

1 **DiTing: A Pipeline to Infer and Compare Biogeochemical Pathways from**
2 **Metagenomic and Metatranscriptomic Data**

3

4 **Running title:** DiTing for Biogeochemical Pathways

5

6 Chun-Xu Xue^{1†}, Heyu Lin^{2†}, Xiao-Yu Zhu¹, Jiwen Liu^{1,3,4}, Gary Rowley⁵, Jonathan D Todd⁵, Meng
7 Li⁶, Xiao-Hua Zhang^{1,3,4*}

8 ¹College of Marine Life Sciences, and Institute of Evolution & Marine Biodiversity, Ocean
9 University of China, Qingdao 266003, China

10 ²School of Earth Sciences, University of Melbourne, Parkville, Victoria 3010, Australia

11 ³Laboratory for Marine Ecology and Environmental Science, Qingdao National Laboratory for
12 Marine Science and Technology, Qingdao 266071, China

13 ⁴Frontiers Science Center for Deep Ocean Multispheres and Earth System, Ocean University of
14 China, Qingdao 266100, China

15 ⁵School of Biological Sciences, University of East Anglia, Norwich Research Park, Norwich,
16 United Kingdom.

17 ⁶Shenzhen Key Laboratory of Marine Microbiome Engineering, Institute for Advanced Study,
18 Shenzhen University, Shenzhen 518060, China.

19

20 [†]Chun-Xu Xue and Heyu Lin contributed equally to this work.

21 ^{*}To whom correspondence should be addressed.

22

23 **Abstract**

24 Metagenomics and metatranscriptomics are powerful tools to uncover key microbes
25 and processes driving biogeochemical cycling in natural ecosystems. Currently
26 available databases depicting metabolic functions from
27 metagenomic/metatranscriptomic data are not dedicated to biogeochemical cycles.

28 There are no databases encompass genes involved in the cycling of
29 dimethylsulfoniopropionate (DMSP), an abundant organosulfur compound.
30 Additionally, a recognized normalization mode to estimate and compare the relative
31 abundance and environmental importance of pathways from metagenomic and
32 metatranscriptomic data has not been available. These limitations impact the ability to
33 accurately relate key microbial driven biogeochemical processes to differences in
34 environmental conditions. Thus, an easy to use specialized tool that infers and
35 visually compares the potential for biogeochemical processes, including DMSP
36 cycling, is urgently required. To solve these issues, we developed DiTing, a tool
37 wrapper to infer and compare biogeochemical pathways among a set of given
38 metagenomic or metatranscriptomic reads in one step, based on the KEGG (Kyoto
39 Encyclopedia of Genes and Genomes) and a manually created DMSP cycling gene
40 database. Accurate and specific formulas for over 100 pathways were developed to
41 calculate their relative abundance. Output reports detail the relative abundance of
42 biogeochemically-relevant pathways in both text and graphical format. We applied
43 DiTing to metagenomes from simulated data, hydrothermal vents and the *Tara* Ocean
44 project. The DiTing outputs were consistent with genetic feature of genomes used in
45 simulated benchmark data, and also demonstrated that the predicted functional
46 profiles correlated strongly with changes in environmental conditions. DiTing can
47 now be confidently applied to wider metagenomic and metatranscriptomic datasets.

48 **Availability and implementation:** <https://github.com/xuechunxu/DiTing>

49 **Contact:** xhzhang@ouc.edu.cn

50 **Supplementary information:** Supplementary data are available at Molecular
51 Ecology Resources online.

52

53

54

55 **Introduction**

56 Microbial communities play integral and unique roles in mediating global
57 biogeochemical cycles. Sequencing techniques, such as amplicon sequencing
58 (Bokulich et al., 2013), whole-genome sequencing (Jones et al., 2016), genome-
59 resolved metagenomics (Parks et al., 2017; Xue et al., 2020b) and shotgun
60 metagenomic sequencing (Sharpton et al., 2014; Xue et al., 2020a), are widely used to
61 characterize the genetic potential of microbial communities. Metagenomics is an
62 important tool to unravel the diversity, function and ecology of complex microbial
63 ecosystems via quantification of the genetic potential for various biogeochemical
64 pathways within microbial communities (Riesenfeld et al., 2004). Moreover,
65 metatranscriptomic data present a more accurate scenario of processes occurring
66 within ecosystems because these methodologies move past genetic potential and
67 report on the transcription of biogeochemical pathways (Aguiar-Pulido et al., 2016;
68 Shakya et al., 2019). Previous studies have predicted community functions according
69 to gene annotation against several established databases, e.g., KEGG (Ogata et al.,
70 2000), COG (Tatusov et al., 2000), MetaCyc (Caspi et al., 2006), Pfam (Finn et al.,
71 2014), TIGRFam (Selengut et al., 2007), SEED (Ross et al., 2014), and eggNOG
72 (Huertacepas et al., 2016). However, these functional annotations are not dedicated to
73 biogeochemical cycling and lack comprehensive lists of annotated genes for important
74 cycles. FOAM (Functional Ontology Assignments for Metagenomes; Prestat et al.,
75 2014) is a functional gene database for environmental datasets that includes
76 biogeochemical cycles, however, this database lacks visualization, and annotates all
77 protein sequences with a universal threshold value, which may lead to prediction
78 biases. Furthermore, some metabolic pathways, e.g. the cycling of
79 dimethylsulfoniopropionate (DMSP), a key marine osmolyte, nutrient, and signaling
80 molecule, with important roles in sulfur cycling (Curson *et al.*, 2011; Zhang *et al.*,
81 2019), lack an accurate and reviewed database for annotating the key metabolic genes.
82 These limitations force researchers to undertake often tricky and time-consuming
83 gathering of gene sequences from primary research and collate them into robust local

84 databases (Acinas et al., 2019; Dombrowski et al., 2018; Llorens-Marès et al., 2015;
85 Zhang et al., 2018). This can also lead to challenges for downstream interpretation,
86 organization and visualization.

87 Additionally, there is no recognized and prepared normalization method to
88 estimate and compare the relative abundance of a pathway in metagenomic and
89 metatranscriptomic data. In some studies, the relative abundance of every gene in a
90 biogeochemical pathway was added together to estimate the relative abundance of the
91 pathway (Ganesh et al., 2014; Petter et al., 2013; Smedile et al., 2013), which is
92 unsuitable to infer and compare pathways. For example, thiosulfate disproportionation
93 (thiosulfate \rightarrow sulfide & sulfite) is catalyzed by thiosulfate reductase, which is
94 encoded by three genes (*phsABC*). Thus, the relative abundance of thiosulfate
95 disproportionation pathway should be equal to the mean relative abundance of
96 *phsABC* instead of the sum of *phsABC* relative abundance together. This
97 normalization mode was applied in somerecent studies (Llorens-Marès et al., 2015,
98 Graham et al., 2018), but no simple tool to achieve this is currently available. In
99 addition, there are few easy methods for high throughput comparison and
100 visualization of samples. Therefore, new automated tools to identify, quantify, and
101 compare the abundance and/or transcription of genes and pathways for
102 biogeochemical cycles, including the DMSP cycle, are needed.

103 Here we developed the software DiTing, a pipeline to infer and compare
104 biogeochemical pathways in metagenomic and metatranscriptomic data. DiTing is
105 named after a Chinese mythical creature who knows everything when he puts his ears
106 on the Earth's surface. Similarly, scientists can gain robust knowledge on microbial
107 driven biogeochemical cycles from environmental 'omic data after analysis with
108 DiTing. DiTing annotates protein sequences based on the KEGG database (Ogata et
109 al., 2000) for most microbial-mediated biogeochemical cycles, and an in-house
110 database specifically for cycling of DMSP, and then estimates the relative abundance
111 of corresponding functional genes. More accurate specific formula for each pathway
112 were developed to calculate the relative abundance of multiple pathways. The output

113 results consist of user-friendly tables containing a summary of over 100
114 biogeochemically-relevant pathways and corresponding genes, and their relative
115 abundances in individual metagenomic/metatranscriptomic samples, alongside
116 graphical outputs consisting of heatmaps and multiple sketch plots for easier
117 visualization.

118

119 **2 Methods**

120 **The main procedure of DiTing**

121 DiTing was written in Python 3 and runs on Linux/Unix platforms. The pre-requisites
122 required for running the software are described on the DiTing GitHub page
123 (<https://github.com/xuechunxu/DiTing>). The input source was a set of metagenomic
124 and/or metatranscriptomic clean reads where low-quality reads, primer and adaptor
125 sequences had been trimmed beforehand (Fig. 1), which were then assembled by
126 Megahit v1.1.2 (Li et al., 2016) (with default parameters) or metaSPAdes v3.12.0
127 (Nurk et al., 2017) (with default parameters). Users can set distinct parameter to
128 choose which software for reads assembly. Compared to Megahit, MetaSPAdes
129 performs better in recovering long contigs, it has a higher assembly quality index and
130 is the recommended assembler for high-complex metagenomes (Forouzan et al., 2018,
131 Pasolli et al., 2019). However, Megahit has a low error rate, is highly memory-
132 efficient and is ideal for large datasets (Forouzan et al., 2018). Genes were predicted
133 and translated from the assembled contigs by Prodigal v2.6.3 with the ‘-p meta’
134 option (Hyatt et al., 2010). To determine the relative abundance of each gene, the
135 input metagenomic reads were mapped against predicted genes (nucleotide) by BWA-
136 MEM (Li, 2013) (bwa v0.7.15, default settings) to generate sequence alignment map
137 (SAM) files. We used the unsorted SAM files as input for pileup.sh (bbmap v38.22)
138 (Bushnell, 2014) (with default parameters) to calculate the average coverage of each
139 gene or transcript. The TPM methodology was used to indicate the relative abundance
140 of a gene by the following formula.

$$141 \quad TPM_i = \frac{b_i}{\sum_j b_j} \cdot 10^6 = \frac{\frac{X_i}{L_i}}{\sum_j \frac{X_j}{L_j}} \cdot 10^6$$

142 Where TPM_i is the relative abundance of gene i , b_i is the copy number of gene i , L_i is
143 the length of gene i , X_i is the number of times that gene i is detected in a sample (that
144 is, the number of reads in alignment), and j is the number of genes in a sample. The
145 translated protein sequences were queried against KOfam database (HMM database of
146 KEGG Orthologs; KOs) (Aramaki et al., 2019) using hmmsearch implemented within
147 HMMER (Finn et al., 2011) (parameter: `hmmsearch -T <threshold> --tblout <output>`
148 `<hmm database> <input protein sequence>` when score type is full; `hmmsearch --`
149 `domT <threshold> --domtblout <output> <hmm database> <input protein sequence>`
150 when score type is domain), which employs methods detecting remote homologs
151 sensitively and efficiently. Kofam suggested values
152 (<ftp://ftp.genome.jp/pub/db/kofam/>) were used as the cutoff threshold values for
153 `hmmsearch`, in which each KEGG Orthology (KO) entry had its unique cutoff
154 threshold values (Aramaki et al., 2019). To test the accuracy of the gene annotation
155 from DiTing, we also submitted translated protein sequences to the KofamKOALA
156 web server (<https://www.genome.jp/tools/kofamkoala/>). KofamKOALA assigns KOs
157 numbers to protein sequences with its accuracy being comparable to the best existing
158 KO assignment tools (Aramaki et al., 2019). For genes assigned into multiple KOs
159 numbers, all the corresponding functions were associated to the genes. To specifically
160 probe DMSP catabolism, 20 verified gene sequences (DMSP lyase genes *dddD*, *dddK*,
161 *dddL*, *dddP*, *dddQ*, *dddY*, *dddW*, *Alma1*; DMSP synthesis genes *dsyB*, *DSYB*, *mmtN*;
162 DMSP demethylation pathway genes *dmdA*, *dmdB*, *dmdC*, *dmdD*; acryloyl-CoA
163 hydratase *acuH*, methanethiol *S*-methylase *mddA*, DMS monooxygenase *dmoA*,
164 methanethiol oxidase *MTO*, and DMSO reductase *dorA*) were collected manually to
165 create the profile HMM (Song et al., 2020). A table with the relative abundance and
166 annotation of genes is used to estimate the relative abundance of approximately one

167 hundred biogeochemical pathways in each sample.

168 The formula for each pathway is specifically designed to estimate the relative
169 abundance of the pathway according to the definitions ([https://github.com/xuechunxu/
170 DiTing/blob/master/Pathway_formulas.txt](https://github.com/xuechunxu/DiTing/blob/master/Pathway_formulas.txt)). For example, assimilatory sulfite
171 reduction (ASR) that converts sulfite to sulfide has two known possible pathways: (1)
172 Sir protein (K00392) mediated pathway (Gisselmann et al., 1993; Bork et al., 1998),
173 and (2) CysJI protein (K00380 + K00381) mediated pathway (Ostrowski et al., 1989a,
174 b; Zeghouf et al., 2000). Thus, the relative abundance of assimilatory sulfite reduction
175 pathway is estimated by the following formula:

$$176 \quad A_{ASR} = a_{K00392} + \frac{a_{K00380} + a_{K00381}}{2}$$

177 Where A_{ASR} is the relative abundance of the ASR pathway, a_{KO} is the relative
178 abundance of KO in each sample. Dissimilatory nitrite reduction (DNRA), which
179 converts nitrite to ammonia, can occur via two different enzymatic reactions: (1)
180 *NirBD* proteins (K00362 + K00363) to convert nitrite to ammonia, or (2) *NrfAH*
181 protein (K03385 + K15876) to convert nitrite to ammonia. Thus, the relative
182 abundance of dissimilatory nitrite reduction to ammonia is estimated by the following
183 formula:

$$184 \quad A_{DNRA} = \frac{a_{K00362} + a_{K00363}}{2} + \frac{a_{K03385} + a_{K15876}}{2}$$

185 Where A_{DNRA} is the relative abundance of DNRA pathway, a_{KO} is the relative
186 abundance of KO in each sample. For other pathways, a customized formula for each
187 pathway was utilized (see Supplemental Table S1).

188 DiTing produces a table in the specified output directory. This table contains
189 approximately 100 biogeochemical pathways and their relative abundance in each
190 input sample. Another table of the relative abundances of corresponding KO/genes
191 within these pathways in each sample is also generated (like Supplemental Table S2).
192 Researchers can evaluate the completeness of pathways from this table. For improved

193 visualization, heatmaps and sketch plots for comparing the relative abundances of
194 biogeochemical pathways in different samples are drawn by a Python script. DiTing
195 can be installed via Conda (<https://docs.conda.io>).

196 **Construction of the DMSP database and other selected genes**

197 DMSP is a marine organosulfur compound with important roles in global sulfur cycle
198 and may affect climate (Zhang *et al.*, 2019), yet genes involved in the cycling of this
199 compound are missing in currently available databases. Profile HMM were manually
200 generated for eight pathways related to the cycling of DMSP (Song *et al.*, 2020),
201 including DMSP biosynthesis (methionine → DMSP), DMSP demethylation (DMSP -
202 > MMPA), DMSP demethylation (MMPA → MeSH), DMSP cleavage (DMSP →
203 DMS), DMS oxidation (DMS → MeSH), DMS oxidation (DMS → DMSO), DMSO
204 reduction (DMSO → DMS), MddA pathway (MeSH → DMS), MeSH oxidation
205 (MeSH → Formaldehyde). 20 verified gene sequences encoding key enzymes of these
206 pathways were used to create the profile HMM (Song *et al.*, 2020).

207 (i) *DMSP biosynthesis (methionine → DMSP)*. Three gene families participating
208 in DMSP biosynthesis from methionine (Met), including DSYB, DsyB and MmtN are
209 included in DiTing. DSYB and DsyB are methylthiohydroxybutyrate *S*-
210 methyltransferase enzymes found in marine eukaryotes and prokaryotes, respectively
211 (Curson *et al.*, 2018; Curson *et al.*, 2017). The MmtN Met *S*-methyltransferase is
212 found in some Gram-positive bacteria, alpha- and gamma-proteobacteria (Liao *et al.*,
213 2019; Williams *et al.*, 2019). The cut-off E-values of DSYB, DsyB and MmtN are $1 \times$
214 10^{-30} , 1×10^{-67} and 1×10^{-98} , respectively.

215 (ii) *DMSP demethylation (DMSP → MMPA)*. The first step of DMSP
216 demethylation pathway that results in the production of methylmercaptopropionate
217 (MMPA) is initiated by the DmdA enzyme (Reisch *et al.*, 2011). The cut-off E-values
218 of the DmdA is 1×10^{-130} .

219 (iii) *DMSP demethylation (MMPA → MeSH)*. Further degradation of MMPA
220 generating gaseous methanethiol (MeSH) catalyzed by the DmdBCD/AcuH enzymes

221 (Reisch et al., 2011; Shao et al., 2019). The cut-off E-values of DmdB, DmdC, DmdD
222 and AcuH are 1×10^{-75} , 1×10^{-100} , 1×10^{-30} and 1×10^{-56} , respectively.

223 (iv) *DMSP cleavage (DMSP -> DMS)*. Eight distinct DMSP lyase enzymes
224 (DddD, DddK, DddL, DddP, DddQ, DddW, DddY and Alma1) can cleave DMSP to
225 generate dimethylsulfide (DMS) (Curson et al., 2011; Alcolombri et al., 2015;
226 Johnston et al., 2016; Sun et al., 2016). The cut-off E-values of DddD, DddK, DddL,
227 DddP, DddQ, DddW, DddY and Alma1 are 1×10^{-97} , 1×10^{-35} , 1×10^{-33} , 1×10^{-83} , $1 \times$
228 10^{-20} , 1×10^{-49} , 1×10^{-64} and 1×10^{-26} , respectively.

229 (v) *DMS oxidation (DMS -> MeSH)*. DMS can be oxidized to generate MeSH
230 via the DMS monooxygenase enzyme DmoA (Boden et al., 2011). The cut-off E-
231 values of the DmoA is 1×10^{-34} .

232 (vi) *DMS oxidation (DMS -> DMSO)*. DMS can be oxidized to generate
233 dimethylsulfoxide (DMSO) by the DMS dehydrogenase complex (DdhABC)
234 (McDevitt et al., 2002) or trimethylamine monooxygenase (Tmm); (Lidbury et al.,
235 2016). The cut-off E-values of both DdhABC, DdhB and Tmm are 1×10^{-30} .

236 (vii) *MddA pathway (MeSH -> DMS)*. MeSH can be S-methylated to generate
237 DMS by the MddA enzyme (Carrión et al., 2017). The cut-off E-values of MddA is 1
238 $\times 10^{-30}$.

239 (viii) *MeSH oxidation (MeSH -> Formaldehyde)*. MeSH can also be modified
240 through another pathway catalyzed by the MeSH oxidase MTO (Eyice et al., 2018).
241 The cut-off E-values of MTO is 1×10^{-20} .

242 The sugar 6-deoxy-6-sulfoquinovose (sulfoquinovose, SQ) produced by plants,
243 algae, and cyanobacteria, is an important component of carbon and sulfur cycles
244 (Frommeyer et al., 2020). Microbial community can completely degrade SQ into
245 inorganic sulfate or hydrogen sulfide through three pathways, i.e., sulfo-Embden-
246 Meyerhof-Parnas (sulfo-EMP) (Denger et al., 2014), sulfo-Entner-Doudoroff (sulfo-
247 ED) (Felux et al., 2015), and 6-deoxy-6-sulfofructose-transaldolase (SFT) pathways

248 (Frommeyer et al., 2020).

249 (i) *sulfo-EMP pathway*. SQ is converted to 6-deoxy-6-sulfofructose (SF) through
250 an aldose/ketose isomerase YihS. The SF is phosphorylated to 6-deoxy6-
251 sulfofructosephosphate (SFP) by an ATP-dependent SF kinase YihV. The SFP is then
252 cleaved into 3-sulfolactaldehyde (SLA) and dihydroxyacetone phosphate (DHAP) by
253 an SFP aldolase YihT. Finally, the SLA is reduced via an NADH-dependent SLA
254 reductase (YihU) to DHPS, which is excreted from microorganisms. These four genes
255 *YihSVTU* were annotated through K18479, K18478, K01671 and K08318 Orthology
256 in KEGG, respectively.

257 (ii) *sulfo-ED pathway*. This pathway starts with an NAD⁺-dependent SQ
258 dehydrogenase (EC:1.1.1.390) oxidizing SQ to 6-sulfo gluconolactone (SGL). The
259 SGL is hydrolyzed to 6-deoxy-6-sulfo gluconate (SG) by an SGL lactonase
260 (EC:3.1.1.99). The SG is then converted by an SG dehydratase (EC:4.2.1.162) to 2-
261 keto-3,6-deoxy-6-sulfo-gluconate (KDSG). The KDSG is cleaved by a KDSG
262 aldolase (EC:4.1.2.58) into pyruvate and 3-sulfolactaldehyde (SLA). The SLA can be
263 oxidized by a NAD⁺-dependent SLA dehydrogenase (EC:1.2.1.97) to SL. The
264 reference sequences of these enzymes were collected manually from Uniprot database
265 (<https://www.uniprot.org/>).

266 (iii) *SFT pathway*. Three key enzymes take part in this pathway. The SQ is
267 converted to SF by an aldose/ketose isomerase, which is the same enzyme as the first
268 step of sulfo-EMP pathway. SF is cleaved to 3-sulfolactaldehyde (SLA) by SF
269 transaldolase enzyme. Finally, The SLA is oxidized by a NAD⁺-dependent SLA
270 dehydrogenase to SL. The SLA dehydrogenase is same enzyme as the last step of
271 sulfo-ED pathway. The reference sequence of SF transaldolase enzyme was collected
272 from IMG (<https://img.jgi.doe.gov/>) according to Frommeyer et al., 2020.

273 Isoprene (2-methyl-1, 3-butadiene) is an important volatile organic compound
274 emitted to the atmosphere, and has significant effect on the climate (Carrión et al.,
275 2018). Isoprene can be degraded by microbial communities with the isoprene

276 monooxygenase (IsoMO). The gene *isoA* encoding the α -subunit of IsoMO was
277 selected as marker gene for distribution, diversity and abundance of isoprene-
278 degrading pathway in environment (Carrión et al., 2018; Carrión et al., 2020). The
279 reference sequences of IsoA enzyme were collected manually from NCBI according
280 to Carrión et al., 2018.

281 **3 Results and discussion**

282 **General information of DiTing**

283 We developed a new metagenomics/metatranscriptomic analysis pipeline, DiTing, to
284 infer and compare the prevalence of genes and pathways of key biogeochemical
285 cycles. DiTing consists of four main features: (i) automated assembly, CDS
286 prediction, mapping and annotation from reads; (ii) a manually created
287 dimethylsulfoniopropionate (DMSP) cycling related gene database; (iii) accurate and
288 specific formula for DMSP and other biogeochemical pathway to calculate the
289 relative abundance of biogeochemically-relevant pathways and genes; (iv)
290 visualization of results comparing biogeochemical cycling potential between different
291 input samples. These features make DiTing a flexible and versatile tool wrapper for
292 studying biogeochemical cycles, or just as a platform to tackle metagenomic shotgun
293 sequencing data. The speed of DiTing is relatively fast. Five samples (from the
294 hydrothermal vent case study below) that are about 500 Gb in total were used to
295 evaluate the speed. The total run time for all analyses from reads to visualization was
296 ~ 33 hours using 60 CPU threads on a Linux version 4.15.0-20-generic server
297 (Ubuntu 18.04; CPU, Intel(R) Xeon(R) Gold 6140 CPU @ 2.30GHz; RAM, 256 GB).

298 **Accuracy testing of DiTing using simulated benchmark datasets**

299 To verify the accuracy of DiTing in evaluating the abundance of biogeochemical
300 pathways, CAMISIM (Fritz et al., 2019) was used to simulate three group of
301 metagenomic shotgun sequenced samples (photoautotrophs, chemoautotrophs and
302 heterotrophs group). Metagenomes from the photoautotrophic group were simulated
303 by ten *Cyanobacteria* genomes. Metagenomes from the chemoautotrophic group were

304 simulated by 10 ammonia-oxidising archaea (AOA) genomes. Metagenomes from the
305 heterotrophic group were simulated by 10 SAR11 genomes. Each group comprised
306 five metagenomic samples sequenced by Illumina 2×150 bp paired-end reads, and
307 each generated sample had a size of 5 Gb. These 15 simulated samples were fed into
308 DiTing. The overall relative abundance of biogeochemical pathways in simulated
309 samples was consistent with features of genomes used in each group (Fig. 2). For
310 example, metagenomes in the photoautotroph group possessed a high relative
311 abundance of photosynthesis related pathway genes (photosystem I, II and
312 cytochrome *b₆f* complex), which were absent in other two groups (Fig. 2). AOA are
313 the typical known bacterial ammonia oxidisers, which possesses *amoABC* genes
314 encoding the ammonia monooxygenase complex. Correspondingly, in the
315 chemoautotroph group simulated by AOA, the ammonia oxidation pathway was found
316 but was absent in other two groups analysed by DiTing (Fig. 2). In other nitrogen
317 cycle pathways, *nirKS* encoding nitrite reductase and *hzs* encoding hydrazine synthase
318 were only seen the chemoautotroph group of the DiTing results. Consistently, these
319 genes were annotated in ammonia-oxidising archaea genomes through RAST
320 annotation manually. Additionally, bacteria and archaea use F-type ATPase and V/A-
321 type ATPases, respectively, to hydrolyze ATP to ADP, respectively (Pisa et al., 2007;
322 Fillingame et al., 1997). Thus, F-type ATPase was detected in groups simulated by
323 *Cyanobacteria* and SAR11 genomes, and V/A-type ATPase was only detected in the
324 chemoautotroph group simulated by ammonia-oxidising archaea genomes. The
325 translated gene sequences (amino acid) from simulated metagenomes were submitted
326 to the KofamKOALA web server for annotation. The gene annotation results derived
327 from DiTing were the same as those from KofamKOALA web server, verifying the
328 accuracy of gene annotation.

329 **Application of DiTing on five real hydrothermal vent datasets and 15 Tara Ocean** 330 **project datasets**

331 DiTing was used to analyze the biogeochemical potential of five marine metagenomic
332 samples (Table 1; NCBI accession number: ERR1679394-1679398) generated from

333 hydrothermal vent samples taken at PACManus and North Su fields in the Manus
334 Basin (Meier et al., 2017; Table 1). The metagenomic clean reads ranged in size from
335 81 to 112 Gbp from each sample. The reads were assembled into 799,269 to
336 1,182,847 contigs with the total assembly sizes ranging from 0.58 to 1.00 Gbp. A total
337 of 5,639,558 Open Reading Frames (ORFs) within these contigs were then predicted.
338 ~18.9% (1,065,097) ORFs were annotated against KEGG databases and affiliated to
339 8128 KO entries. The relative abundances of ~100 biogeochemically-relevant
340 pathways were calculated according to our new formulas (Supplementary Table S1).
341 The relative abundance of genes within these pathways was also prepared for further
342 analyses at the gene level (Supplementary Table S2). The summary sketch and
343 heatmap plots for visualization of these pathways were generated, and these reflect
344 the different patterns of community function within metagenomic samples (Fig. 3, 4).

345 Of the five metagenomes collected in diffuse hydrothermal vent fluids, NSu-F2b
346 and NSu-F5 originated from acidic samples with sulfide (1.6 mmol l^{-1} and 0.7 mmol l^{-1}
347 H_2S , respectively) and methane (0.2 mmol l^{-1} and $0.01 \text{ mmol l}^{-1} \text{CH}_4$, respectively)
348 levels detected. The Fw-F1b, Fw-F3 and RR-F1b metagenomes originated from sites
349 with no detectable H_2S and CH_4 . Reassuringly, the NSu-F2b and NSu-F5 samples,
350 with similar environmental parameters, showed the most similar distribution patterns
351 for genes and pathways involved in the cycling of nitrogen, carbon and sulfur (Fig. 3,
352 4). Indeed, hierarchical clustering of samples according to their microbial function
353 composition showed NSu-F2b and NSu-F5 fall into one cluster and the other three
354 samples into another cluster (Supplementary Fig. S1).

355 At hydrothermal vents, chemolithoautotrophic microorganisms carry out carbon
356 fixation coupled with oxidation of reduced sulfur compounds (Meier et al., 2017). In
357 accordance, we found the relative abundance of thiosulfate oxidation, sulfite
358 oxidation, and first step of dissimilatory sulfate reduction pathways (reversible
359 conversion of sulfate to sulfite) to be more highly represented compared to other
360 sulfur cycle pathways in all five samples (Fig. 3, 4), indicating sulfate reduction and
361 sulfur oxidation as major processes in microbial sulfur cycling. This finding is

362 supported by the presence of sulfate-reducing *Nitrospirae* and sulfur-oxidizing
363 *Gammaproteobacteria* dominating microbial communities at these hydrothermal
364 vents (Meier *et al.*, 2017, 2019). In addition, assimilatory sulfate reduction and
365 thiosulfate disproportionation pathways were almost only found in NSu-F2b and NSu-
366 F5 (Fig. 3), the only samples with detectable sulfide levels, indicating microbes in
367 these samples may incorporate sulfide into the amino acids cysteine (Cys) or homo-
368 Cys. Here, the relative abundance of thiosulfate disproportionation was estimated by
369 dividing the sum of relative abundance of *phsABC* by the number ($n = 3$) of essential
370 subunits. The relative abundances of each subunit of thiosulfate reductase were often
371 not equal to each other in the metagenomes (Supplementary Table S2). For example,
372 *phsA* (encoding thiosulfate reductase subunit A) was always far more abundant than
373 *phsC* (thiosulfate reductase cytochrome B subunit) and *phsB* (thiosulfate reductase
374 electron transport protein) was not detected in any sample. This may be due to
375 insufficient sequencing depth and/or protein redundancy. Whatever the reason for
376 these discrepancies it cannot be easily solved by bioinformatics alone and culture-
377 dependent work is necessary. This phenomenon highlighted for the thiosulfate
378 disproportionation genes may also occur in other pathways, thus further analyses at
379 the gene level, not only at the pathway level, are essential in predicting the
380 biogeochemical potential of microbial communities after DiTing analysis.

381 In previously tested seawater and sediment samples, known DMSP synthesis
382 genes are always much less abundant than those for its catabolism (Curson *et al* 2017,
383 Curson *et al* 2018, Williams *et al.*, 2019). This was not the case in previously studied
384 hydrothermal samples (Song *et al.*, 2020), with the DMSP lyase gene *dddP* being the
385 only detected DMSP catabolic gene. In three out of five hydrothermal samples
386 interrogated here, the genetic potential to synthesize DMSP, through prokaryotic *dsyB*
387 and *mmtN* genes, is far less than that for DMSP catabolism (DMSP synthesis:DMSP
388 catabolism = 1:16.9) and not so dissimilar to ratios seen in seawater samples (Curson
389 *et al* 2017, Curson *et al.*, 2018, Williams *et al.*, 2019). Reasons for this discrepancy
390 between the distinct samples are unknown. The *DsyB* sequences retrieved from this

391 data were clustered with ratified DsyB proteins, not with DSYB and non-functional
392 DsyB-like proteins from *Streptomyces varsoviensis*, which support their function in
393 DMSP synthesis (Supplementary Fig. S2). Interestingly, sample NSu-F2b has higher
394 DMSP synthesis potential than any other samples due to relatively high levels of
395 bacteria with *mmtN*. As in Song et al 2020, the potential for DMSP cleavage was more
396 prominent than for DMSP demethylation (*dmdA*) in all hydrothermal samples,
397 although catabolism of MMPA, the initial product of DMSP demethylation by DmdA
398 (Howard *et al* 2006), was very abundant. This data supports DMSP cleavage being the
399 dominant DMSP catabolic pathway in hydrothermal sediments, as proposed in Song
400 et al., 2020. Alternatively, there could be novel DMSP demethylase enzymes. This
401 would explain why there were such low *dmdA* levels in hydrothermal sediment, yet
402 very high MMPA degradation potential. The potential for oxidation and reduction of
403 DMSP catabolites, DMS and methanethiol, was similar to that described in Song et
404 al., 2020, with sites NSU-F2b and F5 showing the greatest potential. Thus, some
405 interesting predictions of DMSP cycling were enabled by DiTing analysis on the
406 metagenomes analyzed here. Note any predictions made from genetic potential alone
407 require further investigation regarding function and expression and, importantly,
408 substantiation for synthesis and turnover rate analysis.

409 The samples NSu-F2b and NSu-F5 had lower oxygen concentration than Fw-
410 F1b, Fw-F3 and RR-F1b samples, especially NSu-F2b (0.07 and 0.14 mmol l⁻¹ for
411 NSu-F2b and NSu-F5, respectively; 0.17 - 0.2 mmol l⁻¹ for other three). Indeed,
412 compared to the other three samples, NSu-F2b and NSu-F5 had significantly more
413 genes encoding *bd* ubiquinol cytochrome oxidases ($p < 0.01$) that are associated with
414 low oxygen concentrations (Fig. 4). It is worth noting that the *bd* oxidase was
415 enriched most in NSu-F2b under the highest sulfide concentration (1.6 mmol l⁻¹) and
416 lowest oxygen concentration. A previous study found that *bd* oxidase could promote
417 sulfide-resistant O₂ consumption and growth in *E. coli* (Forte et al., 2016), implying
418 the important role of *bd* oxidases in the low oxygen NSu-F2b environment.

419 The NSu-F2b and NSu-F5 samples showed enrichment for denitrification,

420 nitrification and nitrogen fixation potential, which may be due to the lower oxygen
421 levels of these samples or is possibly reflecting the nitrogen availability at higher
422 temperatures. Notably, in NSu-F5, genes encoding for the denitrification enzymes
423 responsible for reduction of the cytotoxic gaseous intermediates, nitric oxide (NO),
424 *norBC*, and nitrous oxide (N₂O), *nosZ*, are significantly enriched, alongside the
425 nitrifying genes responsible for aerobic conversion of nitrite to nitrate (*nxrAB*). The
426 importance of nitrification and denitrification to nitrogen cycling of hydrothermal
427 vents has previously been reported (Bourbonnais et al., 2012), but not at the resolution
428 allowed by DiTing. The transcriptional and enzymatic activity of these systems at
429 these pH levels would certainly need experimental validation. These metagenomes
430 highlight metabolic importance of nitrogen cycling with the potential for all other
431 pathways being at similarly high levels (Supplementary Table S2) in all samples, with
432 the exception of nitrite assimilation (nitrite to ammonia) and hydroxylamine oxidation
433 to nitrite (*hao*) was not detected. Again, this may reflect nitrogen availability but is
434 also indicative of nitrogen source preference of the microbiomes under the highly
435 reactive physicochemical constraints of the vent environment. This study illustrates
436 the need for comprehensive measurements of nitrogen flux, metatranscriptional
437 analyses to ascertain the most active pathways and identification of the dominant
438 organisms responsible for nitrogen cycling in these ecosystems. Overall, these results
439 highlight potential microbial metabolic differences in communities from different
440 hydrothermal samples that likely reflect changes in environmental conditions.

441 DiTing was also applied to analyze 15 metagenomic samples from chlorophyll *a*
442 (Chl*a*) maximum layer in Mediterranean Sea from *Tara* Ocean project. The
443 metagenomic clean reads ranged in size from 1.24 to 52.53 Gbp from each sample.
444 The reads were assembled into 71,183 to 1,601,956 contigs with the total assembly
445 sizes ranging from 0.045 to 1.38 Gbp. A total of 18,431,131 ORFs within these
446 contigs were then predicted. ~24% (1,065,097) ORFs were annotated against KEGG
447 databases and affiliated to 8759 KO entries. The 74 pathways related biogeochemical
448 cycles were found (Supplementary Table S3). Compared to the sample derived

449 hydrothermal vents, the *Chla* maximum layer contains remarkable high relative
450 abundance of photosystem pathway as expected (Supplementary Table S3 and S4).
451 Additionally, eukaryotic DMSP synthesis gene, *DSYB* was detected in 10 out of 15
452 *Chla* maximum samples, which were absent in the hydrothermal vent samples. The
453 relative abundance of *DSYB* was comparable to that of prokaryotic DMSP synthesis
454 gene *dsyB* in *Chla* maximum layers (Supplementary S4), indicating that the DMSP
455 was produced by both prokaryotes and eukaryotes in these environments. For DMSP
456 degradation, in six out of 15 samples, the genetic potential to DMSP demethylation,
457 through the *dmdA* gene, was higher than that for DMSP cleavage (*dddS* and *alma1*)
458 (DMSP demethylation:DMSP cleavage = 1.69:1). This is contrasted with the
459 hydrothermal vent samples. In other nine samples, the potential for DMSP
460 demethylation was comparable to that for DMSP cleavage (DMSP
461 demethylation:DMSP cleavage = 0.82:1). These data support both DMSP
462 demethylation and cleavage being the dominant DMSP catabolic pathways in the
463 *Chla* maximum layer.

464 **Conclusion**

465 In summary, this study developed a pipeline (DiTing) to infer and compare
466 biogeochemical pathways from metagenomic and metatranscriptomic data. DiTing is
467 a portable tool for metagenomic and metatranscriptomic datasets, providing
468 automatic, multi-threaded bioinformatic workflows for data handling, including read
469 assembly, ORF prediction, annotation, and more accurate specific formulas for
470 calculating the relative abundance of biogeochemical pathways. The visualization
471 module is designed to more easily compare functions between samples via graphical
472 outputs. In addition, a verified database was built manually for the annotation of
473 genes involved in the production and cycling of DMSP. As validation of the outputs
474 produced by DiTing, comparisons of the relative abundance of biogeochemical
475 pathways in published metagenomes and metatranscriptomes to those calculated by
476 DiTing were consistent. By applying DiTing to analyze five hydrothermal shotgun
477 metagenomes, we showed that the functional profile could accurately reflect changes

478 in environmental conditions (H₂S and O₂ concentrations). DiTing can be readily
479 applied to metagenomic and/or metatranscriptomic studies, with relatively
480 straightforward user intervention. This bioinformatics framework will facilitate our
481 understanding of spatial and temporal changes in microbiome-mediated
482 biogeochemical cycles.

483

484 **Acknowledgements**

485 The authors thank Saiyi Zhu from Shenzhen Nanshan Foreign Language School for logo
486 designing.

487

488 **Funding**

489 This work was supported by Marine S & T Fund of Shandong Province for Pilot National
490 Laboratory for Marine Science and Technology (Qingdao) (No. 2018SDKJ0406-4), the National
491 Key Research and Development Program of China (No. 2018YFE0124100), and the National
492 Natural Science Foundation of China (Nos. 41730530 and 91751202) in X-HZ's laboratory, and
493 Natural Environmental Research Council standard, UK, standard grants (NE/N002385,
494 NE/P012671, and NE/S001352) in JDT's laboratory.

495

496 *Conflicts of Interest:* none declared.

497

498 **References**

499 Acinas, S.G., Sánchez, P., Salazar, G., Cornejo-Castillo, F.M., Sebastián, M., Logares, R., ... Gasol,
500 J.M. (2019) Metabolic architecture of the deep ocean microbiome. *bioRxiv*, 635-680.

501 Aguiar-Pulido, V., Huang, W., Suarez-Ulloa, V., Cickovski, T., Mathee, K., & Narasimhan, G. (2016)
502 Metagenomics, metatranscriptomics, and metabolomics approaches for microbiome analysis.
503 *Evolutionary Bioinformatics*, 12, 5-16.

504 Alcolombri, U., Ben-Dor, S., Feldmesser, E., Levin, Y., Tawfik, D.S., & Vardi, A. (2015)
505 Identification of the algal dimethyl sulfide-releasing enzyme: a missing link in the marine sulfur
506 cycle. *Science*, 348, 1466-1469.

507 Aramaki, T., Blanc-Mathieu, R., Endo, H., Ohkubo, K., Kanehisa, M., Goto, S., & Ogata, H. (2019)
508 KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive score threshold.
509 *Bioinformatics*, doi: 10.1093/bioinformatics/btz859.

- 510 Boden, R., Borodina, E., Wood, A. P., Kelly, D. P., Murrell, J. C., & Schäfer, H. (2011) Purification
511 and characterization of dimethylsulfide monooxygenase from *Hyphomicrobium sulfonivorans*.
512 *Journal of Bacteriology*, 193, 1250–1258.
- 513 Bokulich, N.A., Subramanian S., Faith J.J., Gevers, D., Gordon J.I., Knight, R., Mills, D.A., &
514 Caporaso, J.G. (2013) Quality-filtering vastly improves diversity estimates from Illumina
515 amplicon sequencing. *Nature Method*, 10, 57-59.
- 516 Bork, C., Schwenn, J.D., & Hell, R. (1998) Isolation and characterization of a gene for assimilatory
517 sulfite reductase from *Arabidopsis thaliana*. *Gene*, 212, 147-153.
- 518 Bourbonnais, A., Lehmann, M.F., Butterfield, D.A., & Juniper, S.K. (2012) Subseafloor nitrogen
519 transformations in diffuse hydrothermal vent fluids of the Juan de Fuca Ridge evidenced by the
520 isotopic composition of nitrate and ammonium. *Geochemistry Geophysics Geosystems*, 13, 1-23.
- 521 Bushnell, B. (2014) BBMap: a fast, accurate, splice-aware aligner. Lawrence Berkeley National Lab.
522 (LBNL), Berkeley, CA (United States). <https://www.osti.gov/biblio/1241166>
- 523 Caspi, R., Foerster, H., Fulcher, C.A., Hopkinson, R., Ingraham, J., Kaipa, P., ... Karp, P.D. (2006)
524 MetaCyc: a multiorganism database of metabolic pathways and enzymes. *Nucleic Acids Research*,
525 34, D511-D516.
- 526 Carrión, O., Larke-Mejía, N.L., Gibson, L., Haque, M.F.U., Ramiro-García, J., McGenity, T.J., &
527 Murrell, J.C. (2018) Gene probing reveals the widespread distribution, diversity and abundance of
528 isoprene-degrading bacteria in the environment. *Microbiome*, 6, 219.
- 529 Carrión, O., McGenity, T. J., & Murrell, J. C. (2020) Molecular ecology of isoprene-degrading
530 bacteria. *Microorganisms*, 8, 967.
- 531 Carrión, O., Pratscher, J., Curson, A. R., Williams, B. T., Rostant, W. G., Murrell, J. C., & Todd, J.D.
532 (2017) Methanethiol-dependent dimethylsulfide production in soil environments. *The ISME*
533 *Journal*, 11, 2379–2390.
- 534 Curson, A.R., Todd, J.D., Sullivan, M.J., & Johnston, A.W.B. (2011) Catabolism of
535 dimethylsulphoniopropionate: microorganisms, enzymes and genes. *Nature Reviews*
536 *Microbiology*, 9, 849–859.
- 537 Curson, A.R., Liu, J., Martínez, A.B., Green, R.T., Chan, Y., Carrión, O., ... Todd, J.D. (2017).
538 Dimethylsulfoniopropionate biosynthesis in marine bacteria and identification of the key gene in
539 this process. *Nature Microbiology*, 2, 17009. doi:10.1038/nmicrobiol.2017.9
- 540 Curson, A.R., Williams, B.T., Pinchbeck, B.J., Sims, L.P., Martínez, A.B., Rivera, P.P.L., ... Todd,
541 J.D. (2018). DSYB catalyses the key step of dimethylsulfoniopropionate biosynthesis in many
542 phytoplankton. *Nature Microbiology*, 3, 430–439. doi:10.1038/s41564-018-0119-5.
- 543 Denger, K., Weiss, M., Felux, A.K., Schneider, A., Mayer, C., Spitteller, D., ... Schleheck, D. (2014)
544 Sulphoglycolysis in *Escherichia coli* K-12 closes a gap in the biogeochemical sulphur cycle.

- 545 *Nature* 507, 114–117.
- 546 Dombrowski, N., Teske, A.P., & Baker, B.J. (2018) Expansive microbial metabolic versatility and
547 biodiversity in dynamic Guaymas Basin hydrothermal sediments. *Nature Communications*, 9, 1-
548 13.
- 549 Eyice, O., Myronova, N., Pol, A., Carrión, O., Todd, J.D., Smith, T.J., ... Schäfer, H. (2018) Bacterial
550 SBP56 identified as a Cu-dependent methanethiol oxidase widely distributed in the biosphere. *The*
551 *ISME Journal*, 12, 145–160.
- 552 Felux, A-K., Spitteller, D., Klebensberger, J., & Schleheck, D. (2015) Entner–Doudoroff pathway for
553 sulfoquinovose degradation in *Pseudomonas putida* SQ1. *Proceedings of the National Academy of*
554 *Sciences*, 112, 4298-4305.
- 555 Fillingame, R.H. (1997) Coupling H⁺ transport and ATP synthesis in F1F0-ATP synthases: glimpses
556 of interacting parts in a dynamic molecular machine. *Journal of Experimental Biology*, 200, 217-
557 224.
- 558 Finn, R.D., Clements, J., & Eddy, S.R. (2011) HMMER web server: interactive sequence similarity
559 searching. *Nucleic Acids Research*, 39, D29-D37.
- 560 Finn, R.D., Mistry, J., Tate, J., Coghill, P., Heger, A., Pollington, J.E., ... Bateman, A. (2014) Pfam:
561 the protein families database. *Nucleic Acids Research*, 42, D222-D230.
- 562 Firtz, A., Hofmann, P., Majda, S., Dahms, E., Dröge, J., Fiedler, J., ... McHardy, A.C. (2019).
563 CAMISIM: simulating metagenomes and microbial communities. *Microbiome*, 7, 17.
- 564 Forouzan, E., Shariati, P., Maleki, M.S.M., Karkhane, A.A., & Yakhchali, B. (2018) Practical
565 evaluation of 11 de novo assemblers in metagenome assembly. *Journal of Microbiological*
566 *Methods*, 151, 99-105.
- 567 Fortunato, C.S., Larson, B., Butterfield, D.A., & Huber, J.A. (2018) Spatially distinct, temporally
568 stable microbial populations mediate biogeochemical cycling at and below the seafloor in
569 hydrothermal vent fluids. *Environmental Microbiology*, 20, 769-784.
- 570 Friedrich, C.G. (1998) Physiology and genetics of sulfur-oxidizing bacteria. *Advances in Microbial*
571 *Physiology*, 39, 235-289.
- 572 Frommeyer, B., Fiedler, A.W., Oehler, S.R., Hanson, B.T., Loy, A., Franchini, P., ... Schleheck, D.
573 (2020) Environmental and intestinal phylum Firmicutes bacteria metabolize the plant sugar
574 sulfoquinovose via a 6-deoxy-6-sulfofructose transaldolase pathway. *iScience*, 23, 101510.
- 575 Ganesh, S., Parris, D.J., DeLong, E.F., & Stewart, F.J. (2014) Metagenomic analysis of size-
576 fractionated picoplankton in a marine oxygen minimum zone. *The ISME Journal*, 8, 187-211.
- 577 Gisselmann, G., Klausmeier, P., & Schwenn, J.D. (1993) The ferredoxin: sulphite reductase gene from
578 *Synechococcus* PCC7942. *Biochimica et Biophysica Acta (BBA)-Bioenergetics*, 1144, 102-106.

579 Graham, E.D., Heidelberg, J.F., & Tully, B.J. (2018) Potential for primary productivity in a globally-
580 distributed bacterial phototroph. *The ISME Journal*, 350, 1–6.

581 Howard, E.C., Henriksen, J.R., Buchan, A., Reisch, C.R., Bürgmann, H., Welsh, R., ... Moran, M.A.
582 (2006). Bacterial taxa that limit sulfur flux from the ocean. *Science*, 314, 649–652.
583 doi:10.1126/science.1130657

584 Huertacepas, J., Szklarczyk, D., Forslund, K., Cook, H., Heller, D., Walter, M.C., ... Bork, P. (2016)
585 eggNOG 4.5: a hierarchical orthology framework with improved functional annotations for
586 eukaryotic, prokaryotic and viral sequences. *Nucleic Acids Research*, 44, D286-D293.

587 Hyatt, D., Chen, G., Locascio, P.F., Land, M.L., Larimer, F.W., & Hauser, L.J. (2010) Prodigal:
588 prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics*, 11,
589 119.

590 Johnston, A.W., Green, R.T., & Todd, J.D. (2016) Enzymatic breakage of dimethylsulfoniopropionate
591 —a signature molecule for life at sea. *Current Opinion Chemical Biology*, 31, 58–65.

592 Jones, M.R., & Good, J.M. (2016) Targeted capture in evolutionary and ecological genomics.
593 *Molecular. Ecology*, 25, 185-202.

594 Li, D., Luo R., Liu, C-M., Leung, C-M, Ting, H-F., Sadakane, K., ... Lam, T-W. (2016) MEGAHIT
595 v1.0: A fast and scalable metagenome assembler driven by advanced methodologies and
596 community practices. *Methods*, 102, 3-11.

597 Li, H. (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv*
598 *preprint arXiv:1303.3997*.

599 Liao, C., & Seebeck, F.P. (2019) In vitro reconstitution of bacterial DMSP biosynthesis. *Angewandte*
600 *Chemie International Edition*, 58, 3553-3556.

601 Lidbury, I., Kröber, E., Zhang, Z., Zhu, Y., Murrell, J. C., Chen, Y., & Schäfer, H. (2016) A
602 mechanism for bacterial transformations of DMS to DMSO: a missing link in the marine organic
603 sulfur cycle. *Environmental Microbiology*. 18, 2754–2765.

604 Llorens-Marès, T., Yooseph, S., Goll, J., Hoffman, J., Vila-Costa, M., Borrego, C.M., ... Casamayor,
605 E. O. (2015) Connecting biodiversity and potential functional role in modern euxinic
606 environments by microbial metagenomics. *The ISME Journal*, 9, 1648-1661.

607 McDevitt, C.A., Hanson, G.R., Noble, C.J., Cheesman, M.R., & McEwan, A.G. (2002)
608 Characterization of the redox centers in dimethyl sulfide dehydrogenase from *Rhodovulum*
609 *sulfidophilum*. *Biochemistry*, 41, 15234–15244.

610 Meier, D.V., Pjevac, P., Bach, W., Hourdez, S., Girguis, P.R., Vidoudez, C., ... Meyerdierks, A. (2017)
611 Niche partitioning of diverse sulfur-oxidizing bacteria at hydrothermal vents. *The ISME Journal*,
612 11, 1545-1558.

- 613 Meier, D.V., Pjevac, P., Bach, W., Markert, S., Schweder, T., Jamieson, J., ... Meyerdierks, A. (2019)
614 Microbial metal-sulfide oxidation in inactive hydrothermal vent chimneys suggested by
615 metagenomic and metaproteomic analyses. *Environmental Microbiology*, 21, 682-701.
- 616 Nurk, S., Meleshko, D., Korobeynikov, A., & Pevzner, P. (2017) metaSPAdes: a new versatile
617 metagenomic assembler. *Genome Research*, 27, 824-834.
- 618 Ogata, H., & Goto, S. (2000) KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids*
619 *Research*, 27, 29-34.
- 620 Ostrowski, J., Barber, M.J., Rueger, D.C., Miller, B.E., Siegel, L.M., & Kredich, N.M. (1989a)
621 Characterization of the flavoprotein moieties of NADPH-sulfite reductase from *Salmonella*
622 *typhimurium* and *Escherichia coli*. Physicochemical and catalytic properties, amino acid sequence
623 deduced from DNA sequence of *cysJ*, and comparison with NADPH-cytochrome P-450 reductase.
624 *Journal of Biological Chemistry*, 264, 15796-15808.
- 625 Ostrowski, J., Wu, J-Y., Rueger, D.C., Miller, B.E., Siegel, L.M., & Kredich, N.M. (1989b)
626 Characterization of the *cysJIIH* regions of *Salmonella typhimurium* and *Escherichia coli* B. DNA
627 sequences of *cysI* and *cysH* and a model for the siroheme-Fe₄S₄ active center of sulfite reductase
628 hemoprotein based on amino acid homology with spinach nitrite reductase. *Journal of Biological*
629 *Chemistry*, 264, 15726-15737.
- 630 Parks, D.H., Rinke, C., Chuvochina, M., Chaumeil, P-A., Woodcroft, B.J., Evans, P.N., ... Tyson, G.W.
631 (2017) Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree
632 of life. *Nature Microbiology*, 2, 1533-1542.
- 633 Pasolli, E., Asnicar, F., Manara, S., Zolfo, M., Karcher, N., Armanini, F., ... Segata, N. (2019)
634 Extensive unexplored human microbiome diversity revealed by over 150,000 genomes from
635 metagenomes spanning age, geography, and lifestyle. *Cell*, 176, 649-662.
- 636 Petter, T., Lundin, D., Plathan, J., Poole, A.M., Sjöberg, B-M., & Sjöling, S. (2013) A metagenomics
637 transect into the deepest point of the Baltic Sea reveals clear stratification of microbial functional
638 capacities. *Plos One*, 8, e74983.
- 639 Pisa, K.Y., Huber, H., Thomm, M., & Muller, V. (2009) A sodium ion-dependent A1AO ATP synthase
640 from the hyperthermophilic archaeon *Pyrococcus furiosus*. *The FEBS Journal*, 274, 3928-3938.
- 641 Prestat, E., David, M.M., Hultman, J., Taş N., Lamendella, R., Dvornik, J., ... Jansson, J.K. (2014)
642 FOAM (Functional Ontology Assignments for Metagenomes): a Hidden Markov Model (HMM)
643 database with environmental focus. *Nucleic Acids Research*, 42, e145.
- 644 Reisch, C.R., Moran, M.A., & Whitman, W.B. (2011) Bacterial catabolism of
645 dimethylsulfoniopropionate (DMSP). *Frontier in Microbiology*, 2, 172.
- 646 Reisch, C.R., Stoudemayer, M.J., Varaljay, V.A., Amster, I.J., Moran, M.A., & Whitman, W.B. (2011)
647 Novel pathway for assimilation of dimethylsulphoniopropionate widespread in marine bacteria.

648 *Nature*, 473, 208–211.

649 Riesenfeld, C.S., Schloss, P.D., & Handelsman, J. (2004) Metagenomics: genomic analysis of
650 microbial communities. *Annual Reviews of Genetics*, 38, 525-552.

651 Ross, O., Olson, R., Pusch, G.D., Olsen, G.J., Davis, J.J., Disz, T., ... Stevens, R. (2014) The SEED
652 and the Rapid Annotation of microbial genomes using Subsystems Technology (RAST). *Nucleic
653 Acids Research*, 42, D206-D214.

654 Selengut, J.D., Haft, D.H., Davidsen, T., Ganapathy, A., Gwinn-Giglio, M., Nelson, W.C., ... White,
655 O. (2007) TIGRFAMs and Genome Properties: tools for the assignment of molecular function and
656 biological process in prokaryotic genomes. *Nucleic Acids Research*, 35, D260-D264.

657 Shakya, M., Lo, C-C., & Chain, P.S.G. (2019) Advances and challenges in metatranscriptomic
658 analysis. *Frontier in Microbiology*, 10, 904.

659 Shao, X., Cao, H-Y., Zhao, F., Ming, P., Wang, P., and Li, C-Y., ... Zhang, Y-Z. (2019) Mechanistic
660 insight into 3-methylmercaptopropionate metabolism and kinetical regulation of demethylation
661 pathway in marine dimethylsulfoniopropionate-catabolizing bacteria. *Molecular Microbiology*,
662 111, 1057–1073.

663 Sharpton, T.J. (2014) An introduction to the analysis of shotgun metagenomic data. *Frontiers in Plant
664 Science*, 5, 00209.

665 Smedile, F., Messina, E., Cono, V.L., Tsoy, O., Monticelli, L.S., ... Yakimov, M.M. (2013)
666 Metagenomic analysis of hadopelagic microbial assemblages thriving at the deepest part of
667 Mediterranean Sea, Matapan-Vavilov Deep. *Environmental Microbiology*, 15, 167-182.

668 Song, D., Zhang, Y., Liu, J., Zhong, H., Zheng, Y., Zhou, S., ... Zhang, X-H. (2020) Metagenomic
669 insights into the cycling of dimethylsulfoniopropionate and related molecules in the Eastern China
670 Marginal seas. *Frontiers in Microbiology*, 11, 157.

671 Sun, J., Todd, J.D., Thrash, J.C., Qian, Y., Qian, M.C., Temperton, B., ... Giovannoni, S.J. (2016) The
672 abundant marine bacterium *Pelagibacter* simultaneously catabolizes dimethylsulfoniopropionate
673 to the gases dimethyl sulfide and methanethiol. *Nature Microbiology*, 1, e16065.

674 Tatusov, R.L., Galperin, M.Y., Natale, D.A., & Koonin, E.V. (2000) The COG database: a tool for
675 genome-scale analysis of protein functions and evolution. *Nucleic Acids Research*, 28, 33-36.

676 Williams, B.T., Cowles, K., Martínez, A.B., Curson, A.R.J., Zheng Y., Liu, J., ... Todd, J.D. (2019)
677 Bacteria are important dimethylsulfoniopropionate producers in coastal sediments. *Nature
678 Microbiology*, 4, 1815-1825.

679 Xue, C-X., Liu, J., Lea-Smith, D.J., Rowley, G., Lin, H., Zheng, Y., ... Zhang, X-H. (2020a) Insights
680 into the vertical stratification of microbial ecological roles across the deepest seawater column on
681 Earth. *Microorganisms*, 8, 1309.

682 Xue, C-X. Zhang, H., Lin, H., Sun, Y., Luo, D., Huang, Y., ... Luo, H. (2020b) Ancestral niche
683 separation and evolutionary rate differentiation between sister marine Flavobacteria lineages.
684 *Environmental Microbiology*, 22, 3234-3247.

685 Zeghouf, M., Fontecave, M., & Coves, J. (2000) A simplified functional version of the *Escherichia coli*
686 sulfite reductase. *Journal of Biological Chemistry*, 275, 37651-37656.

687 Zhang, X., Xu, W., Liu, Y., Cai, M., Luo, Z., & Li, M. (2018) Metagenomics reveals microbial
688 diversity and metabolic potentials of seawater and surface sediment from a hadal biosphere at the
689 Yap Trench. *Frontiers in Microbiology*, 9, 2402.

690 Zhang, X.H., Liu, J., Liu, J.L., Yang, G., Xue, C-X., Curson, A.R.J., & Todd, J.D. (2019). Biogenic
691 production of DMSP and its degradation to DMS—their roles in the global sulfur cycle. *Science*
692 *China Life Science*, 62, 1296–1319.

693

694 **Fig. 1.** A flowchart of the major steps involved in running DiTing. First (A), clean
695 reads of metagenomes or/metatranscriptomes are assembled, annotated and mapped.
696 Second (B), a table for relative abundances of KO number in KEGG among samples
697 is constructed and relative abundances of biogeochemical pathways are estimated
698 according to unbiased specific formulas. Third (C), heatmap and sketch plots are
699 drawn to aid visualization.

700

701 **Fig. 2.** Bubble plots depicting the DiTing result of the relative abundance of pathways
702 in simulated metagenomes. Photoautotroph group contains sample1-5 that simulated
703 by *Cyanobacteria* genomes. Chemoautotroph group contains sample6-10 that
704 simulated by ammonia-oxidizing archaea genomes. Heterotroph group contains
705 sample11-15 that simulated by SAR11 genomes.

706

707 **Fig. 3.** Pie charts representing the relative abundance of carbon (A), nitrogen (B),
708 sulfur (C) and DMSP (D) cycle related pathways for five metagenomic samples from
709 the Manus Basin. Normalized relative abundance was calculated through dividing the
710 relative abundance of a pathway in an individual sample by the sum of this pathway's
711 relative abundance in all samples. Pie chart area reflects the relative abundance of the

712 process according to the scale shown in pink. The dashed line in panel D means the
 713 data was not shown. (A) CBB, Calvin-Benson-Bassham cycle; rTCA, reductive citric
 714 acid cycle; WL, Wood-Ljungdahl pathway; 3HB, 3-hydroxypropionate bicycle. (B)
 715 ANRA, assimilatory nitrate reduction to ammonia; DNRA, Dissimilatory nitrate
 716 reduction to ammonia; Anammox, anaerobic ammonia oxidation. (C) ASR,
 717 assimilatory sulfate reduction; DSR, dissimilatory sulfate reduction. (D) DMSP,
 718 dimethylsulfoniopropionate; MMPA, methylmecaptopropionate; MeSH,
 719 methanethiol; DMSO, dimethylsulfoxide; *L*-Met, *L*-methionine.

720

721 **Fig. 4.** Bubble plots depicting the relative abundance of pathways for carbon (A),
 722 sulfur (B), nitrogen (C) and other selected (D) processes. The key marker genes used
 723 to report on the genetic potential for pathways (as the relative abundances) are
 724 indicated in brackets. ASR, assimilatory sulfate reduction; DSR, dissimilatory sulfate
 725 reduction. The full name of these key marker genes can be found in Supplementary
 726 Table S1. For better visualization, we multiply the relative abundance by 10^{-3} and
 727 transformed by $\log(10)$.

728

729

730 **Table 1:** A summary of sampling sites and environmental parameters for collected
 731 samples

Sample name	Sample type	Latitude	Longitude	Depth [m]	T [°C]	pH	H ₂ S [mM]	CH ₄ [mM]	DIC [mM]	O ₂ [mM]
NSu-F2b	water/fluid	S 03°47.995'	E 152°06.052'	1155	51.7	4.3	1.61	0.2	3.07	0.07
NSu-F5	water/fluid	S 03°47.955'	E 152°06.080'	1199	31.4	5.1	0.7	0.01	0.18	0.14
Fw-F1b	water/fluid	S 03°43.700'	E 151°40.344'	1709	3.7	6.5	0	0	0.24	0.17
Fw-F3	water/fluid	S 03°43.698'	E 151°40.350'	1705	3.2	7.2	ND	ND	ND	ND
RR-F1b	water/fluid	S 03°43.238'	E 151°40.519'	1685	6.6	7.5	0	0	2.34	0.2

732 ND – 'not determined'. 0 – below detection limit

Table 2 The relative abundance of biogeochemical pathways in metagenomes from the Manus Basin

Pathway	NSu-F2b	NSu-F5	Fw-F1b	Fw-F3	RR-F1b
Photosystem II (psbABCDEF)	0.0687	2.83	0.0315	0	0.264
Photosystem I (psaABCDEF)	0	0.105	0	0	0
Cytochrome b6/f complex (petABCDGLMN)	0.95	0.728	0.448	0.456	1.15
Anoxygenic photosystem II (pufML)	0	0	0	0	0
Anoxygenic photosystem I (pscABCD)	0	0	0	0	0
RuBisCo	13	34.1	31.3	40.9	57.9
CBB cycle (prkB)	12.6	74.8	45.7	55	55.9
rTCA cycle (aclAB, ccsAB, ccl)	74.4	53.7	4.2	2.09	0.871
Wood-Ljungdahl pathway (acsABCDE)	15.5	2.32	0	0	0
3-Hydroxypropionate Bicycle	2.02	2	0.39	0.335	0.661
Glycolysis (glk, pfk, pyk)	123	158	49.8	64.2	94.7
Entner-Doudoroff pathway, glucose-6P -> glyceraldehyde-3P + pyruvate	19.6	31.4	3.39	3.92	8.68
Gluconeogenesis (fbp, pck)	383	281	66.4	75	103
TCA cycle	178	184	38.5	45	76.5
Methanogenesis (mcrABG)	0	0	0	0	0
Methanogenesis, methanol -> methane (mtaABC)	0	0	0	0	0
Methanogenesis, amines -> methane (mtbA, mtmC, mtbC, mttC)	0	0	0	0	0
Methanogenesis, acetate -> methane (cdhCDE)	2.63	0.565	0	0	0
Methanogenesis, CO ₂ -> methane	4.49	1.74	0.438	0.646	2.65
Methane oxidation, methane -> methanol (mmoBCDXYZ, amoABC)	22	6.15	7.86	6.53	5.7
Methane oxidation, methanol -> formaldehyde (mxoFI, xoxF)	0.101	0	0	0	0
Fermentation to lactate, pyruvate -> lactate (LDH)	4.96	0.19	0	0	0
Fermentation to formate, pyruvate -> formate (pflD)	0.563	3.88	0	0	0
Fermentation to formate -> CO ₂ & H ₂ (fdh)	14.7	14.4	2.88	2.87	4.64
Fermentation to acetate, pyruvate -> acetate (poxB, poxL, acyP)	76.8	40.7	10.2	13.5	21.1
Fermentation to acetate, acetyl-CoA -> acetate (ach1, eutD, pta, acyP)	83.5	80.5	11.3	13.9	21.4
Fermentation to acetate, lactate -> acetate (EC:1.13.12.4)	0	0	0	0	0
Fermentation to ethanol, acetate to acetaldehyde (ald)	16	39	7.52	9.45	6.19

Fermentation to ethanol, acetyl-CoA to acetaldehyde (reversible)	2.01	10.1	0.409	0.453	0.403
Fermentation to ethanol, acetaldehyde to ethanol (adh, mdh)	40.1	69.4	11.8	17.2	16.5
Fermentation to succinate	245	216	32.8	38.8	62.9
Anaplerotic genes (pyruvate -> oxaloacetate)	699	627	88.6	107	137
Dissimilatory nitrate reduction, nitrate -> nitrite (narGHI or napAB)	187	160	28.1	33.4	33.3
Dissimilatory nitrate reduction, nitrite -> ammonia (nirBD or nrfAH)	19.7	86.9	31.6	36.5	65.1
Assimilatory nitrate reduction, nitrate -> nitrite (narB or NR or nasAB)	4.56	7.21	0	0	0.132
Assimilatory nitrate reduction, nitrite -> ammonia (NIT-6 or nirA)	0	0	0	0	0.911
Denitrification, nitrite -> nitric oxide (nirK or nirS)	9.05	70.9	9.87	2.98	4.02
Denitrification, nitric oxide -> nitrous oxide (norBC)	68	338	31	34.7	8.62
Denitrification, nitrous oxide -> nitrogen (nosZ)	31.4	111	11.8	9.11	1.9
Nitrogen fixation, nitrogen -> ammonia (nifKDH)	0.981	2.12	0	0	0
Nitrification, ammonia -> hydroxylamine (amoABC)	22	6.15	7.86	6.53	5.7
Nitrification, hydroxylamine -> nitrite (hao)	0	0	0	0	0
Nitrification, nitrite -> nitrate (nxrAB)	62.9	23.3	4.29	4.79	4.88
Anammox, nitric oxide + ammonia -> hydrazine (hzs)	0.976	5.89	0.286	0.131	1.35
Anammox, hydrazine -> nitrogen (hdh)	0.219	0.254	0	0	0
Assimilatory sulfate reduction, sulfate -> sulfite	89.3	91.9	12.9	15.7	23.6
Assimilatory sulfate reduction, sulfite -> sulfide (cysJI or sir)	10.2	17.5	0.353	0.187	1.81
Dissimilatory sulfate reduction, sulfate -> sulfite (reversible) (sat and aprAB)	103	134	53.7	64.8	92.8
Dissimilatory sulfate reduction, sulfite -> sulfide (reversible) (dsrAB)	6.34	82.6	73.7	83	71.7
Thiosulfate oxidation by SOX complex, thiosulfate -> sulfate	20.3	183	77.5	90.6	147
Alternative thiosulfate oxidation (doxAD)	2.5	3.8	0.871	0.784	1.49
Alternative thiosulfate oxidation (tsdA)	18.8	39.8	1.79	1.05	0.667
Sulfur reduction, sulfur -> sulfide (sreABC)	0	0	0	0	0
Thiosulfate disproportionation, thiosulfate -> sulfide & sulfite (phsABC)	31.9	16.2	1.48	0.818	0.0703
Sulfhydrogenase, (sulfide) _n -> (sulfide) _{n-1}	6.64	0.367	0	0	0
Sulfur disproportionation, sulfur -> sulfide & sulfite	0	0	0	0	0
Sulfur dioxygenase	15.9	65.7	47.1	44.3	84.6
Sulfite oxidation, sulfite -> sulfate (sorB, SUOX, soeABC)	64.9	287	56	64.9	80.9
Sulfide oxidation, sulfide -> sulfur (fccAB)	3.05	22.5	15.5	18.6	14.6

DMSP biosynthesis, Met -> DMSP (DSYB or dsyB or mmtN)	5.16	0.286	0	0	0.168
DMSP demethylation, DMSP -> MMPA (dmdA)	9.31	5.7	2.44	0.933	5.57
DMSP demethylation, MMPA -> MeSH (dmdBCD or acuH)	85.2	102	13.1	14	26.8
DMSP cleavage, DMSP -> DMS (dddS or alma1)	25.3	13.9	7.74	8.96	15
DMS oxidation, DMS -> MeSH (dmoA)	8.51	16.4	2.96	1.34	6.59
DMS oxidation, DMS -> DMSO (ddhABC or tmm)	75	53.1	8.73	10.6	13.7
DMSO reduction, DMSO -> DMS (dms or dorA)	32.8	87.4	26.1	43.6	39.2
MddA pathway, MeSH -> DMS (mddA)	1.33	31.1	0.488	0.505	0.501
MeSH oxidation, MeSH -> Formaldehyde (MTO)	0.5	2.55	0	0	0.49
F-type ATPase	232	248	59.3	72.9	113
V/A-type ATPase	23.2	11.1	5.85	5.8	7.34
NADH-quinone oxidoreductase	89.7	173	52.5	57.4	96.3
NAD(P)H-quinone oxidoreductase	0.0815	0.232	0.103	0.17	0.151
Succinate dehydrogenase (ubiquinone)	0	0	0	0	0
Cytochrome c oxidase, cbb3-type	34.8	147	34.3	40.3	74
Cytochrome <i>bd</i> ubiquinol oxidase	239	128	5.67	3.47	5.09
Cytochrome <i>o</i> ubiquinol oxidase	11.2	13.5	0.188	0.184	0.944
Cytochrome c oxidase, prokaryotes, aa3-type	39.5	106	42	54.9	104
Cytochrome aa3-600 menaquinol oxidase	0	0	0	0	0
Cytochrome bc1 complex	13.8	35.3	11.7	14.3	21.8
Type I Secretion	7.82	18.8	1.29	1.09	2.93
Type III Secretion	0.0069	0.00643	0.19	0.118	0.0265
Type II Secretion	40.3	57.8	9.38	10.3	9.41
Type IV Secretion	9.82	10.9	0.252	0.218	0.172
Type VI Secretion	3.81	23.7	1.71	1.29	1.3
Sec-SRP	196	200	50.6	56.9	90.5
Twin arginine targeting	183	199	49.7	59.7	83
Type Vabc secretion	0	0	0	0	0
Bacterial chemotaxis	119	70.5	5.09	3.93	8.31
Flagellum assembly	112	51.9	5.06	4.64	10.4
Dissimilatory arsenic reduction	105	181	24.5	26.4	37

735

736