in the carbon cycle of aquatic  
76 habitats [17, 18]. Four genera are so far validly described (*Methylotenera*, *Methylobacillus*,

77 *Methylophilus*, *Methylovorus*) that mainly inhabit the sediment of freshwater lakes [19-22].  
78 Axenic strains have been also isolated from the pelagial of lakes ('*Ca. Methylopumilus*') [14]  
79 and oceans (lineage OM43) [15, 16, 23]. Freshwater '*Ca. Methylopumilus planktonicus*' are  
80 ubiquitous and highly abundant in lakes [24] with distinct maxima during diatom and/or  
81 cyanobacterial blooms [14, 25, 26], indicating that C<sub>1</sub> substrates supporting their growth are  
82 released from primary producers. Members of the coastal marine OM43 lineage display  
83 similar temporal patterns with highest numbers during phytoplankton blooms [27-29].

84 In this work, we analysed the evolutionary history of the family Methylophilaceae by  
85 comparative genomics. While sediment dwellers have a larger cell and genome size, pelagic  
86 lineages are genome-streamlined. We hypothesize that the evolutionary origin of the family  
87 can be traced back to freshwater sediments, from where these microbes emerged to colonize  
88 the plankton of lakes and eventually also crossing the freshwater-marine boundary. The  
89 transition from sediments to the pelagial resulted in a pronounced genome reduction and  
90 adaptive gene loss has mainly affected functions that are not necessarily required or  
91 advantageous in the pelagial or are encoded in redundant pathways. Likewise, genes  
92 providing adaptations to oligotrophic conditions might have been transmitted horizontally from  
93 indigenous pelagic microbes.

94

## 95 **Material and Methods**

### 96 *Isolation of planktonic freshwater methylotrophs.*

97 Novel strains of '*Ca. Methylopumilus*' and other Methylophilaceae were isolated from the  
98 pelagial of Lake Zurich (CH), Římov Reservoir (CZ), and Lake Medard (CZ). Dilution-to-  
99 extinction using filtered (0.2 µm) and autoclaved water amended with vitamins and amino  
100 acids as a medium was used for Lake Zurich [14]. A full-cycle isolation approach [30] was  
101 employed for samples from Římov Reservoir and Lake Medard, with filtered water samples  
102 (0.45 µm filters), being diluted 1:10 with Artificial Lake Water (ALW [31]) containing vitamins

103 (0.593  $\mu$ M thiamine, 0.08  $\mu$ M niacin, 0.000074  $\mu$ M cobalamine, 0.005  $\mu$ M para-amino  
104 benzoic acid, 0.074  $\mu$ M pyridoxine, 0.081  $\mu$ M pantothenic acid, 0.004  $\mu$ M biotin, 0.004  $\mu$ M  
105 folic acid, 0.555  $\mu$ M myo-inosito, 0.01  $\mu$ M riboflavin), 30  $\mu$ M LaCl<sub>3</sub>, 1 mM methanol and 0.1  
106 mM methylamine and incubated for 1-2 days at *in situ* temperatures. This step resulted in a  
107 preadaptation of methylotrophs only (C<sub>1</sub> compounds as sole carbon source) without causing  
108 a shift in the assemblage of 'Ca. Methylopusillus' as these microbes displays slow growth  
109 with doubling times of approx. two days. Thereafter, a dilution to extinction technique was  
110 employed [14] with approx. 1 cell per cultivation well in 24-well-plates. Plates were incubated  
111 for 4-6 weeks at *in situ* temperature and growth in individual wells was checked  
112 microscopically and by PCR and Sanger sequencing of 16S rRNA genes.

### 113 *Whole-genome sequencing, assembly, and functional annotation.*

114 Thirty-eight pure cultures of 'Ca. Methylopusillus sp.' and three *Methylophilus* sp. were  
115 grown in 400 ml ALW medium supplemented with vitamins, LaCl<sub>3</sub>, methanol and  
116 methylamine for 6-8 weeks, pelleted by centrifugation, and DNA was isolated with a  
117 MagAttract® HMW DNA Kit (Qiagen). 550-bp libraries were constructed with the KAPA  
118 Hyper Prep Kit (Roche) and paired-end sequences (2 x 250-bp) were generated on an  
119 Illumina MiSeq instrument with a 500-cycle MiSeq Reagent v2 kit (Illumina). Library  
120 preparation and sequencing was done at the Genetic Diversity Center Zurich (GDC). Raw  
121 reads were quality trimmed with trimmomatic [32], assembled with SPAdes [33] and  
122 subsequently mapped to the resulting assemblies with Geneious 9 ([www.geneious.com](http://www.geneious.com)) in  
123 order to identify potential assembly errors. Assembly usually resulted in 1-2 large contigs  
124 with overlapping ends that mostly could be circularized *in silico*. In the case of non-  
125 overlapping contigs, genomes were closed by designing specific primers for PCR and  
126 Sanger sequencing. Moreover, regions containing low coverage ( $\leq 10$  fold), ambiguities, or  
127 anomalies in the mapping were verified by designing specific primers for PCR and Sanger  
128 sequencing to produce high-quality reference genomes. Gene prediction was done with  
129 PROKKA [34] and annotation was done with an in-house pipeline [10] based on BLAST

130 searches to NCBI NR, COG [35], TIGRFAM [36] and KEGG databases [37]. Metabolic  
131 pathways were inferred from KEGG [37] and MetaCyc [38] and manually examined for  
132 completeness. Pathways involved in methylotrophy were identified by collecting 1016  
133 reference protein sequences from published genomes of methylotrophs [14, 17, 39-45] and  
134 for the sake of completeness, also pathways not common to Methylophilaceae were included  
135 (e.g., methane oxidation [46-48], methylovory [49, 50]). These proteins were classified into  
136 25 modules representing distinct (or sometimes alternative) biochemical transformations  
137 relevant to a methylotrophic lifestyle (e.g. M01-methanol oxidation, M02-pyrroloquinoline  
138 quinone biosynthesis, etc.; a complete list is provided in Supplementary Table S5). Protein  
139 sequences were clustered at 90% identity and 90% coverage with cd-hit [51] and the clusters  
140 were aligned using muscle [52]. The alignments were converted to HMMs (Hidden Markov  
141 Models) using the hmmbuild program in the HMMER3 package [53]. The program  
142 hmmsearch was used to scan complete genomes using these HMMs using e-value cut-off of  
143 1e-3. The entire set of HMMs is available as Supplementary Data Set.

#### 144 *Fragment recruitment from metagenomes.*

145 Publicly available metagenomes gained from freshwater sediments (n=131), the pelagial of  
146 lakes (n=345), rivers (n=43), estuaries, brackish and coastal oceanic sites (n=53) as well as  
147 open oceans (n=201) were used for fragment recruitment (see Table S2 for sampling sites  
148 and SRR accessions). rRNA sequences in genomes were identified with barrnap  
149 (<http://www.vicbioinformatics.com/software.barrnap.shtml>) and masked to avoid biases, and  
150 metagenomic reads were queried against the genomes using BLASTN [54] (cut-offs: length  
151  $\geq 50$  bp, identity  $\geq 95\%$ , e-value  $\leq 1e-5$ ). These hits were used to compute RPKG values  
152 (number of reads recruited per kb of genome per Gb of metagenome), which provide a  
153 normalized value that is comparable across different genomes and metagenomes.

#### 154 *Comparative genomic analyses.*

155 All publicly available genomes of high quality ( $>95\%$  completeness,  $<20$  scaffolds) affiliated  
156 with the family Methylophilaceae (Table S1, n=37) were downloaded from NCBI and re-

157 annotated for comparative analyses. Average nucleotide identities (ANI [55]) and average  
158 amino acid identities (AAI [56]) were calculated to discriminate different species and genera.  
159 Phylogenomic trees based on conserved concatenated protein sequences (351,312 amino  
160 acid sites from 878 proteins for all Methylophilaceae, Fig. 1; 337,501 amino acid sites from  
161 983 proteins for all 'Ca. Methylopumilus' spp., Fig. S1) was generated with FastTree [57] (100  
162 bootstraps) after alignment with kalign [58]. *Methyloversatilis* sp. RAC08 (NZ\_CP016448) and  
163 'Ca. Methylosempumilus turicensis' MMS-10A-171 (NZ\_LN794158) served as outgroup for  
164 the trees displayed in Fig. 1 and Fig. S1, respectively. The core- and pangenome of the family  
165 was computed using all-vs-all comparisons of all proteins for each genome using BLASTP  
166 ( $\geq 50\%$  identity and  $\geq 50\%$  coverage cut-offs to define an ortholog). Paralogs in each genome  
167 were identified with BLASTP (cut-offs:  $\geq 80\%$  coverage,  $\geq 70\%$  similarity,  $\geq 50\%$  identity).  
168 Closest relatives for proteins putatively transferred horizontally were identified with BLASTP  
169 against the NCBI Protein Reference Sequences database (cut-off:  $E$  values  $\leq 1e-5$ ). Trees for  
170 individual proteins or concatenated proteins for specific pathways were constructed with  
171 RAxML (GAMMA BLOSUM62 model [59]) after alignment with MAFFT v7.388 [60].

#### 172 *Availability of data*

173 All genomes have been submitted to NCBI under BioProject XXX, BioSamples XY-XYX  
174 (Please note: Submission is in progress, accession numbers will be provided as soon as  
175 possible). The entire set of HMMs related to methylotrophic functions is available as  
176 Supplementary Data Set.

177

## 178 **Results and discussion**

### 179 *Phylogenomics and global occurrence of Methylophilaceae*

180 Currently, 31 Methylophilaceae genomes of high quality (i.e.,  $>99\%$  completeness,  $<20$   
181 scaffolds) are publicly available, mostly from axenic isolates from freshwater sediments (Fig.  
182 1, Table S1). We additionally sequenced the genomes of 41 strains of planktonic freshwater

183 strains affiliated with ‘*Ca. Methylopusillus planktonicus*’ (38 strains) and *Methylophilus* sp. (3  
184 strains). These microbes were isolated from the pelagial of three different freshwater habitats  
185 (Lake Zurich, CH; Římov Reservoir, CZ; Lake Medard, CZ) by dilution-to-extinction [14, 30].  
186 All novel genomes are of very high quality, i.e., they are complete, with one circular  
187 chromosome (Table S1). The 39 strains classified as ‘*Ca. M. planktonicus*’ by 16S rRNA gene  
188 sequences analysis (99.94-100% sequence identity), constitute at least three different species  
189 according to average nucleotide and amino acid identity (ANI and AAI [61]) (Fig. S1-S4). We  
190 tentatively name these three taxa ‘*Ca. Methylopusillus rimovensis*’ (two strains isolated from  
191 Římov Reservoir), ‘*Ca. Methylopusillus universalis*’ (29 strains from Lake Zurich and Římov  
192 Reservoir) and the originally described ‘*Ca. Methylopusillus planktonicus*’ (eight strains from  
193 Lake Zurich; Fig. 1, Fig. S1)[14]. AAI values also suggest that ‘*Ca. M. turicensis*’ is a different  
194 genus (62% AAI with ‘*Ca. Methylopusillus*’, Fig. S3) that we tentatively rename to ‘*Ca.*  
195 *Methylosemipusillus turicensis*’. This reclassification is in line with the recently released  
196 Genome Taxonomy Database (GTDB [62]). Moreover, the genus *Methylopusillus* might be split  
197 in different genera and the GTDB suggests a reclassification of several strains to the genus  
198 *Methylophilus*. Our analysis notes a polyphyletic pattern of *Methylopusillus* with three different  
199 genera (*Methylopusillus-1*, *Methylopusillus-2*, and *Methylopusillus-3*; Fig. 1a, Fig. S3, >70% AAI).  
200 However, further work is necessary to clarify the formal naming of these strains as AAI values  
201 are inconclusive and the proposed hard cut-off of 65% AAI for genus delineation are not met  
202 for most members of Methylophilaceae. Three novel pelagic *Methylophilus* sp. isolates (MMS-  
203 M-34, MMS-M-37, MMS-M-51) constitute a novel species that we tentatively named *M.*  
204 *medardicus*, with closest hits to isolates from freshwater sediment. These strains might  
205 originate from the same clone, as they were gained from the same sample from Lake Medard  
206 and were 100% identical in their genome sequence. *M. medardicus* seem to be not  
207 abundant in the pelagial of lakes, as indicated by recruitments from 345 different pelagic  
208 freshwater metagenomic datasets, however they could be readily detected in relatively high  
209 proportions in sediment metagenomes (Fig. 1b, Table S2). Sediments also appear to be the  
210 main habitat of other *Methylophilus* and *Methylopusillus*. The three strains isolated from marine



211 systems, that were referred to as OM43-lineage [15, 16, 23], form two different genera based  
212 on AAI (Fig. S3). However, none appear to be abundant in the open ocean (Fig. 1b), and only  
213 strain HTCC2181 could be detected in estuarine/coastal metagenomes, although lineage  
214 OM43 has been repeatedly reported in coastal oceans by CARD-FISH, where they can reach  
215 up to 4% or  $0.8 \times 10^5$  cells ml<sup>-1</sup> during phytoplankton blooms [28, 29]. It is thus likely that other,  
216 more abundant strains of OM43 still await isolation. ‘*Ca. Methylopumilus* spp.’ on the other  
217 hand, were found in moderate proportions in estuarine/coastal systems, but their main habitat  
218 is clearly the pelagial of lakes, where they are highly abundant (Fig. 1b), as previously  
219 reported based on CARD-FISH analyses [14, 63], 16S rRNA gene amplicon sequencing [24,  
220 64-66], and metagenomics [67, 68]. All ‘*Ca. Methylopumilus*’, especially ‘*Ca. M. rimovensis*’  
221 were also prevalent in rivers (Fig. 1b).

#### 222 *Genome-streamlining in pelagic strains*

223 The genomes of pelagic freshwater ‘*Ca. Methylopumilus* sp.’ (n=39) and the marine  
224 OM43 lineage (n=3) are characterized by very small sizes (1.26-1.36 Mbp) and a low genomic  
225 GC content (35.3-37.7%) (Table S1, Fig. 2, Fig. S5). ‘*Ca. Methylosempumilus turicensis*’  
226 MMS-10A-171 has a slightly larger genome (1.75 Mbp) with higher GC content (44.5%), while  
227 all other Methylophilaceae have genome sizes >2.37 Mbp (max. 3.25 Mbp) and a higher GC  
228 content (41.9-55.7%, average 47.3%). A highly significant relationship between genome size  
229 and GC content, length of intergenic spacers, coding density, mean CDS length, number of  
230 overlapping CDS, paralogs, and numbers of genes involved in sensing of the environment  
231 (i.e., histidine kinases and sigma factors) was evident (Fig. 2). All these features have been  
232 proposed to be relevant for genome-streamlining [3] with freshwater ‘*Ca. Methylopumilus*’ and  
233 marine OM43 displaying the most reduced forms and ‘*Ca. Methylosempumilus turicensis*’  
234 presenting an intermediate state (Table S1). Moreover, we observed a negative relationship  
235 between genomic GC content and stop-codon usage of TAA instead of TAG, as well as a  
236 preferred amino acid usage of lysine instead of arginine (Fig. 2), both suggested to be  
237 involved in nitrogen limitation [4]. Furthermore, amino acids with less nitrogen and sulphur and

238 more carbon atoms were favourably encoded by the genome-streamlined microbes (Fig. 2,  
239 S5, S6).

#### 240 *Adaptive gene loss during habitat transition from the sediment to the pelagial*

241 The core genome of the family Methylophilaceae consists of 664 protein families (4.3%  
242 of the pangenome) and an open pangenome of >15,000 protein families, while the  
243 streamlined genomes of 'Ca. Methylopumilus' have a highly conserved core (48%, Fig. S7).  
244 By contrast, sediment Methylophilaceae have a larger pangenome with a more modular  
245 assortment featuring several redundant pathways for methylotrophic functions [18] and a  
246 large fraction of proteins overlapping with 'Ca. M. turicensis', indicating a high evolutionary  
247 relatedness (Fig. S7). It appears that both extant pelagic and sediment Methylophilaceae  
248 shared a common sediment-dwelling methylotrophic ancestor. While one lineage  
249 (*Methylotenera* and *Methylophilus*) retains the ancestral character (large genomes) of the  
250 common ancestor, the other lineage diversified towards a pelagic lifestyle ('Ca. M. turicensis',  
251 'Ca. Methylopumilus' and OM43). Remarkably, 'Ca. M. turicensis' appears to constitute an  
252 early diverging lineage that displays somewhat mixed characteristics of both sediment and  
253 truly pelagic forms ('Ca. Methylopumilus' and OM43), not only in its phylogenetic position, but  
254 also in genomic characteristics. Monitoring data from Lake Zurich showed consistently high  
255 cell densities of 'Ca. Methylopumilus', while 'Ca. M. turicensis' were mostly below detection  
256 limits except for a 3-month phase in one year where they reached high numbers in the  
257 hypolimnion [14]. Moreover, also fragment recruitment from freshwater metagenomes  
258 showed a global occurrence of 'Ca. Methylopumilus' in very high relative proportions, while  
259 'Ca. M. turicensis' was less prevalent (Fig. 1b). This hints again at the somewhat transitional  
260 character of 'Ca. M. turicensis' (occasional 'bloomer') that is not as perfectly adapted to the  
261 pelagial as 'Ca. Methylopumilus'.

262 We tested the hypothesis of adaptive gene loss driven by evolutionary selection during  
263 the transition from sediment to the pelagial by comparative genomics of metabolic modules of  
264 *Methylophilus*, *Methylotenera*, 'Ca. M. turicensis', 'Ca. Methylopumilus' and marine OM43

265 strains. *Methylobacillus* and *Methylovorus* were excluded as they seem too distantly related  
266 and also not very abundant in lake sediments (Fig. 1b, Table S2). The most pronounced  
267 differences in the genetic make-up of sediment vs. pelagic strains were detected in motility  
268 and chemotaxis (Figs. 3, 4), with all *Methylophilus* and all but two *Methylostenella* strains  
269 having flagella and type IV pili, while the planktonic strains have lost mobility and also greatly  
270 reduced the number of two-component regulatory systems and sigma factors. A large number  
271 of membrane transporters for inorganic compounds was detected exclusively in sediment  
272 *Methylophilaceae*, while this number is reduced in 'Ca. *M. turicensis*' and even more in 'Ca.  
273 *Methylopumilus*' and OM43 (Fig. 4, Table S3). Moreover, *Methylophilus* and *Methylostenella*  
274 encode multiple pathways for nitrogen acquisition, with transporters for ammonia, nitrate,  
275 nitrite, taurine, cyanate or urea, and pathways for urea or cyanate utilization (Fig. 3, Table  
276 S4)[69]. 'Ca. *M. turicensis*', on the other hand, has only transporters for nitrate/taurine and  
277 ammonia (Amt family), and 'Ca. *Methylopumilus*' and marine OM43 only carry ammonia  
278 transporters. All sediment *Methylophilaceae* further possess genes for assimilatory nitrate  
279 reduction to ammonia, some for dissimilatory nitrate reduction to nitrous oxide or detoxification  
280 of nitric oxide (quinol type of *norB*) [70] and *Methylostenella mobilis* 13 is a complete denitrifier  
281 [69, 71, 72], while none of the pelagic strains have any genes involved in nitrate reduction  
282 (Figs. 3, 4, Table S4). Ammonia is the main microbial nitrogen source in the epilimnion of  
283 lakes and oceans, while nitrate and other compounds like urea, taurine, or cyanate are more  
284 abundant in deeper, oxygenated layers and the sediment [14, 73, 74]. Therefore, an  
285 adaptation to ammonium uptake might be advantageous for pelagic microbes.

286 Furthermore, a high diversity of pathways involved in sulfur metabolism was detected in  
287 *Methylophilaceae*, with the genome-streamlined strains representing the most reduced forms  
288 again. All *Methylophilus*, *Methylostenella*, 'Ca. *Methylosemipumilus turicensis*' and one strain of  
289 'Ca. *Methylopumilus rimovensis*' encode ABC transporters for sulfate uptake, and a sulfate  
290 permease was annotated in OM43 and several sediment *Methylophilaceae*, while the majority  
291 of 'Ca. *Methylopumilus*' lack these transporters (Fig. 3, 4). Canonical assimilatory sulfate  
292 reduction seems to be incomplete in most *Methylophilaceae*, as adenylyl sulfate kinase *cysC*

293 was annotated only in a few strains (Table S4). Thus, the mode of sulfite generation remains  
294 unclear, with unknown APS kinases or other links from APS to sulfite. *Methylophilus*  
295 *rhizosphaera* encodes genes for dissimilatory sulfate reduction and most sediment strains  
296 possess ABC transporters for alkanesulfonates, most likely transporting methylsulfonate, that  
297 can be oxidised to sulfite by methanesulfonate monooxygenases generating formaldehyde as  
298 by-product. Dimethyl sulfide (DMS) seems to be a source for sulfur and formaldehyde as well,  
299 as dimethylsulfoxide and dimethyl monooxygenases are present in several sediment  
300 Methylophilaceae, but absent in all pelagic strains. It is thus still unclear how 'Ca.  
301 Methylopumilus' fuel their sulfur demand, especially as they grow in a defined medium  
302 containing sulfate (200  $\mu\text{M}$   $\text{MgSO}_4$ ; 160  $\mu\text{M}$   $\text{CaSO}_4$ ) and vitamins as sole sulfur sources.

303 All Methylophilaceae have complete pathways for the biosynthesis of amino acids and  
304 vitamins, with the exception of cobalamin (vitamin B12) that was lacking in the pelagic lineage  
305 ('Ca. *M. turicensis*', 'Ca. *Methylopumilus*', marine OM43), while either the complete  
306 biosynthesis or the salvage pathway was present in the sediment isolates (Figs. 3, 4, Table  
307 S4). However, putative cobalamin transporters were annotated in all isolates.

308 The methylcitric acid (MCA) cycle for oxidising propionate via methylcitrate to pyruvate is  
309 present in *Methylothenera*, *Methylophilus*, 'Ca. *M. turicensis*', and the marine OM43, but absent  
310 in all 'Ca. *Methylopumilus*' strains, suggesting it has been selectively lost in these organisms.  
311 All genes were arranged in a highly conserved fashion, with the exception of 'Ca. *M.*  
312 *turicensis*' having a bifunctional aconitate hydratase 2/2-methylisocitrate dehydratase  
313 (*prpD/acnB*) gene and *acnB* genes being located in different genomic regions in OM43 and  
314 'Ca. *M. turicensis*' (possessing two copies), however, with high synteny of flanking genes (Fig.  
315 S8). Phylogenetic analysis of the MCA gene cluster resulted in genus-specific branching, and  
316 notably, the MCA pathway of OM43 is most closely related to that of 'Ca. *M. turicensis*' (Fig.  
317 S8), suggesting it was retained in the OM43 lineage after divergence from a common ancestor  
318 of OM43 and 'Ca. *M. turicensis*'.

319 *Genome-streamlining leading to a loss of redundant methylotrophic pathways*

320 Some of the sediment dwellers seem to be facultative methylotrophs, as ABC  
321 transporters for amino acids were annotated (Figs. 3, 4, Table S3). *Methylotenera versatilis*  
322 301 additionally encodes a fructose-specific phosphotransferase system (PTS) and a 1-  
323 phosphofruktokinase, as well as transporters for putrescine uptake and the subsequent  
324 pathway for its degradation. 'Ca. M. turicensis' might also be a facultative methylotroph, as it  
325 possesses a PTS system for cellobiose, while this (as well as amino acid transporters) is  
326 lacking in all 'Ca. Methylopumilus' and OM43 strains, making them obligate methylotrophs.  
327 These observations suggest that the ancestor of both pelagic and sediment lineages was also  
328 a facultative methylotroph and that obligate methylotrophy emerged only in the truly pelagic  
329 strains.

330 Remarkably, also pathways involved in methylotrophy were reduced in the course of  
331 genome streamlining with the sediment dwelling *Methylophilus* and *Methylotenera* having the  
332 most complete modules for C<sub>1</sub> compound oxidation, demethylation and assimilation (Fig. 3, 4,  
333 Table S4). They also encode multiple types of methanol dehydrogenases (up to five different  
334 types in single strains), while the pelagic forms possess only XoxF4-1 (Fig. S9). Moreover, the  
335 latter encode neither traditional methylamine-dehydrogenases nor the N-methylglutamate  
336 (NMG) pathway for methylamine oxidation. Thus, the mode of methylamine uptake is still  
337 unclear, although it has been experimentally demonstrated that some pelagic strains can  
338 utilize this C<sub>1</sub> substrate [14, 75]. However, also nearly half of the sediment strains lack these  
339 well-described pathways in a patchy manner only partly reflected by phylogeny, therefore it is  
340 likely that methylamine utilization is not a common feature within Methylophilaceae, or that  
341 alternative routes of its oxidation still await discovery [69]. Formaldehyde oxidation can be  
342 achieved via three alternative routes, and only four *Methylophilus* strains encode all of them,  
343 i.e., all others lack a formaldehyde-dehydrogenase. All *Methylophilus* and *Methylotenera* as  
344 well as 'Ca. M. turicensis' carry genes for the tetrahydromethanopterin (H<sub>4</sub>MPT) pathway, but  
345 none of the 'Ca. Methylopumilus' and OM43 strains. Therefore, the only route for  
346 formaldehyde oxidation in these genome-streamlined microbes is the tetrahydrofolate (H<sub>4</sub>F)  
347 pathway which includes the spontaneous reaction of formaldehyde to H<sub>4</sub>F and is thought to be

348 relatively slow [14, 17]. The ribulose monophosphate (RuMP) cycle for formaldehyde  
349 assimilation/oxidation and formate oxidation via formate dehydrogenases was annotated in all  
350 Methylophilaceae, while none of them possess other potential methylotrophic modules such  
351 as the serine cycle, the ethylmalonyl-CoA-pathway for glyoxylate regeneration, a glyoxylate  
352 shunt, nor the Calvin-Benson-Bassham cycle for CO<sub>2</sub> assimilation, as already previously noted  
353 to be lacking in Methylophilaceae [17]. Thus, the core methylotrophic modules in  
354 Methylophilaceae contain methanol oxidation via XoxF methanol dehydrogenases,  
355 formaldehyde oxidation via the H<sub>4</sub>F pathway, the RuMP cycle, and formate oxidation (Fig. 3,  
356 Table S4)[17, 69, 76]. The majority of genes encoding these pathways were organized in  
357 operon structures or found in close vicinity to each other with high synteny and  
358 phylogenetically reflecting the overall phylogeny of the family (Fig. S10, S11).

359 *Photoheterotrophy as adaptation to oligotrophic pelagic conditions.*

360 Rhodopsins are light-driven proton pumps producing ATP that fuel e.g., membrane  
361 transporters [77] and play important roles during carbon starvation [78] in oligotrophic aquatic  
362 environments. Therefore, the acquisition of rhodopsins are proposed to be powerful  
363 adaptations to the pelagial. ‘Ca. Methylosemipumilus turicensis’ acquired a proteorhodopsin  
364 and the complete pathway for retinal biosynthesis via horizontal gene transfer (HGT) from the  
365 abundant freshwater microbe *Polynucleobacter cosmopolitanus* (81% amino acid similarity,  
366 Fig. 5, Fig. S12a). Interestingly, ‘Ca. Methylopumilus spp.’ carry, in addition to a  
367 proteorhodopsin highly similar to ‘Ca. M. turicensis’ (78.3-79.5% similarity), a second  
368 rhodopsin gene inserted between the proteorhodopsin and the retinal biosynthesis cluster  
369 (Fig. 5). This xantho-like rhodopsin was most likely gained from rare freshwater  
370 Betaproteobacteria (*Janthinobacterium lividum*, *Massilia psychrophila*; 49.4-53.5% similarity,  
371 Fig. S12b), however, the binding of a carotenoid antenna seems unlikely due to the  
372 replacement of a glycine with tryptophan in position 156, suggesting it also functions as a  
373 proteorhodopsin (Fig. S12c)[79, 80]. Both rhodopsins are tuned to green light, which is  
374 common in freshwaters [81] and possess the canonical DxxxK retinal binding motif in helix-7

375 that is characteristic of proton pumping rhodopsins [82]. The marine OM43 lineage only carry  
376 the xantho-like rhodopsin (59.0-63.9% similarity to 'Ca. Methylopusillus'). It is unclear if the  
377 proteorhodopsin was never present in the marine lineage or was lost subsequently, and if so,  
378 the reasons for a secondary loss remain enigmatic as two rhodopsins would provide an even  
379 better adaptation to oligotrophic waters than one.

380 *The second transition from freshwater pelagial to the marine realm is characterized by*  
381 *adaptations to a salty environment*

382 The second habitat transition across the freshwater-marine boundary does not appear  
383 to involve genome streamlining, as genomes of pelagic freshwater and marine methylotrophs  
384 are of similar small size and low GC content (Figs. 1, 2). We hypothesize that this transition  
385 had less impact on the lifestyle (purely planktonic, oligotrophic) but required specific  
386 adaptations to the marine realm that were mainly acquired by HGT, and as suggested by the  
387 long branches in the phylogenetic tree (Fig. 1), multiple, rapid changes in existing genes.  
388 Besides the MCA pathway and the proteorhodopsin, no major rearrangements or reductions  
389 in general metabolic pathways were detected in marine OM43 in comparison to freshwater  
390 'Ca. Methylopusillus'. However, several adaptations to higher salt concentrations could be  
391 identified (Figs. 3, 4). Salinity is one of the most important obstacle in freshwater-marine  
392 colonization, and successful transitions have occurred rarely during the evolution of  
393 Proteobacteria [83, 84]. Main adaptations to higher salinities involve genes for  
394 osmoregulation and inorganic ion metabolism that might have been acquired from the  
395 indigenous community by HGT. For example, genes regulating the Na<sup>+</sup>-dependent  
396 respiratory chain (Na<sup>+</sup>-translocating NADH:quinone oxidoreductase, NQR, Fig. S13) have  
397 been transmitted from the marine *Roseobacter* lineage to strain HTCC2181 [16, 84]. The  
398 NQR system provides energy by generating a sodium motif force, yet, the sodium pumping  
399 might also be an adaptation to enhanced salinities [85]. All other Methylophilaceae possess  
400 the energetically more efficient H<sup>+</sup>-translocating type (NDH), which works better under low  
401 salinity conditions and is thus common in freshwater microbes [85].

402 Ectoine, a compatible solute along with glycine betaine, helps organisms survive  
403 extreme osmotic stress by acting as an osmolyte [86]. Ectoine is synthesized from L-aspartate  
404 4-semialdehyde, the central intermediate in the synthesis of amino acids of the aspartate  
405 family. Two marine OM43 strains (KB13 and MBRS-H7) encode this pathway followed by  
406 sodium:proline symporter *putP* arranged in high synteny and protein similarity with  
407 marine/hypersaline sediment microbes, thus it is likely that both components were gained via  
408 HGT (Fig. S14). A second copy of the *putP* symporter was common to all Methylophilaceae  
409 (data not shown). Also a dipeptide/tripeptide permease (DtpD) unique for the marine OM43  
410 lineage seems to be transferred horizontally, either from marine Bacteroidetes or sediment-  
411 dwelling *Sulfurifustis* (Gammaproteobacteria, Fig. S15). Other putative membrane compounds  
412 involved in sodium transport in marine OM43 include a sodium:alanine symporter (AlsT, Fig.  
413 S16a), a sodium:acetate symporter (ActP, Fig. S16b), a sodium:dicarboxylate symporter (GltT,  
414 Fig. S16c), a sodium:proton antiporter (NhaP, Fig. S16d), and another putative sodium:proton  
415 antiporter (NhaE-like, Fig. S16e). Although also several other Methylophilaceae carry some of  
416 these sodium transporters, they are only distantly related to OM43, thus they might be  
417 acquired horizontally. Conversely, ActP and GltT of OM43 are most closely related to three  
418 ‘*Ca. M. universalis*’ strains and the two ‘*Ca. M. rimovensis*’ strains, respectively (Fig. S16b,  
419 S16c). Both symporters are related to microbes from freshwater and marine habitats, hinting  
420 to some yet unknown lineages related to both OM43 and ‘*Ca. Methylopumilus*’ most likely  
421 thriving in the freshwater-marine transition zone.

## 422 *Conclusions*

423 Our study provides first genomic evidence that the ancestors of genome-streamlined  
424 pelagic Methylophilaceae can be traced back to sediments with two habitat transitions  
425 occurring in the evolutionary history of the family. The first from sediments to the pelagial is  
426 characterized by pronounced genome reduction driven by selection pressure for relatively  
427 more oligotrophic environmental conditions. This adaptive gene loss has mainly affected  
428 functions that (i) are not necessarily required in the pelagial (e.g., motility, chemotaxis), (ii) are



429 not advantageous for survival in an oligotrophic habitat (e.g., low substrate affinity  
430 transporters), and (iii) are encoded in redundant pathways (e.g., formaldehyde oxidation).  
431 Likewise, (iv) genes providing adaptations to oligotrophic conditions have been transmitted  
432 horizontally from indigenous pelagic microbes (e.g., rhodopsins). The second habitat transition  
433 across the freshwater-marine boundary did not result in further genome-streamlining, but is  
434 characterized by adaptations to higher salinities acquired by HGT. 'Ca. M. turicensis' was  
435 identified as transitional taxon, retaining multiple ancestral characters while also gaining  
436 adaptations to the pelagial. In this regard, the family Methylophilaceae is an exceptional model  
437 for tracing the evolutionary history of genome-streamlining as such a collection of  
438 evolutionarily related microbes from different habitats is practically unknown for other similarly  
439 abundant genome-streamlined microbes (e.g., 'Ca. Pelagibacterales', 'Ca. Nanopelagicales').

440

#### 441 **Acknowledgements**

442 We thank the team of the Genetic Diversity Center Zurich (GDC) for providing sequencing  
443 facilities and help with library preparation. Thomas Posch and Eugen Loher are acknowledged  
444 for help in sampling of Lake Zurich, Petr Znachor, Pavel Rychtecký and Jiří Nedoma for help  
445 in sampling of Řimov Reservoir and Lake Medard. MMS was supported by the research grant  
446 19-23469S (Grant Agency of the Czech Republic). RG was supported by the research grant  
447 17-04828S (Grant Agency of the Czech Republic). Sampling for the isolation of novel strains  
448 from Lake Zurich was supported by the SNF D-A-CH project 310030E-160603/1 awarded to  
449 Thomas Posch.

#### 450 **Authors' contributions**

451 MMS conceived the project, isolated and sequenced the strains, analysed the data and wrote  
452 the manuscript. DS, MK and SMN sequenced the strains and contributed to data analyses.  
453 RG developed programs for analysis and contributed to data analyses. All authors helped to  
454 interpret the results and contributed to writing the manuscript.

455

456 **Competing interests**

457 The authors declare that they have no competing interests.

458

459 **References**

- 460 1. Button DK. Biochemical basis for whole-cell uptake kinetics: specific affinity, oligotrophic  
461 capacity, and the meaning of the Michaelis constant. *Appl Environ Microbiol* 1991; **57**: 2033-  
462 2038.
- 463 2. Pernthaler J. Predation on prokaryotes in the water column and its ecological implications.  
464 *Nat Rev Microbiol* 2005; **3**: 537-546.
- 465 3. Giovannoni SJ, Cameron Thrash J, Temperton B. Implications of streamlining theory for  
466 microbial ecology. *ISME J* 2014; **8**: 1553–1565.
- 467 4. Luo H, Thompson LR, Stingl U, Hughes AL. Selection maintains low genomic GC content in  
468 marine SAR11 lineages. *Mol Biol Evol* 2015; **32**: 2738-2748.
- 469 5. Giovannoni SJ, Tripp HJ, Givan S, Podar M, Vergin KL, Baptista D *et al.* Genome streamlining  
470 in a cosmopolitan oceanic bacterium. *Science* 2005; **309**: 1242-1245.
- 471 6. Salcher MM, Pernthaler J, Posch T. Seasonal bloom dynamics and ecophysiology of the  
472 freshwater sister clade of SAR11 bacteria 'that rule the waves' (LD12). *ISME J* 2011; **5**: 1242-  
473 1252.
- 474 7. Eiler A, Mondav R, Sinclair L, Fernandez-Vidal L, Scofield DG, Schwientek P *et al.* Tuning fresh:  
475 radiation through rewiring of central metabolism in streamlined bacteria. *ISME J* 2016; **10**:  
476 1902-1914.
- 477 8. Ghai R, Mizuno CM, Picazo A, Camacho A, Rodriguez-Valera F. Metagenomics uncovers a new  
478 group of low GC and ultra-small marine Actinobacteria. *Sci Rep* 2013; **3**.
- 479 9. Newton RJ, Jones SE, Helmus MR, McMahon KD. Phylogenetic ecology of the freshwater  
480 *Actinobacteria* acl lineage. *Appl Environ Microbiol* 2007; **73**: 7169-7176.
- 481 10. Neuenschwander SM, Ghai R, Pernthaler J, Salcher MM. Microdiversification in genome-  
482 streamlined ubiquitous freshwater Actinobacteria. *ISME J* 2018; **12**: 185.
- 483 11. Hahn MW, Schmidt J, Taipale SJ, Doolittle WF, Koll U. *Rhodoluna ladicola* gen. nov., sp. nov., a  
484 planktonic freshwater bacterium with stream-lined genome. *Int J Syst Evol Microbiol* 2014;  
485 **64**: 3254-3263.
- 486 12. Kang I, Lee K, Yang S-J, Choi A, Kang D, Lee YK *et al.* Genome sequence of “*Candidatus*  
487 *Aquiluna*” sp. strain IMCC13023, a marine member of the Actinobacteria isolated from an  
488 Arctic fjord. *J Bacteriol* 2012; **194**: 3550-3551.
- 489 13. Santoro AE, Dupont CL, Richter RA, Craig MT, Carini P, McIlvin MR *et al.* Genomic and  
490 proteomic characterization of “*Candidatus Nitrosopelagicus brevis*”: An ammonia-oxidizing  
491 archaeon from the open ocean. *Proc Natl Acad Sci USA* 2015; **112**: 1173-1178.
- 492 14. Salcher MM, Neuenschwander SM, Posch T, Pernthaler J. The ecology of pelagic freshwater  
493 methylotrophs assessed by a high-resolution monitoring and isolation campaign. *ISME J*  
494 2015; **9**: 2442-2453.
- 495 15. Giovannoni SJ, Hayakawa DH, Tripp HJ, Stingl U, Givan SA, Cho JC *et al.* The small genome of  
496 an abundant coastal ocean methylotroph. *Environ Microbiol* 2008; **10**: 1771-1782.
- 497 16. Jimenez-Infante F, Ngugi DK, Vinu M, Alam I, Kamau AA, Blom J *et al.* Comprehensive  
498 genomic analyses of the OM43 clade, including a novel species from the Red Sea, indicate  
499 ecotype differentiation among marine methylotrophs. *Appl Environ Microbiol* 2016; **82**:  
500 1215-1226.
- 501 17. Chistoserdova L. Modularity of methylotrophy, revisited. *Environ Microbiol* 2011; **13**: 2603-  
502 2622.

- 503 18. Chistoserdova L. Methyloprophs in natural habitats: current insights through metagenomics. *Appl Microbiol Biotechnol* 2015; **99**: 5763-5779.
- 504
- 505 19. Kalyuzhnaya MG, Bowerman S, Lara JC, Lidstrom ME, Chistoserdova L. *Methylotenera mobilis*  
506 gen. nov., sp. nov., an obligately methylamine-utilizing bacterium within the family  
507 *Methylophilaceae*. *Int J Syst Evol Microbiol* 2006; **56**: 2819-2823.
- 508 20. Govorukhina NI, Trotsenko YA. Methylovorus, a new genus of restricted facultatively  
509 methyloprophic bacteria. *Int J Syst Bacteriol* 1991; **41**: 158-162.
- 510 21. Yordy JR, Weaver TL. Methylobacillus: a new genus of obligately methyloprophic bacteria. *Int*  
511 *J Syst Bacteriol* 1977; **27**: 247-255.
- 512 22. Jenkins O, Byrom D, Jones D. *Methylophilus* - a new genus of methanol-utilizing bacteria. *Int J*  
513 *Syst Bacteriol* 1987; **37**: 446-448.
- 514 23. Huggett M, Hayakawa D, Rappe M. Genome sequence of strain HIMB624, a cultured  
515 representative from the OM43 clade of marine Betaproteobacteria. *Standards in Genomic*  
516 *Sciences* 2012; **6**.
- 517 24. Newton RJ, Jones SE, Eiler A, McMahon KD, Bertilsson S. A guide to the natural history of  
518 freshwater lake bacteria. *Microbiol Mol Biol R* 2011; **75**: 14-49.
- 519 25. Woodhouse JN, Kinsela AS, Collins RN, Bowling LC, Honeyman GL, Holliday JK *et al*. Microbial  
520 communities reflect temporal changes in cyanobacterial composition in a shallow ephemeral  
521 freshwater lake. *ISME J* 2016; **10**: 1337-1351.
- 522 26. Li J, Zhang J, Liu L, Fan Y, Li L, Yang Y *et al*. Annual periodicity in planktonic bacterial and  
523 archaeal community composition of eutrophic Lake Taihu. *Sci Rep* 2015; **5**: 15488.
- 524 27. Ramachandran A, Walsh DA. Investigation of XoxF methanol dehydrogenases reveals new  
525 methyloprophic bacteria in pelagic marine and freshwater ecosystems. *FEMS Microbiol Ecol*  
526 2015; **91**: fiv105.
- 527 28. Morris RM, Longnecker K, Giovannoni SJ. Pirellula and OM43 are among the dominant  
528 lineages identified in an Oregon coast diatom bloom. *Environ Microbiol* 2006; **8**: 1361-1370.
- 529 29. Sekar R, Fuchs BM, Amann R, Pernthaler J. Flow sorting of marine bacterioplankton after  
530 fluorescence in situ hybridization. *Appl Environ Microbiol* 2004; **70**: 6210-6219.
- 531 30. Salcher MM, Šimek K. Isolation and cultivation of planktonic freshwater microbes is essential  
532 for a comprehensive understanding of their ecology. *Aquat Microb Ecol* 2016; **77**: 183-196.
- 533 31. Zotina T, Köster O, Jüttner F. Photoheterotrophy and light-dependent uptake of organic and  
534 organic nitrogenous compounds by *Planktothrix rubescens* under low irradiance. *Freshwater*  
535 *Biol* 2003; **48**: 1859-1872.
- 536 32. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data.  
537 *Bioinformatics* 2014; **30**: 2114-2120.
- 538 33. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS *et al*. SPAdes: A new  
539 genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol*  
540 2012; **19**: 455-477.
- 541 34. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 2014; **30**: 2068-  
542 2069.
- 543 35. Tatusov RL, Natale DA, Garkavtsev IV, Tatusova TA, Shankavaram UT, Rao BS *et al*. The COG  
544 database: new developments in phylogenetic classification of proteins from complete  
545 genomes. *Nucl Acid Res* 2001; **29**: 22-28.
- 546 36. Haft DH, Loftus BJ, Richardson DL, Yang F, Eisen JA, Paulsen Ian T *et al*. TIGRFAMs: a protein  
547 family resource for the functional identification of proteins. *Nucl Acid Res* 2001; **29**: 41-43.
- 548 37. Kanehisa M, Sato Y, Morishima K. BlastKOALA and GhostKOALA: KEGG tools for functional  
549 characterization of genome and metagenome sequences. *J Mol Biol* 2016; **428**: 726-731.
- 550 38. Caspi R, Billington R, Fulcher CA, Keseler IM, Kothari A, Krummenacker M *et al*. The MetaCyc  
551 database of metabolic pathways and enzymes. *Nucl Acid Res* 2017; **46**: D633-D639.
- 552 39. Chistoserdova L, Lapidus A, Han C, Goodwin L, Saunders L, Brettin T *et al*. Genome of  
553 *Methylobacillus flagellatus*, molecular basis for obligate methyloprophic, and polyphyletic  
554 origin of methyloprophic. *J Bacteriol* 2007; **189**: 4020-4027.

- 555 40. Lapidus A, Clum A, LaButti K, Kaluzhnaya MG, Lim S, Beck DAC *et al.* Genomes of three  
556 methylotrophs from a single niche reveal the genetic and metabolic divergence of the  
557 *Methylophilaceae*. *J Bacteriol* 2011; **193**: 3757-3764.
- 558 41. Vuilleumier S, Chistoserdova L, Lee M-C, Bringel F, Lajus A, Zhou Y *et al.* Methylobacterium  
559 Genome Sequences: A Reference Blueprint to Investigate Microbial Metabolism of C1  
560 Compounds from Natural and Industrial Sources. *PLOS ONE* 2009; **4**: e5584.
- 561 42. Good N, Lamb A, Beck D, Martinez-Gomez N, Kalyuzhnaya M. C1-Pathways in  
562 *Methyloversatilis universalis* FAM5: Genome Wide Gene Expression and Mutagenesis  
563 Studies. *Microorganisms* 2015; **3**: 175.
- 564 43. Kalyuzhnaya MG, Hristova KR, Lidstrom ME, Chistoserdova L. Characterization of a Novel  
565 Methanol Dehydrogenase in Representatives of Burkholderiales: Implications for  
566 Environmental Detection of Methylotrophy and Evidence for Convergent Evolution. *J*  
567 *Bacteriol* 2008; **190**: 3817-3823.
- 568 44. Brautaset T, Jakobsen ØM, Flickinger MC, Valla S, Ellingsen TE. Plasmid-Dependent  
569 Methylotrophy in Thermotolerant *Bacillus methanolicus*. *J Bacteriol* 2004; **186**: 1229-1238.
- 570 45. Vorholt JA, Kalyuzhnaya MG, Hagemeyer CH, Lidstrom ME, Chistoserdova L. MtdC, a Novel  
571 Class of Methylene Tetrahydromethanopterin Dehydrogenases. *J Bacteriol* 2005; **187**: 6069-  
572 6074.
- 573 46. Ward N, Larsen Ø, Sakwa J, Bruseth L, Khouri H, Durkin AS *et al.* Genomic Insights into  
574 Methanotrophy: The Complete Genome Sequence of *Methylococcus capsulatus* (Bath). *PLoS*  
575 *Biol* 2004; **2**: e303.
- 576 47. Hou S, Makarova KS, Saw JH, Senin P, Ly BV, Zhou Z *et al.* Complete genome sequence of the  
577 extremely acidophilic methanotroph isolate V4, *Methylacidiphilum infernorum*, a  
578 representative of the bacterial phylum Verrucomicrobia. *Biology Direct* 2008; **3**: 26.
- 579 48. Wu ML, Wessels HJCT, Pol A, Op den Camp HJM, Jetten MSM, van Niftrik L *et al.* XoxF-type  
580 methanol dehydrogenase from the anaerobic methanotroph "*Candidatus Methyloimrabilis*  
581 *oxyfera*". *Appl Environ Microbiol* 2015; **81**: 1442-1451.
- 582 49. Sun J, Steindler L, Thrash JC, Halsey KH, Smith DP, Carter AE *et al.* One carbon metabolism in  
583 SAR11 pelagic marine bacteria. *PLoS ONE* 2011; **6**: e23973.
- 584 50. Denef VJ, Mueller RS, Chiang E, Liebig JR, Vanderploeg HA. Chloroflexi CL500-11 populations  
585 that predominate deep-lake hypolimnion bacterioplankton rely on nitrogen-rich dissolved  
586 organic matter metabolism and C1 compound oxidation. *Appl Environ Microbiol* 2016; **82**:  
587 1423-1432.
- 588 51. Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein or  
589 nucleotide sequences. *Bioinformatics* 2006; **22**: 1658-1659.
- 590 52. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
591 *Nucl Acid Res* 2004; **32**: 1792-1797.
- 592 53. Eddy SR. Accelerated Profile HMM Searches. *PLOS Computational Biology* 2011; **7**: e1002195.
- 593 54. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W *et al.* Gapped BLAST and PSI-  
594 BLAST: a new generation of protein database search programs. *Nucl Acid Res* 1997; **25**: 3389-  
595 3402.
- 596 55. Goris J, Konstantinidis KT, Klappenbach JA, Coenye T, Vandamme P, Tiedje JM. DNA-DNA  
597 hybridization values and their relationship to whole-genome sequence similarities. *Int J Syst*  
598 *Evol Microbiol* 2007; **57**: 81-91.
- 599 56. Rodriguez-R LM, Konstantinidis KT. Bypassing cultivation to identify bacterial species. *ASM*  
600 *Microbe Magazine* 2014; **9**: 111-118.
- 601 57. Price MN, Dehal PS, Arkin AP. FastTree 2 – Approximately maximum-likelihood trees for large  
602 alignments. *PLOS ONE* 2010; **5**: e9490.
- 603 58. Lassmann T, Sonnhammer EL. Kalign – an accurate and fast multiple sequence alignment  
604 algorithm. *BMC Bioinformatics* 2005; **6**: 298.
- 605 59. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
606 phylogenies. *Bioinformatics* 2014; **30**: 1312-1313.

- 607 60. Katoh K, Standley DM. MAFFT Multiple Sequence Alignment Software Version 7:  
608 Improvements in Performance and Usability. *Mol Biol Evol* 2013; **30**: 772-780.
- 609 61. Konstantinidis KT, Rossello-Mora R, Amann R. Uncultivated microbes in need of their own  
610 taxonomy. *ISME J* 2017; **11**: 2399-2406.
- 611 62. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarszewski A, Chaumeil P-A *et al.* A  
612 standardized bacterial taxonomy based on genome phylogeny substantially revises the tree  
613 of life. *Nature Biotechnology* 2018.
- 614 63. Salcher MM, Pernthaler J, Frater N, Posch T. Vertical and longitudinal distribution patterns of  
615 different bacterioplankton populations in a canyon-shaped, deep prealpine lake. *Limnol  
616 Oceanogr* 2011; **56**: 2027-2039.
- 617 64. Bock C, Salcher M, Jensen M, Pandey RV, Boenigk J. Synchrony of eukaryotic and prokaryotic  
618 planktonic communities in three seasonally sampled Austrian lakes. *Front Microbiol* 2018; **9**.
- 619 65. Linz AM, Cray BC, Shade A, Owens S, Gilbert JA, Knight R *et al.* Bacterial community  
620 composition and dynamics spanning five years in freshwater bog lakes. *mSphere* 2017; **2**.
- 621 66. Okazaki Y, Nakano S-i. Vertical partitioning of freshwater bacterioplankton community in a  
622 deep mesotrophic lake with a fully oxygenated hypolimnion (Lake Biwa, Japan). *Environ  
623 Microbiol Rep* 2016; **8**: 780-788.
- 624 67. Cabello-Yeves PJ, Zemskaya TI, Rosselli R, Coutinho FH, Zakharenko AS, Blinov VV *et al.*  
625 Genomes of novel microbial lineages assembled from the sub-ice waters of Lake Baikal. *Appl  
626 Environ Microbiol* 2018; **84**: e02132-02117.
- 627 68. Bendall ML, Stevens SLR, Chan L-K, Malfatti S, Schwientek P, Tremblay J *et al.* Genome-wide  
628 selective sweeps and gene-specific sweeps in natural bacterial populations. *ISME J* 2016; **10**:  
629 1589-1601.
- 630 69. Beck DAC, McTaggart TL, Setboonsarng U, Vorobev A, Kalyuzhnaya MG, Ivanova N *et al.* The  
631 expanded diversity of *Methylophilaceae* from Lake Washington through cultivation and  
632 genomic sequencing of novel ecotypes. *PLoS ONE* 2014; **9**: e102458.
- 633 70. Hendriks J, Oubrie A, Castresana J, Urbani A, Gemeinhardt S, Saraste M. Nitric oxide  
634 reductases in bacteria. *Biochimica et Biophysica Acta (BBA) - Bioenergetics* 2000; **1459**: 266-  
635 273.
- 636 71. Mustakhimov I, Kalyuzhnaya MG, Lidstrom ME, Chistoserdova L. Insights into denitrification  
637 in *Methylostenella mobilis* from denitrification pathway and methanol metabolism mutants. *J  
638 Bacteriol* 2013; **195**: 2207-2211.
- 639 72. Kalyuzhnaya MG, Lapidus A, Ivanova N, Copeland AC, McHardy AC, Szeto E *et al.* High-  
640 resolution metagenomics targets specific functional types in complex microbial communities.  
641 *Nature Biotechnology* 2008; **26**: 1029-1034.
- 642 73. Kitzinger K, Padilla CC, Marchant HK, Hach PF, Herbold CW, Kidane AT *et al.* Cyanate and urea  
643 are substrates for nitrification by Thaumarchaeota in the marine environment. *Nat Microbiol*  
644 2019; **4**: 234-243.
- 645 74. Wetzel R: *Limnology. Lake and River Ecosystems*. 3<sup>rd</sup> edn: Elsevier Academic Press; 2001.
- 646 75. Halsey KH, Carter AE, Giovannoni SJ. Synergistic metabolism of a broad range of C<sub>1</sub>  
647 compounds in the marine methylotrophic bacterium HTCC2181. *Environ Microbiol* 2011; **14**:  
648 630-640.
- 649 76. Chistoserdova L. Methylotrophy in a lake: from metagenomics to single-organism physiology.  
650 *Appl Environ Microbiol* 2011; **77**: 4705-4711.
- 651 77. Pinhassi J, DeLong EF, Béjà O, González JM, Pedrós-Alió C. Marine bacterial and archaeal ion-  
652 pumping rhodopsins: Genetic diversity, physiology, and ecology. *Microbiol Mol Biol R* 2016;  
653 **80**: 929-954.
- 654 78. Steindler L, Schwalbach MS, Smith DP, Chan F, Giovannoni SJ. Energy starved *Candidatus*  
655 Pelagibacter ubique substitutes light-mediated ATP production for endogenous carbon  
656 respiration. *PLoS ONE* 2011; **6**: e19725.

- 657 79. Luecke H, Schobert B, Stagno J, Imasheva ES, Wang JM, Balashov SP *et al.* Crystallographic  
658 structure of xanthorhodopsin, the light-driven proton pump with a dual chromophore.  
659 *Proceedings of the National Academy of Sciences* 2008; **105**: 16561-16565.
- 660 80. Imasheva ES, Balashov SP, Choi AR, Jung K-H, Lanyi JK. Reconstitution of Gloeobacter  
661 violaceus Rhodopsin with a Light-Harvesting Carotenoid Antenna. *Biochemistry* 2009; **48**:  
662 10948-10955.
- 663 81. Rusch DB, Halpern AL, Sutton G, Heidelberg KB, Williamson S, Yooseph S *et al.* The *Sorcerer II*  
664 Global Ocean Sampling Expedition: Northwest Atlantic through Eastern Tropical Pacific. *PLoS*  
665 *Biol* 2007; **5**: e77.
- 666 82. Bèjà O, Aravind L, Koonin EV, Suzuki MT, Hadd A, Nguyen LP *et al.* Bacterial Rhodopsin:  
667 Evidence for a New Type of Phototrophy in the Sea. *Science* 2000; **289**: 1902-1906.
- 668 83. Logares R, Bråte J, Bertilsson S, Clasen JL, Shalchian-Tabrizi K, Rengefors K. Infrequent  
669 marine-freshwater transitions in the microbial world. *Trends Microbiol* 2009; **17**: 414-422.
- 670 84. Walsh DA, Lafontaine J, Grossart H-P: On the eco-evolutionary relationships of fresh and salt  
671 water bacteria and the role of gene transfer in their adaptation. In *Lateral Gene Transfer in*  
672 *Evolution*. Edited by Gophna U: Springer New York; 2013: 55-77
- 673 85. Zhang H, Yoshizawa S, Sun Y, Huang Y, Chu X, González JM *et al.* Repeated evolutionary  
674 transitions of flavobacteria from marine to non-marine habitats. *Environ Microbiol* 2019; **0**.
- 675 86. Roberts MF. Organic compatible solutes of halotolerant and halophilic microorganisms.  
676 *Saline Systems* 2005; **1**: 5.

677

678

679 **Figure legends:**

680 **Figure 1: Phylogeny of Methylophilaceae and their occurrence in different**  
681 **environments.** (a) Phylogenomic tree based on 878 common concatenated proteins  
682 (351,312 amino acid sites) with *Methyloversatilis* sp. RAC08 as outgroup. The 39 complete  
683 genomes of 'Ca. Methylopumilus sp.' are collapsed to the species level, see Fig. S1 for a  
684 complete tree. Different genera (70% AAI cut-off, Fig. S2) are marked by grey boxes.  
685 Isolation sources of strains are indicated by different colours and incomplete genomes  
686 consisting of <17 contigs (estimated completeness >99%) are marked with asterisks.  
687 Bootstrap values (100 repetitions) are indicated at the nodes, the scale bar at the bottom  
688 indicates 20% sequence divergence. The genome sizes for all strains are shown with circles  
689 of proportional size and GC content is depicted within each circle. (b) Fragment recruitment  
690 of public metagenomes from freshwater sediments (n=131), lake pelagial (n=345), rivers  
691 (n=43), estuaries and coastal oceans (n=53), and open oceans (n=201). Maximum RPKG  
692 values (number of reads recruited per kb of genome per Gb of metagenome) for each  
693 ecosystem are shown for each genome.

694 **Figure 2: Genome streamlining in Methylophilaceae.** Significant relationships between  
695 genome sizes (Mbp) and genomic GC content (%), lengths of intergenic spacers (bp), coding  
696 density (%), mean CDS length (bp), overlapping CDS (%), number of paralogs, histidine  
697 kinases and sigma factors and significant relationships between genomic GC content (%)  
698 and TAA stop codon, TAG stop codon, lysine, and arginine usage (%) and C, N, and S  
699 atoms per amino acid (mol%).

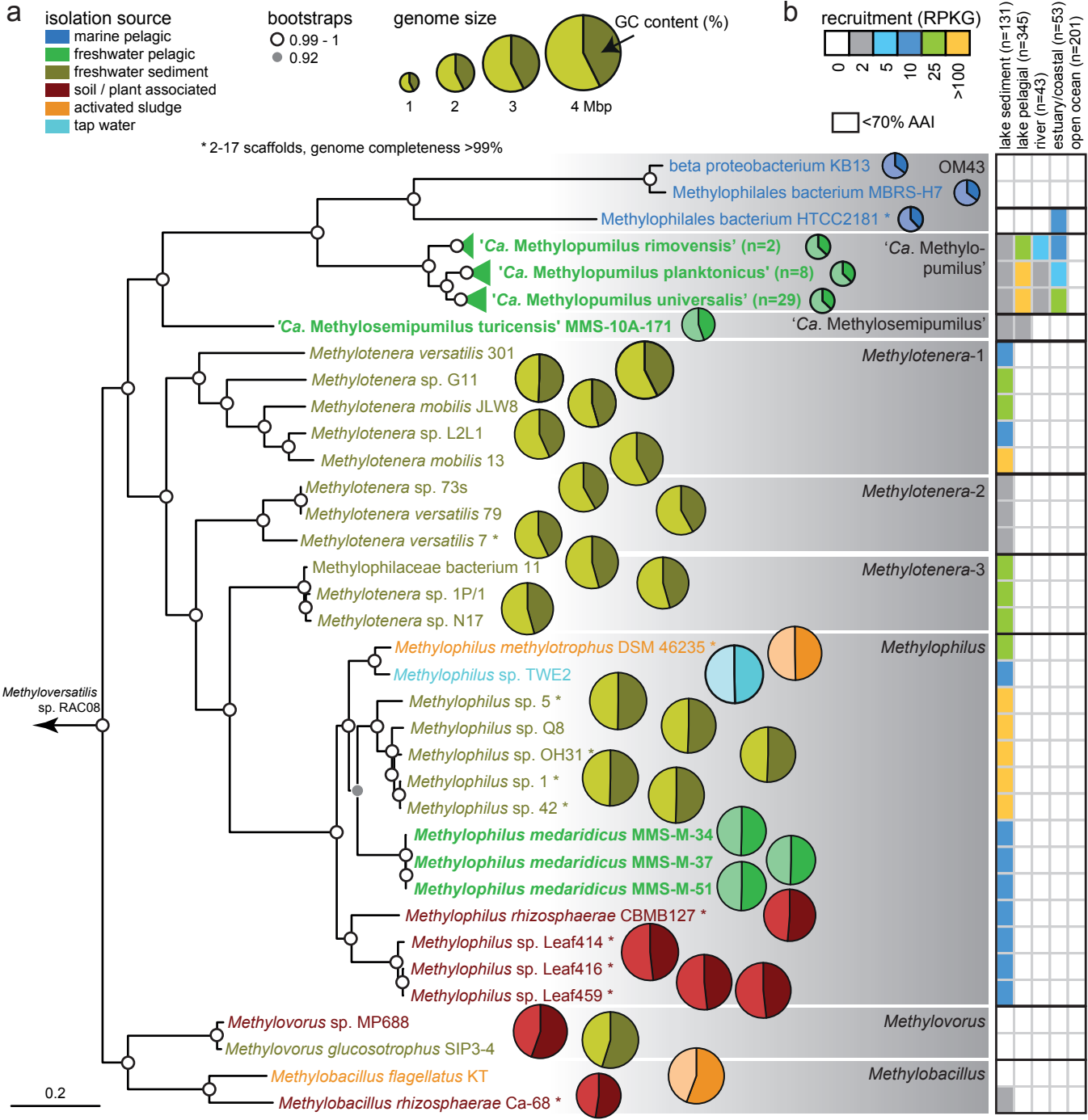
700 **Fig. 3: Metabolic modules in Methylophilaceae.** Presence of selected metabolic modules  
701 in Methylophilaceae strains. For details on phylogenomic tree see Fig. 1 and Fig. S1, for  
702 details on pathways see Table S4. T4SS: type 4 secretion system; MOX-PQQ: methanol  
703 oxidation and pyrroloquinoline quinone biosynthesis; XoxF/MxaF: methanol dehydrogenase  
704 XoxF/MxaF; NMG: N-methylglutamate pathway; MADH: methylamine dehydrogenase;  
705 FaDH: formaldehyde dehydrogenase; H4F-foID: H<sub>4</sub>-linked formaldehyde oxidation, foID form;

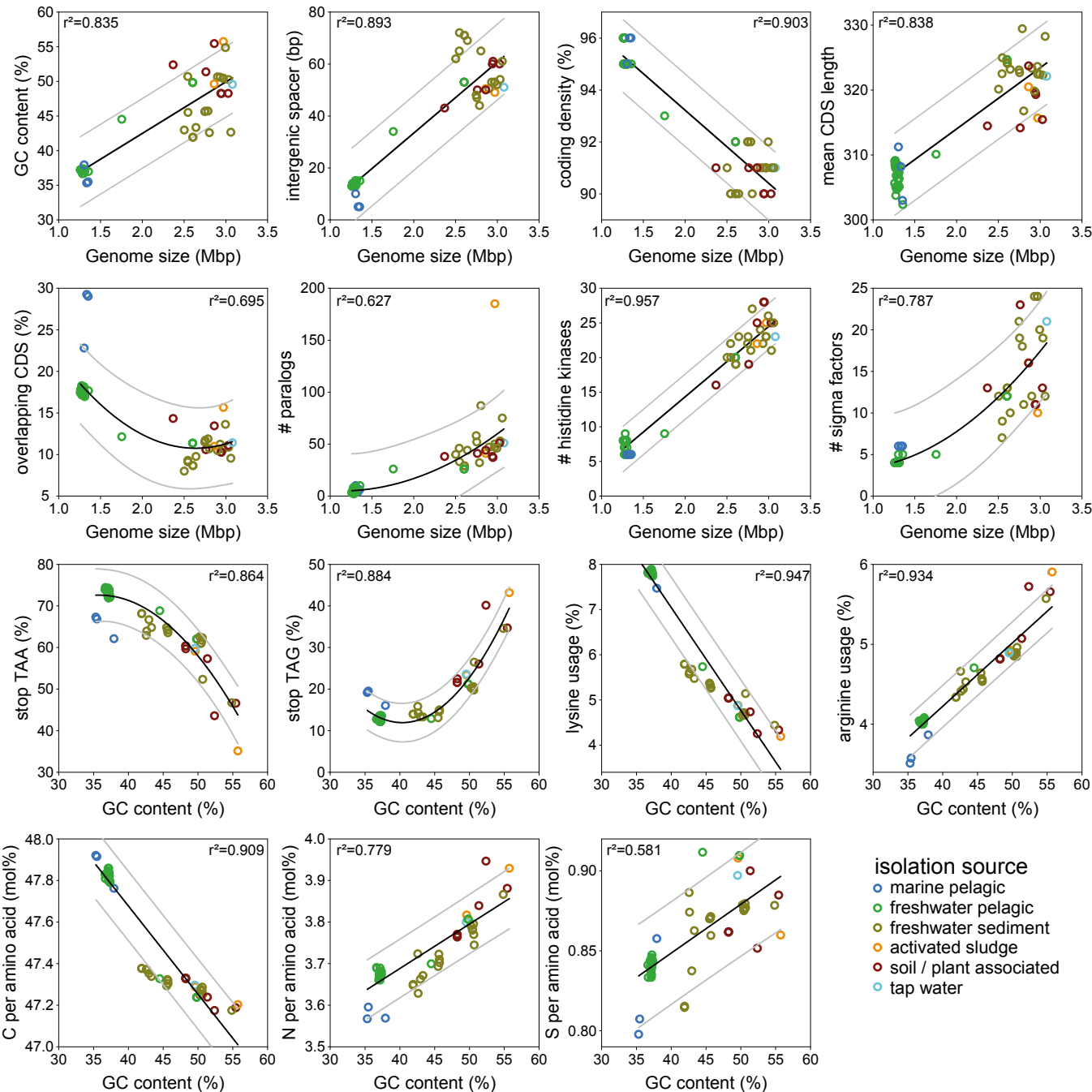
706 H4MPT: H<sub>4</sub>MPT-linked formaldehyde oxidation; FOX: formate oxidation; RuMP: ribulose  
707 monophosphate cycle; ass. NO<sub>3</sub> reduction: assimilatory nitrate reduction; diss. NO<sub>3</sub>  
708 reduction: dissimilatory nitrate reduction; ass. SO<sub>4</sub> reduction: assimilatory sulfate reduction;  
709 NDH-dehydrogenase: H<sup>+</sup>-translocating NADH:quinone oxidoreductase; NQR-  
710 dehydrogenase: Na<sup>+</sup>-translocating NADH:quinone oxidoreductase; DtpD: dipeptide/tripeptide  
711 permease DtpD; AlsT: sodium:alanine symporter AlsT; ActP: sodium:acetate symporter  
712 ActP; GltT: sodium:dicarboxylate symporter GltT; NhaP2: sodium:proton antiporter NhaP2;  
713 NhaE-like: sodium:proton antiporter NhaE; AA transporter: amino acid transporters; fructose  
714 PTS: fructose-specific phosphotransferase system; MCA: methylcitric acid cycle.

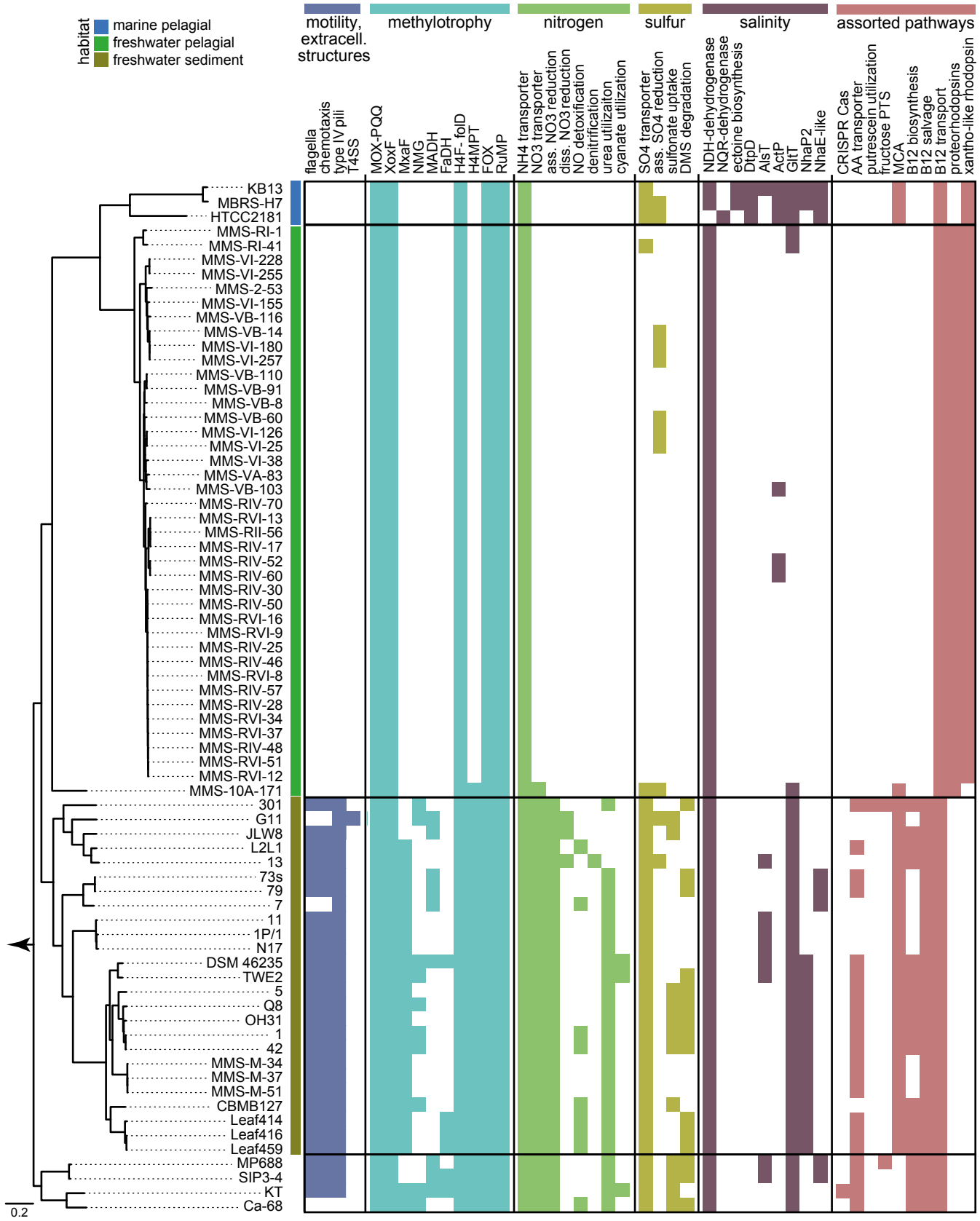
715 **Figure 4: Comparative metabolic maps of different taxa of Methylophilaceae.** (a)  
716 Comparison of the core metabolism in sediment Methylophilaceae (*Methylotenera* and  
717 *Methylophilus*) vs. 'Ca. Methylosemipumilus turicensis' (Mtur). (b) Comparison of the core  
718 metabolism in 'Ca. Methylosemipumilus turicensis' (Mtur) vs. 'Ca. Methylopumilus' (Mpum)  
719 vs. marine OM43. For details on pathways see Table S4.

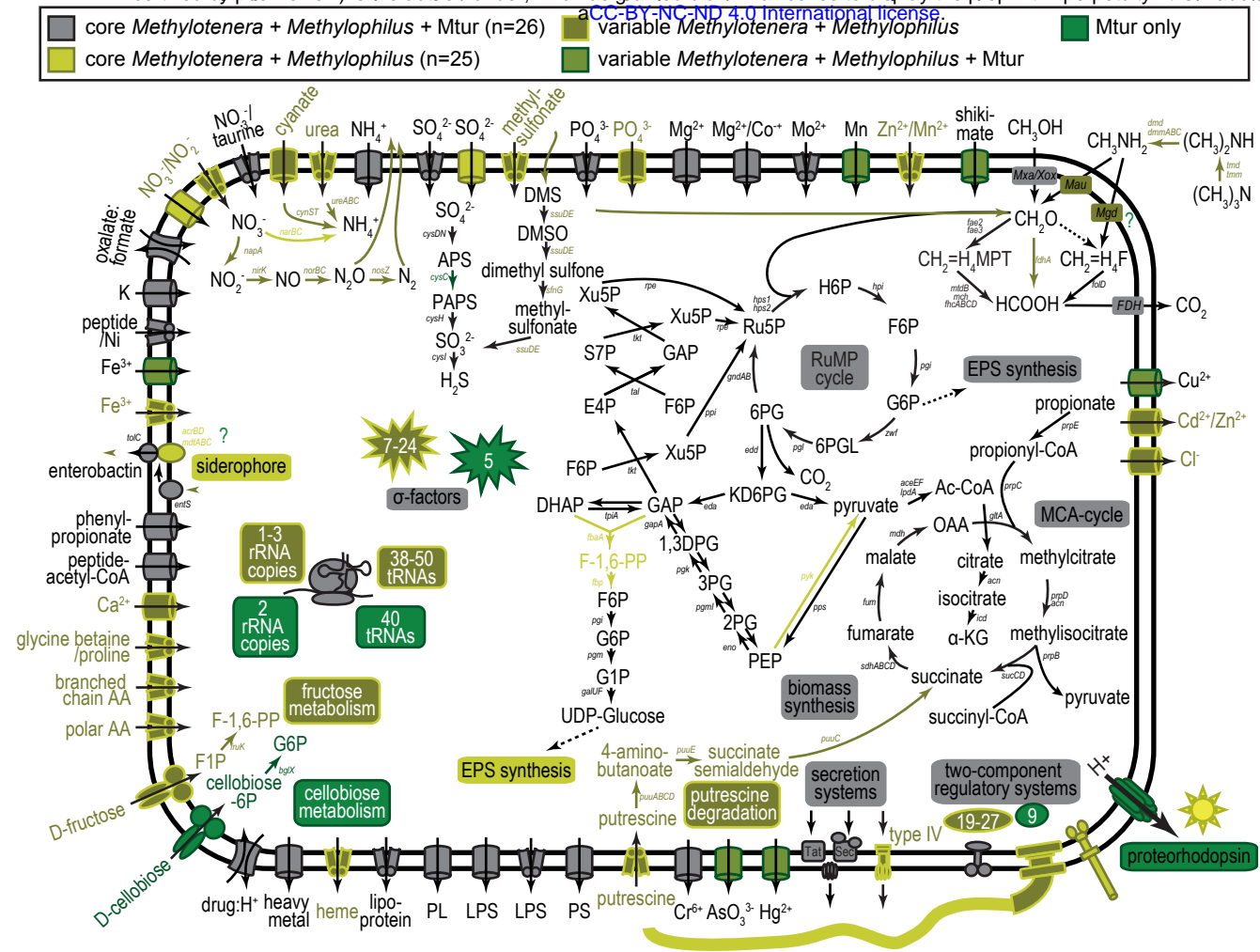
720 **Figure 5: Horizontal gene transfers of two different rhodopsins.** (a) Phylogenetic tree  
721 (RAxML, 100 bootstraps) of different rhodopsin types. See Fig. S12 for details of closely  
722 related proteo- and xantho-like rhodopsins of Methylophilaceae. (b) Arrangement and protein  
723 similarity of rhodopsins and the retinal biosynthesis gene cluster in 'Ca. Methylopumilus  
724 spp.', 'Ca. Methylosemipumilus turicensis', marine OM43 and other freshwater microbes with  
725 closely related rhodopsin types.











**b** 'Ca. Methylosemipumilus turicensis' (Mtur) vs. 'Ca. Methylophilus' (Mpum) vs. marine OM43

