

1 **Supplemental Material**

2 **Genomic evidence for simultaneous optimization of**  
3 **transcription and translation through codon variants in the**  
4 ***pmoCAB* operon of type Ia methanotrophs**

5 Juan C. Villada, Maria F. Duran, and Patrick K. H. Lee<sup>#</sup>

6 School of Energy and Environment, City University of Hong Kong, Kowloon, Hong  
7 Kong SAR, China

8 **Correspondence:** #B5423, Yeung Kin Man Academic Building, School of Energy  
9 and Environment, City University of Hong Kong, Tat Chee Avenue, Kowloon, Hong  
10 Kong SAR, China; E-mail: patrick.kh.lee@cityu.edu.hk; Tel: (852) 3442-4625; Fax:  
11 (852) 3442-0688.

## 12 Supplemental Methods

13 **Computational packages and scripts.** *R* [1], *RStudio* [2], and *ggplot2* [3] were  
14 used to produce all the analyses and figures presented in this study unless otherwise  
15 indicated. All the scripts used in this work are available at GitHub  
16 ([https://github.com/PLeeLab/methane\\_oxidation\\_genetic\\_trait](https://github.com/PLeeLab/methane_oxidation_genetic_trait)).

17 **Metagenomics analysis pipeline.** Our analysis was applied to publicly available  
18 metagenomic data from five potentially methanotrophic environments: 1) Lake  
19 Washington, USA [4]; 2) Serpentinite Springs of the Voltri Massif, Italy [5]; 3) Movile  
20 Cave in Mangalia, Romania [6]; 4) Santa Elena Ophiolite alkaline spring, Costa Rica  
21 [7]; and 5) Coastal basin of Golfo Dulce, Costa Rica [8]. For 1), the publicly available  
22 metagenome-assemble genomes (MAGs)  
23 (<https://gold.jgi.doe.gov/studies?id=Gso114290>) were also examined.

24 The metagenomics analysis pipeline consists of seven main stages and a  
25 preliminary Stage 0 for data and software preparation. In Stage 0, raw metagenomic  
26 reads in FASTQ format were downloaded from NCBI/SRA using *fastq-dump* from the  
27 SRA Toolkit [9]. When a sample was produced using paired-end sequencing, sample  
28 integrity was verified by confirming it contained the same number of forward and  
29 reverse reads. The average and standard deviation read count were then calculated.  
30 The algorithms and packages used in our pipeline are summarized below:

Stage	Function	Software	Ref	Link
0	Preliminary	<i>fastq-dump</i>	[9]	<a href="https://ncbi.github.io/sra-tools/fastq-dump.html">https://ncbi.github.io/sra-tools/fastq-dump.html</a>
1	Quality control	<i>illumina-utils</i>	[10]	<a href="https://github.com/merenlab/illumina-utils">https://github.com/merenlab/illumina-utils</a>
2	Co-assembly	<i>MEGAHIT</i> <i>Anvi'o</i>	[11] [12]	<a href="https://github.com/voutcn/megahit">https://github.com/voutcn/megahit</a> <a href="https://github.com/merenlab/anvi'o">https://github.com/merenlab/anvi'o</a>
3	Binning	<i>MaxBin</i>	[13]	<a href="https://downloads.jbei.org/data/microbial_communities/MaxBin/MaxBin.html">https://downloads.jbei.org/data/microbial_communities/MaxBin/MaxBin.html</a>
4	Refine bins	<i>CheckM</i>	[14]	<a href="http://ecogenomics.github.io/CheckM/">http://ecogenomics.github.io/CheckM/</a>
5	Functional annotation	<i>Prokka</i>	[15]	<a href="https://github.com/tseemann/prokka">https://github.com/tseemann/prokka</a>
6	Taxonomy classification of bins	<i>PhyloPhlan</i>	[16]	<a href="https://bitbucket.org/nsegata/phylophlan/wiki/Home">https://bitbucket.org/nsegata/phylophlan/wiki/Home</a>
7	Refine functional annotation of	<i>eggNOG-mapper</i>	[17]	<a href="https://github.com/jhcepas/eggnoG-mapper">https://github.com/jhcepas/eggnoG-mapper</a>

	methanotrophic MAGs			
--	------------------------	--	--	--

31 For Stage 1 (quality control), high-quality reads were selected using *illumina-utils*  
32 [10] with the Minoche [18] method using the command *iu-filter-quality-minoche* with  
33 the parameter *--ignore-deflines*. As the only exception, we used the method *iu-*  
34 *merge-pairs* for the metagenome from Serpentinite Springs of the Voltri Massif, Italy  
35 as the authors reported that the sequencing in this project yielded partially  
36 overlapped paired-end reads. In Stage 2 (co-assembly), pooled samples from the  
37 same environment were co-assembled following published methods [19–21]. Briefly,  
38 reads that passed quality control were co-assembled into contigs with *MEGAHIT*  
39 using the parameter *--min-contig-len = 1000*. Contigs produced by *MEGAHIT* [11]  
40 were subjected to refinement with *anvi'o* [12] three times using the parameter *--*  
41 *simplify-names* and setting *--min-len* to 1000, 1500, and 2500. In Stage 3 (binning),  
42 refined contigs and high-quality reads were binned with *MaxBin* [13] to produce  
43 MAGs (MAGs and bins refer to the same item in this pipeline). In Stage 4 (refine  
44 bins), the quality of MAGs was assessed with *CheckM* [14] and retained if they  
45 exceeded the quality standards defined in the Minimum Information about a  
46 Metagenome-Assembled Genome (MIMAG) for bacteria [22] for a medium-quality  
47 draft (completeness > 70% and contamination < 10%). MAGs of potential  
48 methanotrophs were selected according to the presence of methane oxidation genes  
49 (*pmoCAB* or *mmoXYZCDB*) in their genomes. In Stage 5 (functional annotation),  
50 MAGs were annotated with *Prokka* [15] using the parameters *--metagenome* and *--*  
51 *kingdom=Bacteria*. In Stage 6 (taxonomic classification), the amino acid sequences  
52 produced by *Prokka* were used as input for taxonomic characterization of MAGs  
53 using *PhyloPhlAn* [16] with parameters *-i* and *-t*. Only MAGs resulting in a taxonomic  
54 classification with high or medium confidence were selected for subsequent  
55 analyses. Finally, in Stage 7 (functional annotation refinement), all MAGs  
56 characterized as probable methanotrophic bacteria were subjected to a second,  
57 more comprehensive annotation procedures with *eggNOG-mapper* [17], in which  
58 KEGG Orthologs (KO), Gene Ontology (GO), and Clusters of Orthologous Groups  
59 (COGs) were assigned to genome features. The retrieved publicly available

60 assembled and binned MAGs of methanotrophic bacteria of Lake Washington, USA  
61 [4] were subjected to our metagenomics analysis pipeline from Stage 4 to Stage 7.

62 **Geographical location of methanotroph genomes.** The geographical coordinates  
63 of the origin for each sample were determined either from manual inspection of  
64 published reports (Table S1) or IMG/JGI [23]. When available, the exact coordinates  
65 of the sampling location were used to place genomes in the map. For nine genomes,  
66 the origin location could not be found using either method. The coordinates were  
67 plotted using the *maps* [24] package in *R*. The *position\_jitter* parameters were set to  
68  $w = 3.1$  and  $h = 3.1$  to avoid overlapping of dots.

69 **Genome-scale phylogenetic tree of genomes and MAGs.** 59 methanotroph  
70 genomes and MAGs and one outgroup genome of the non-methanotrophic  
71 bacterium *Bacteroides ovatus* ATCC 8483 were used to reconstruct the phylogenetic  
72 tree with *PhyloPhlAn* [16] with parameter *-u* (*de novo* phylogenetic tree).

73 **Incorporating metadata and nucleotide content into genome-scale phylogeny.**  
74 The resultant phylogenetic tree (raw tree *1\_proteomes\_tree.nwk* available in GitHub  
75 repository) was imported to *R* using the *ape* package [25] and re-rooted to the  
76 outgroup genome of *B. ovatus* ATCC 8483. The outgroup genome was selected  
77 based on its close placement to known methanotrophs in the microbial tree of life  
78 [26]. Metadata of genomes and MAGs were also imported in order to assign features  
79 to each genome and to differentiate the seven methanotroph types. Methanotroph  
80 types were assigned using the *treeio* [27] *R* package. In the tree, number of coding  
81 sequences (CDSs) and distribution of GC and GC<sub>3</sub> content were plotted using the *R*  
82 packages *ggtree* [28] and *ggridges* [29]. GC and GC<sub>3</sub> compositions of CDSs were  
83 determined using the *gc* and *gc3* functions of the *seqinr* [30] *R* package. The  
84 standalone version of *EMBOSS* [31] was used to corroborate the GC and GC<sub>3</sub>  
85 content of each CDS of our interest. All the data were compiled and manually  
86 curated and are available in the file *1\_QC\_CH4.txt* in our GitHub repository.

87 **Analysis of relative synonymous codon usage (RSCU).** The frequency of  
88 individual codon usage per CDS normalized to the amino acid usage of its

89 corresponding protein was calculated as RSCU [32] using the function *uco* with  
90 parameter *index = rscu* from the *seqinr R* package. The equation used to calculate  
91 the RSCU is:

$$RSCU = \frac{O_{ij}}{\frac{[\sum_j^{n_i} O_{ij}] * 1}{n_i}} \quad (1)$$

92 where  $O_{ij}$  is the occurrence of the  $j$ th codon for the  $i$ th amino acid and  $n_i$  the total  
93 number of synonymous codons coding for the  $i$ th amino acid. We considered a  
94 codon frequently used if  $RSCU \geq 1.6$ , and rarely used if  $RSCU \leq 0.6$ . Principal  
95 component analysis (PCA) was computed using the *R* function *prcomp* to identify  
96 CDSs that share similar preferences codon usage biases based on RSCU values for  
97 59 codons (the conventional set of 64 codons excluding the two non-redundant  
98 codons for methionine and tryptophan, which have a fixed  $RSCU = 1.0$ , and the  
99 three stop codons).

100 **Calculation of the codon adaptation index (CAI).** The CAI [32] was used to  
101 analyze the codon usage of each CDS relative to a reference set of CDSs. Codon  
102 frequency was calculated for each CDS in each of the 67 isolate genomes and  
103 MAGs. Frequencies were calculated for a single reading frame of the CDS and only  
104 ~1% of all CDSs had length not divisible by three. The codon relative adaptiveness  
105 ( $w$ ) was calculated as the frequency of a codon divided by the frequency of the  
106 synonymous codon with the highest frequency [32].  $w$  values were used to compute  
107 the CAI for each codon using the *cai* function from the *seqinr R* package, using  
108 either the full set of CDSs in the genome ( $CAI_{\text{genome}}$ ) or only ribosomal protein genes  
109 ( $CAI_{\text{ribosome}}$ ). The percentile rank of each CDS within the distribution of  
110  $CAI_{\text{genome}}/CAI_{\text{ribosome}}$  was calculated.

111 **Analysis of the effective number of codons (ENC).** ENC [33] is a measure of  
112 CDS codon usage bias based on codon preference per amino acid and has been

113 applied recently to study genomes assembled from environmental samples [34, 35].  
114 ENC values were computed for each CDS using the *chips* program from *EMBOSS*  
115 [31] based on Wright's equation [33]:

$$ENC = 2 + \frac{9}{\hat{F}_2} + \frac{1}{\hat{F}_3} + \frac{5}{\hat{F}_4} + \frac{3}{\hat{F}_6} \quad (2)$$

116 where  $\hat{F}_i$  is the codon homozygosity for the amino acids of degeneracy  $i$ . ENC as a  
117 function of GC<sub>3</sub> content was analyzed in all methanotrophs. A linear model relating  
118 ENC and GC<sub>3</sub> was fitted using the *stat\_smooth* function from the *ggplot2* R package.

119 **tRNA copy number.** tRNA frequencies were analyzed only for isolate genomes  
120 where the total tRNA pool should be known. When available, the tRNA counts of the  
121 genomes were downloaded from the public databases of IMG/JGI and GtRNAdb [36,  
122 37], otherwise they were computed with the local version of tRNAscan-SE 2.0 [38].

123 **tRNA adaptation index (tAI).** tAI was developed to estimate translation efficiency  
124 [39, 40]. The tAI was calculated for all CDSs in each genome using the *R* package  
125 *codonR* [40] with the parameter *sking* set to 1 (Prokaryote super kingdom) and the  
126 default *s* parameter for codon selection penalties. The tAI computation required the  
127 genomic tRNA counts (Table S2) and CDS codon frequencies, which were  
128 calculated using *CodonM*. Within-genome tAI percentile ranks were calculated for  
129 each CDS. The manually curated dataset containing the tAI data for our CDSs of  
130 interest can be found in the file *1\_QC\_V\_manuallycurated.txt* in our GitHub  
131 repository.

132 **Interaction network of codons and tRNAs.** The interaction network was  
133 reconstructed for isolates of type Ia methanotrophs based on RSCU values for six  
134 CDS sets, tRNA copy numbers and codon-anticodon pairing rules. Four CDS sets  
135 (*pmoCAB*, *mmoXYZCDB*, *mxoFI* and *xoxF*) represented the methane oxidation  
136 metabolic module, one set comprised ribosomal protein genes, and one set  
137 comprised all the CDSs in each genome. The median RSCU of each codon for each

138 CDS set was computed from the distribution of RSCU values of all type Ia  
139 methanotrophs. The median copy number of each tRNA anticodon was calculated  
140 from all the copy numbers of all tRNA anticodons in type Ia methanotrophs. The  
141 tRNA anticodon matrix is shown in Fig. S4e. Standard codon-anticodon recognition  
142 rules [40] were used and are detailed in *0\_wobble\_pairing\_rules.txt* available in our  
143 GitHub repository.

144 The integrated dataset was transformed into a network of sources (tRNA anticodons)  
145 and targets (CDS codons). The raw network matrix can be found in the file  
146 *RSCU\_complete\_network.txt* in GitHub. The matrix was imported to Cytoscape [41]  
147 and edited as shown in the Cytoscape file *1\_Fig3D\_net\_cytoscape.cys*. The  
148 complete network containing all amino acids and codons is shown in Fig. S5a. A  
149 quantitative analysis was applied to the raw network matrix  
150 (*RSCU\_complete\_network.txt*). For each CDS, the number of accessible tRNA  
151 copies was calculated for a range of RSCU thresholds (e.g. for RSCU threshold = 0  
152 each CDS can access every possible tRNA). This allowed the number of tRNA  
153 copies that a CDS can access as function of codon bias usage (as determined by  
154 RSCU) to be calculated. For each CDS, access to the tRNA pool can be measured  
155 in absolute term and relative to the tRNA pool available when compared with the  
156 access granted to other CDSs. The significance of the difference between two states  
157 (RSCU = 0 and RSCU = 2) was assessed with a Chi-square test, with accessible  
158 tRNA copies at RSCU = 0.0 as the expected value and at RSCU = 2.0 as the  
159 observed value. The test was applied only to CDSs of the methane oxidation  
160 metabolic module.

161 **CDS amino acid composition.** For each CDS, the codon exhibiting the highest  
162 median RSCU for each amino acid was selected. Methionine and tryptophan were  
163 excluded as they are each encoded by only one codon. The median and standard  
164 deviation were calculated from the distribution of RSCU values of each type of  
165 methanotrophs. The median and standard deviation of amino acid composition of  
166 each translated CDS of each operon was calculated using the distribution of the  
167 amino acid composition of each type of methanotrophs. A linear model (using the *lm*

168 function in  $R$ ) was fitted to determine the relationship between codon preference and  
169 amino acid usage to serve as a proxy to identify selection for optimal codons at  
170 synonymous sites occupied by the most abundant amino acids.

171 **Prebiotic amino acids analysis.** The amino acid content of each protein in the  
172 metabolic module of all methanotrophs was calculated from its translated CDS.  
173 Amino acids were categorized as ‘cheap’/prebiotic (alanine, aspartic acid, glutamic  
174 acid, glycine, isoleucine, leucine, proline, serine, threonine, and valine) or  
175 ‘expensive’/modern amino acids [42–44]. A  $t$ -test was used to compute the statistical  
176 significance of the difference between modern and prebiotic amino acid composition  
177 of each protein in the metabolic module from type Ia methanotrophs. The sample for  
178 this test was the median amino acid composition (Fig. S5b).

179 **Transcriptome analysis.** Transcribed CDSs were analyzed by modifying thymine  
180 (T) for uracil (U) in all CDSs. Three publicly available type Ia methanotroph  
181 transcriptomic datasets (*Methylobacterium buryatense* 5G [45], *Methylobacterium*  
182 *alcaliphilum* 20Z [46] and *Methylobacter tundripaludum* 31/32 [47]) were used.  
183 Normalized mRNA abundance was obtained from each dataset as reported. The  
184 purine (A+G) and pyrimidine (T+C) content of each transcribed CDS was calculated.  
185 The purine and pyrimidine content of each transcriptome was calculated based on  
186 the ribonucleotide composition of each transcribed CDS multiplied by the transcript  
187 abundance, summed across all transcribed CDSs. The effect on transcriptome  
188 composition of removing a set of transcripts was calculated by subtracting the total  
189 transcribed CDS composition (transcribed CDS composition  $\times$  transcript abundance)  
190 from the dataset and re-calculating the total ribonucleotide composition.

191 **Elemental composition of transcribed CDSs.** The carbon (C), hydrogen (H),  
192 oxygen (O) and nitrogen (N) composition of transcribed CDSs was calculated based  
193 on ribonucleotide molecular formulae (adenine  $C_5H_5N_5$ , guanine  $C_5H_5N_5O$ , cytosine  
194  $C_4H_5N_3O$ , uracil  $C_4H_4N_2O_2$ ) and normalized to the number of codons in each CDS.  
195 To provide statistical support for the observations of elemental composition bias in  
196 *pmoCAB* transcripts, the mean per-codon elemental content of 1,000 randomly  
197 selected combinations of three transcribed CDSs was calculated.

198 **Correlation between transcript abundance and elemental composition.** It has  
199 been recently proposed that highly expressed genes tend to decrease per-codon  
200 nitrogen requirements of their RNA transcripts [48, 49]. The relationship between  
201 elemental composition and mRNA abundance was investigated by computing the  
202 Pearson correlation coefficient with 95% confidence levels.

## 203 **References**

- 204 1. R Core Team. R: A language and environment for statistical computing. R  
205 Foundation for Statistical Computing, Vienna, Austria. 2017.
- 206 2. RStudio Team. RStudio: Integrated development for R. RStudio, Inc., Boston, MA.  
207 2015.
- 208 3. Wickham, H. ggplot2: elegant graphics for data analysis. Springer-Verlag New  
209 York. 2016.
- 210 4. Oshkin, I. Y. et al. Methane-fed microbial microcosms show differential community  
211 dynamics and pinpoint taxa involved in communal response. *ISME J.* 2015; 9:1119–  
212 1129.
- 213 5. Brazelton, W. J. et al. Metagenomic identification of active methanogens and  
214 methanotrophs in serpentinite springs of the Voltri Massif, Italy. *PeerJ.* 2017;  
215 5:e2945.
- 216 6. Kumaresan, D. et al. Aerobic proteobacterial methylotrophs in Movile Cave:  
217 genomic and metagenomic analyses. *Microbiome.* 2018; 6:1.
- 218 7. Crespo-Medina, M. et al. Methane dynamics in a tropical serpentinizing  
219 environment: the Santa Elena Ophiolite, Costa Rica. *Front Microbiol.* 2017; 8: 916.
- 220 8. Padilla, C. C. et al. Metagenomic binning recovers a transcriptionally active  
221 gammaproteobacterium linking methanotrophy to partial denitrification in an anoxic  
222 oxygen minimum zone. *Front Mar Sci.* 2017; 4:23.

- 223 9. Leinonen, R., Sugawara, H. & Shumway, M. The sequence read archive. *Nucleic*  
224 *Acids Res.* 2011; 39:D19–D21.
- 225 10. Eren, M. A., Vineis, J. H., Morrison, H. G. & Sogin, M. L. A Filtering method to  
226 generate high quality short reads using Illumina paired-end technology. *PLoS ONE.*  
227 2013; 8:e66643.
- 228 11. Li, D. et al. MEGAHIT v1.0: A fast and scalable metagenome assembler driven  
229 by advanced methodologies and community practices. *Methods.* 2016; 102:3–11.
- 230 12. Eren, M. A. et al. Anvi'o: an advanced analysis and visualization platform for  
231 'omics data. *PeerJ.* 20153;e1319.
- 232 13. Wu, Y.-W., Simmons, B. A. & Singer, S. W. MaxBin 2.0: an automated binning  
233 algorithm to recover genomes from multiple metagenomic datasets. *Bioinformatics.*  
234 2016; 32:605–607.
- 235 14. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W.  
236 CheckM: assessing the quality of microbial genomes recovered from isolates, single  
237 cells, and metagenomes. *Genome Res.* 2015; 25:1043–1055.
- 238 15. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics.* 2014;  
239 30:2068–2069.
- 240 16. Segata, N., Börnigen, D., Morgan, X. C. & Huttenhower, C. PhyloPhlAn is a new  
241 method for improved phylogenetic and taxonomic placement of microbes. *Nat*  
242 *Commun.* 2013; 4.
- 243 17. Huerta-Cepas, J. et al. Fast genome-wide functional annotation through  
244 orthology assignment by eggNOG-mapper. *Mol Biol Evol.* 2017; 34:2115–2122.
- 245 18. Minoche, A. E., Dohm, J. C. & Himmelbauer, H. Evaluation of genomic high-  
246 throughput sequencing data generated on Illumina HiSeq and Genome Analyzer  
247 systems. *Genome Biol.* 2011; 12:1–15.

- 248 19. Parks, D. H. et al. Recovery of nearly 8,000 metagenome-assembled genomes  
249 substantially expands the tree of life. *Nat Microbiol.* 2017; 2:1.
- 250 20. Stewart, R. D. et al. Assembly of 913 microbial genomes from metagenomic  
251 sequencing of the cow rumen. *Nat Commun.* 2018; 9:870.
- 252 21. Delmont, T. O. et al. Nitrogen-fixing populations of Planctomycetes and  
253 Proteobacteria are abundant in surface ocean metagenomes. *Nat Microbiol.* 2018;  
254 3:804–813.
- 255 22. Bowers, R. M. et al. Minimum information about a single amplified genome  
256 (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and archaea.  
257 *Nat Biotechnol.* 2017; 35:725–731.
- 258 23. Markowitz, V. M. et al. IMG/M 4 version of the integrated metagenome  
259 comparative analysis system. *Nucleic Acids Res.* 2014; 42:D568–D573.
- 260 24. Becker, R. A., Wilks, A. R., Brownrigg, R., Minka, T. P. & Deckmyn, A. maps:  
261 draw geographical maps. R package version 3.3.0.2018.
- 262 25. Paradis, E., Claude, J. & Strimmer, K. APE: analyses of phylogenetics and  
263 evolution in R language. *Bioinformatics.* 2004; 20, 289–290.
- 264 26. Hug, L. A. et al. A new view of the tree of life. *Nat Microbiol.* 2016; 1:16048.
- 265 27. Yu, G. & Tsan-Yuk Lam, T. treeio: base classes and functions for phylogenetic  
266 tree input and output. R package version 1.2.2. 2018.
- 267 28. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. ggtree: an R package for  
268 visualization and annotation of phylogenetic trees with their covariates and other  
269 associated data. *Methods Ecol Evol.* 2017; 8:28–36.
- 270 29. Wilke, C. O. ggrydges: ridgeline plots in 'ggplot2'. R package version 0.5.0.2018.

- 271 30. Charif, D. & Lobry, J. R. Biological and medical physics, biomedical engineering.  
272 2007:207–232. doi:10.1007/978-3-540-35306-5\_10
- 273 31. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology  
274 Open Software Suite. *Trends Genet.* 2000; 16:276–277.
- 275 32. Sharp, P. M. & Li, W.-H. The codon adaptation index—a measure of directional  
276 synonymous codon usage bias, and its potential applications. *Nucleic Acids Res.*  
277 1987; 15:1281–1295.
- 278 33. Wright, F. The ‘effective number of codons’ used in a gene. *Gene.* 1990; 87:23–  
279 29.
- 280 34. Alves, R. J., Minh, B., Urich, T., von Haeseler, A. & Schleper, C. Unifying the  
281 global phylogeny and environmental distribution of ammonia-oxidising archaea  
282 based on *amoA* genes. *Nat Commun.* 2018; 9:1517.
- 283 35. Zhalnina, K. et al. Dynamic root exudate chemistry and microbial substrate  
284 preferences drive patterns in rhizosphere microbial community assembly. *Nat*  
285 *Microbiol.* 2018; 3:470–480.
- 286 36. Chan, P. P. & Lowe, T. M. GtRNAdb: a database of transfer RNA genes detected  
287 in genomic sequence. *Nucleic Acids Res.* 2009; 37:D93–D97.
- 288 37. Chan, P. P. & Lowe, T. M. GtRNAdb 2.0: an expanded database of transfer RNA  
289 genes identified in complete and draft genomes. *Nucleic Acids Res.* 2016; 44:D184–  
290 D189.
- 291 38. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: A program for improved detection of  
292 transfer RNA genes in genomic sequence. *Nucleic Acids Res.* 1997; 25:955–964.
- 293 39. dos Reis, M., Wernisch, L. & Savva, R. Unexpected correlations between gene  
294 expression and codon usage bias from microarray data for the whole *Escherichia*  
295 *coli* K-12 genome. *Nucleic Acids Res.* 2003; 31:6976–6985.

296 40. dos Reis, M., Savva, R. & Wernisch, L. Solving the riddle of codon usage  
297 preferences: a test for translational selection. *Nucleic Acids Res.* 2004; 32:5036–  
298 5044.

299 41. Shannon, P. et al. Cytoscape: a software environment for integrated models of  
300 biomolecular interaction networks. *Genome Res.* 2003; 13:2498–2504.

301 42. Li, N., Lv, J. & Niu, D.-K. Low contents of carbon and nitrogen in highly abundant  
302 proteins: evidence of selection for the economy of atomic composition. *J Mol Evol.*  
303 2009; 68:248–255.

304 43. Longo, L. M., Lee, J. & Blaber, M. Simplified protein design biased for prebiotic  
305 amino acids yields a foldable, halophilic protein. *Proc Natl Acad Sci USA.* 2013;  
306 110:2135–2139.

307 44. Blanco, L. P., Payne, B. L., Feyertag, F. & Alvarez-Ponce, D. Proteins of  
308 generalist and specialist pathogens differ in their amino acid composition. *Life Sci*  
309 *Alliance.* 2018; 1:e201800017.

310 45. Torre, A. et al. Genome-scale metabolic reconstructions and theoretical  
311 investigation of methane conversion in *Methylobacterium buryatense* strain 5G(B1).  
312 *Microb Cell Fact.* 2015; 14:1–15.

313 46. Kalyuzhnaya, M. et al. Highly efficient methane biocatalysis revealed in a  
314 methanotrophic bacterium. *Nat Commun.* 2013; 4:2785.

315 47. Krause, S. et al. Lanthanide-dependent cross-feeding of methane-derived carbon  
316 is linked by microbial community interactions. *Proc Natl Acad Sci USA.* 2017;  
317 114:358–363.

318 48. Seward, E. A. & Kelly, S. Dietary nitrogen alters codon bias and genome  
319 composition in parasitic microorganisms. *Genome Biol.* 2016; 17:226.

320 49. Seward, E. A. & Kelly, S. Selection-driven cost-efficiency optimization of  
321 transcripts modulates gene evolutionary rate in bacteria. *Genome Biol.* 2018; 19:102.