

1 **METABOLIC: High-throughput profiling of microbial genomes for functional**
2 **traits, metabolism, biogeochemistry, and community-scale functional networks**

3

4

5 Zhichao Zhou¹, Patricia Q. Tran^{1,2}, Adam M. Breister¹, Yang Liu³, Kristopher Kieft^{1,4}, Elise S.
6 Cowley^{1,4}, Ulas Karaoz⁵, Karthik Anantharaman^{1,*}

7

8

9 ¹Department of Bacteriology, University of Wisconsin-Madison, Madison, WI, 53706, USA

10 ²Department of Integrative Biology, University of Wisconsin-Madison, Madison, WI, 53706,
11 USA

12 ³Institute for Advanced Study, Shenzhen University, Shenzhen, Guangdong Province, 518060,
13 China

14 ⁴Microbiology Doctoral Training Program, University of Wisconsin-Madison, Madison, WI,
15 53706, USA

16 ⁵Earth and Environmental Sciences, Lawrence Berkeley National Laboratory, Berkeley, CA,
17 94720, USA

18

19

20

21 *Correspondence to Karthik Anantharaman, karthik@bact.wisc.edu

22 **ABSTRACT**

23 **Background:** Advances in microbiome science are being driven in large part due to our ability
24 to study and infer microbial ecology from genomes reconstructed from mixed microbial
25 communities using metagenomics and single-cell genomics. Such omics-based techniques
26 allow us to read genomic blueprints of microorganisms, decipher their functional capacities
27 and activities, and reconstruct their roles in biogeochemical processes. Currently available
28 tools for analyses of genomic data can annotate and depict metabolic functions to some extent,
29 however, no standardized approaches are currently available for the comprehensive
30 characterization of metabolic predictions, metabolite exchanges, microbial interactions, and
31 microbial contributions to biogeochemical cycling.

32
33 **Results:** We present METABOLIC (METabolic And BiogeOchemistry analYses In
34 miCobes), a scalable software to advance microbial ecology and biogeochemistry studies
35 using genomes at the resolution of individual organisms and/or microbial communities. The
36 genome-scale workflow includes annotation of microbial genomes, motif validation of
37 biochemically validated conserved protein residues, metabolic pathway analyses, and
38 calculation of contributions to individual biogeochemical transformations and cycles. The
39 community-scale workflow supplements genome-scale analyses with determination of
40 genome abundance in the microbiome, potential microbial metabolic handoffs and metabolite
41 exchange, reconstruction of functional networks, and determination of microbial contributions
42 to biogeochemical cycles. METABOLIC can take input genomes from isolates, metagenome-
43 assembled genomes, or single-cell genomes. Results are presented in the form of tables for
44 metabolism and a variety of visualizations including biogeochemical cycling potential,
45 representation of sequential metabolic transformations, community-scale microbial functional
46 networks using a newly defined metric ‘MW-score’ (metabolic weight score), and metabolic
47 Sankey diagrams. METABOLIC takes ~3 hours with 40 CPU threads to process ~100
48 genomes and corresponding metagenomic reads within which the most compute-demanding
49 part of hmmsearch takes ~45 mins, while it takes ~5 hours to complete hmmsearch for ~3600
50 genomes. Tests of accuracy, robustness, and consistency suggest METABOLIC provides
51 better performance compared to other software and online servers. To highlight the utility and
52 versatility of METABOLIC, we demonstrate its capabilities on diverse metagenomic datasets
53 from the marine subsurface, terrestrial subsurface, meadow soil, deep sea, freshwater lakes,
54 wastewater, and the human gut.

55
56 **Conclusion:** METABOLIC enables the consistent and reproducible study of microbial
57 community ecology and biogeochemistry using a foundation of genome-informed microbial
58 metabolism, and will advance the integration of uncultivated organisms into metabolic and
59 biogeochemical models. METABOLIC is written in Perl and R and is freely available at
60 <https://github.com/AnantharamanLab/METABOLIC> under GPLv3.

61
62 **Keywords:** functional traits, metagenome-assembled genomes, microbiome,
63 biogeochemistry, metabolic potential, microbial functional networks.

64 **INTRODUCTION**

65 Metagenomics and single-cell genomics have transformed the field of microbial ecology by
66 revealing a rich diversity of microorganisms from diverse settings, including terrestrial [1-3]
67 and marine environments [4, 5] and the human body [6]. These approaches can provide an
68 unbiased and insightful view into microorganisms mediating and contributing to
69 biogeochemical activities at a number of scales ranging from individual organisms to
70 communities [7-9]. Recent studies have also enabled the recovery of hundreds to thousands of
71 genomes from a single sample or environment [8, 10, 11]. However, analyses of ever-increasing
72 datasets remain a challenge. For example, there is a lack of scalable and reproducible
73 bioinformatic approaches for characterizing metabolism and biogeochemistry, as well as
74 standardizing their analyses and representation for large datasets.

75
76 Microbially-mediated biogeochemical processes serve as important driving forces for the
77 transformation and cycling of elements, energy, and matter among the lithosphere, atmosphere,
78 hydrosphere, and biosphere [12]. Microbial communities in natural environmental settings exist
79 in the form of complex and highly connected networks that share and compete for metabolites
80 [13-15]. The interdependent and cross-linked metabolic and biogeochemical interactions within
81 a community can provide a relatively high level of plasticity and flexibility [16]. For instance,
82 multiple metabolic steps within a specific pathway are often separately distributed in a number
83 of microorganisms and they are interdependent on utilizing the substrates from the previous
84 step [2, 17, 18]. This scenario, referred to as ‘metabolic handoffs’, is based on sequential
85 metabolic transformations, and provides the benefit of high resilience of metabolic activities
86 which make both the community and function stable in the face of perturbations [17, 18]. It is
87 therefore highly valuable to obtain the information of microbial metabolic function from the
88 perspective of individual genomes as well as the entire microbial community.

89
90 Currently, there are many quantitative software and platforms for reconstructing species and
91 community-level metabolic networks [19-25]. They are largely based on building microbial
92 metabolic models containing reactions for substrate utilization and product generation [15, 19].
93 Based on individual microbial models, metabolic phenotypes for the whole community can be
94 further predicted [15]. These approaches allow providing mechanistic bases for predicting and
95 thus operating community metabolisms based on the given environmental conditions and
96 predicted microbial phenotypes [26]. Thus they are more focused on illustrating the operating
97 principles of community metabolisms and the underlying metabolic networks of connected
98 reactions to achieve better outcomes for metabolite production [21, 22], industrial applications
99 [19], drug discovery [19], etc.

100
101 Yet, seldom have approaches been developed to study the functional role of microorganisms in
102 the context of biogeochemistry and community-level functional networks [27, 28]. Such tools
103 are based on the principles of facilitating the understanding of microbially-mediated
104 biogeochemical activities. The tools ask for identifying and providing metabolic predictions on
105 the functional details, transformations of nutrients and energy, and functional connections for

106 microorganisms within the community [29]. The resulting genome-informed microbial
107 metabolisms are important for understanding the microbial roles within a whole community in
108 mediating the biogeochemical processes. Currently, such quantitative approaches to interpret
109 functional details, reconstruct metabolic relationships, and visualize microbial functional
110 networks are still limited [27, 28].

111

112 Prediction of microbial metabolism relies on the annotation of protein function for
113 microorganisms using a number of established databases, e.g., KEGG [30], MetaCyc [31],
114 Pfam [32], TIGRfam [33], SEED/RAST [34], and eggNOG [35]. However, these results are
115 often highly detailed, and therefore can be overwhelming to users. Obtaining a functional
116 profile and identifying metabolic pathways in a microbial genome can involve manual
117 inspection of thousands of genes [36]. Organizing, interpreting, and visualizing such datasets
118 remains a challenge and is often untenable especially with datasets larger than one microbial
119 genome. There is a critical need for approaches and tools to identify and validate the presence
120 of metabolic pathways, biogeochemical function, and connections in microbial communities in
121 a user-friendly manner. Such tools addressing this gap would also allow standardization of
122 methods and easier integration of genome-informed metabolism into biogeochemical models,
123 which currently rely primarily on physicochemical data and treat microorganisms as black
124 boxes [37]. A recent statistical study indicates that incorporating microbial community structure
125 in biogeochemical modeling could significantly increase model accuracy of processes that are
126 mediated by narrow phylogenetic guilds via functional gene data, and processes that are
127 mediated by facultative microorganisms via community diversity metrics [38]. This highlights
128 the importance of integrating microbial community and genomic information into the
129 prediction and modeling of biogeochemical processes.

130

131 Here we present the software METABOLIC (METabolic And BiogeOchemistry anaLyses In
132 miCrobes), a toolkit to profile metabolic and biogeochemical traits, and functional networks in
133 microbial communities based on microbial genomes. METABOLIC integrates annotation of
134 proteins using KEGG [30], TIGRfam [33], Pfam [32], custom hidden Markov model (HMM)
135 databases [2], dbCAN2 [39], and MEROPS [40], incorporates a protein motif validation step
136 to accurately identify proteins based on prior biochemical validation, and determines presence
137 or absence of metabolic pathways based on KEGG modules. METABOLIC also produces user-
138 friendly outputs in the form of tables and figures including a summary of microbial functional
139 profiles, biogeochemically-relevant pathways, functional networks at the scale of individual
140 genomes and community levels, and microbial contribution to the biogeochemical processes.

141

142 METHODS

143 HMM databases used by METABOLIC

144 To generate a broad range of metabolic gene HMM profiles, we integrated three sets of HMM-
145 based databases, which are KOfam [41] (July 2019 release, containing HMM profiles for
146 KEGG/KO with predefined score thresholds), TIGRfam [33] (Release 15.0), Pfam [32]
147 (Release 32.0), and custom metabolic HMM profiles [2]. In order to achieve a better HMM

148 search result excluding non-specific hits, we have tested and manually curated cutoffs for those
149 HMM databases listed above into the resulting HMMs: KOfam database - KOfam suggested
150 values; TIGRfam/Pfam/Custom databases - manually curated by adjusting noise cutoffs (NC)
151 or trusted cutoffs (TC) to avoid potential false positive hits. For the KOfam suggested cutoffs,
152 we considered both the score type (full length or domain) and the score value to assign whether
153 an individual protein hit is significant or not. HMM databases were used as the reference for
154 hmmsearch [42] to find protein hits of input genomes. Prodigal [43] was used to annotate
155 genomic sequences (the method used to find ORFs by Prodigal can be set by METABOLIC as
156 “meta” or “single”), or a user can provide self-annotated proteins (with extensions of “.faa”) to
157 facilitate incorporation into existing pipelines. Methods on the manual curation of these HMM
158 databases are described in the next section.

159

160 **Curation of cutoff scores for metabolic HMMs**

161 Two curation methods for adjusting NC or TC of TIGRfam/Pfam/Custom databases were used
162 for a specific HMM profile. First, we parsed and downloaded representative protein sequences
163 according to either the corresponding KEGG identifier or UniProt identifier [44]. We then
164 randomly subsampled a small portion of the sequences (10% of the whole collection if this was
165 more than 10 sequences, or at least 10 sequences) as the query to search against the
166 representative protein collections [42]. Subsequently, we obtained a collection of hmmsearch
167 scores by pairwise sequence comparisons. We plotted scores against hmmsearch hits and
168 selected the mean value of the sharpest decreasing interval as the adjusted cutoff
169 (approximately the F1 score). Second, we downloaded a collection of proteins that belong to a
170 specific HMM profile and pre-checked the quality and phylogeny of these proteins by
171 reconstructing and manually inspecting phylogenetic trees. We applied pre-checked protein
172 sequences as the query search against a set of training metagenomes (data not shown). We then
173 obtained a collection of hmmsearch scores of resulting hits from the training metagenomes. By
174 using a similar method as described above, the cutoff was selected as the mean value of the
175 sharpest decreasing interval.

176

177 The following example demonstrates how the method above was used to curate the cutoffs for
178 hydrogenase enzymes. We then expanded this method to all genes using a similar method. We
179 downloaded the individual protein collections for each hydrogenase functional group from the
180 HydDB [45], which included [FeFe] Group A-C series, [Fe] Group, and [NiFe] Group 1-4
181 series. The individual hydrogenase functional groups were further categorized based on the
182 catalyzing directions, which included H₂-evolution, H₂-uptake, H₂-sensing, electron-
183 bifurcation, and bidirection. To define the NC cutoff (‘--cut_nc’ in hmmsearch) for individual
184 hydrogenase groups, we used the protein sequences from each hydrogenase group as the query
185 to hmmsearch against the overall hydrogenase collections. By plotting the resulting hmmsearch
186 hit scores against individual hmmsearch hits, we selected the mean value of the sharpest
187 decreasing interval as the cutoff value.

188

189 **Motif validation**

190 To automatically validate protein hits and avoid false positives, we introduced a motif
191 validation step by comparing protein motifs against a manually curated set of highly conserved
192 residues in important proteins. This manually curated set of highly conserved residues is
193 derived from either reported works or protein alignments from this study. We chose 20 proteins
194 associated with important metabolisms (with a focus on important biogeochemical cycling
195 steps) that are prone to be misannotated into proteins within the same protein family. Details of
196 these proteins are provided in [Additional file 8: Dataset S1](#). For example, DsrC (sulfite
197 reductase subunit C) and TusE (tRNA 2-thiouridine synthesizing protein E) are similar proteins
198 that are commonly misannotated. Both of them are assigned to the family KO:K11179 in the
199 KEGG database. To avoid assigning TusE as a sulfite reductase, we identified a specific motif
200 for DsrC but not TusE (GPXKXXCXXXGXPXPXXCX", where "X" stands for any amino
201 acid) [46]. We used these specific motifs to filter out proteins that have high sequence similarity
202 but functionally divergent homologs.

203

204 **Annotation of carbohydrate-active enzymes and peptidases**

205 For carbohydrate-active enzymes (CAZymes), dbCAN2 [39] was used to annotate proteins with
206 default settings. The hmmcan parser and HMM database (2019-09-05 release) were
207 downloaded from the dbCAN2 online repository (<http://bcb.unl.edu/dbCAN2/download/>) [39].
208 The non-redundant library of protein sequences which contains all the peptidase/inhibitor units
209 from the peptidase (inhibitor) database MEROPS [40] (known as the 'MEROPS pepunit'
210 database) was used as the reference database to search against putative peptidases and inhibitors
211 using DIAMOND. The settings used for the DIAMOND BLASTP search were "-k 1 -e 1e-10
212 --query-cover 80 --id 50" [47]. We used the 'MEROPS pepunit' database due to the fact that it
213 only includes the functional unit of peptidases/inhibitors [40] which can effectively avoid
214 potential non-specific hits.

215

216 **Implementation of METABOLIC-G and METABOLIC-C**

217 To target specific applications in processing omics datasets, we have implemented two versions
218 of METABOLIC: METABOLIC-G (genome version) and METABOLIC-C (community
219 version). METABOLIC-G intakes only genome files and provides analyses for individual
220 genome sequences (including three kinds of genomes, e.g., single-cell genomes, isolate
221 genomes, and metagenome-assembled genomes). The analyzing procedures of METABOLIC-
222 G for all these three kinds of genomes are the same.

223

224 METABOLIC-C includes an option for users to include metagenomic reads for mapping to
225 metagenome-assembled genomes (MAGs). Using Bowtie 2 (version \geq v2.3.4.1) [48],
226 metagenomic BAM files were generated by mapping all input metagenomic reads to gene
227 collections from input genomes. Subsequently, SAMtools (version \geq v0.1.19) [49], BAMtools
228 (version \geq v2.4.0) [50], and CoverM (<https://github.com/wwood/CoverM>) were used to convert
229 BAM files to sorted BAM files and to calculate the gene coverage. To calculate the relative
230 abundance of a specific biogeochemical cycling step, all the coverage of genes that are
231 responsible for this step were summed up and normalized by overall gene coverage. Reads from

232 single-cell and isolate genomes can also be mapped in an identical manner to metagenomes.
233 The gene coverage result generated by metagenomic read mapping was further used in
234 downstream processing steps to conduct community-scale interaction and network analyses.
235

236 **Classifying microbial genomes into taxonomic groups**

237 To study community-scale interactions and networks of each microbial group within the whole
238 community, we classified microbial genomes into individual taxonomic groups. GTDB-Tk
239 v0.1.3 [51] was used to assign taxonomy of input genomes with default settings. GTDB-Tk can
240 provide automated and objective taxonomic classification based on the rank-normalized
241 Genome Taxonomy Database (GTDB) taxonomy within which the taxonomy ranks were
242 established by a sophisticated criterion counting the relative evolutionary divergence (RED)
243 and average nucleotide identity (ANI) [51, 52]. Subsequently, genomes were clustered into
244 microbial groups at the phylum level, except for Proteobacteria which were replaced by its
245 subordinate classes due to its wide coverage. Taxonomic assignment information for each
246 genome was used in the downstream community analyses.
247

248 **Analyses and visualization of metabolic outputs, biogeochemical cycles, MW-scores, 249 functional networks, and metabolic Sankey diagrams**

250 To visualize the outputted metabolic results, the R script “*draw_biogeochemical_cycles.R*” was
251 used to draw the corresponding metabolic pathways for individual genomes. We integrated
252 HMM profiles that are related to biogeochemical activities and assigned HMM profiles to 31
253 distinct biogeochemical cycling steps (See details in “*METABOLIC_template_and_database*”
254 folder on the GitHub page). The script can generate figures showing biogeochemical cycles for
255 individual genomes and the summarized biogeochemical cycle for the whole community. By
256 using the results of metabolic profiling generated from hmmsearch and gene coverage from the
257 mapping of metagenomic reads, we can depict metabolic capacities of both individual genomes
258 and all genomes within a community as a whole. The community-level diagrams, including
259 sequential transformation diagrams, functional network diagrams, and metabolic Sankey
260 diagrams, were generated using both metabolic profiling and gene coverage results. The
261 diagrams are made by the scripts “*draw_sequential_reaction_diagram.R*”,
262 “*draw_metabolic_Sankey_diagram.R*”, and “*draw_functional_network_diagram.R*”,
263 respectively (For details, refer to GitHub wiki pages).
264

265 MW-score (metabolic weight score) is a metric reflecting the functional capacity and
266 abundance of a microbial community in co-sharing functional networks. It was calculated at
267 the community-scale level based on results of metabolic profiling and gene coverage from
268 metagenomic read mapping as described above. We divided metabolic/biogeochemical cycling
269 steps (31 in total) into a finer level – function (51 functions in total) – for better resolution on
270 reflecting functional networks. By using similar methods for determining metabolic
271 interactions (as described above), we selected functions that are shared among genomes. MW-
272 score for each function was calculated by summing up all the coverage values of each function
273 (calculated by summing up all coverage values of genomes that contain this function) and

274 subsequently normalizing it by the overall function coverage. For each function, the
275 contribution percentage of each microbial phylum (the default taxonomic level setting) was
276 also calculated accordingly. One can also change the taxonomic level setting to the resolution
277 of “class”, “order”, “family”, or “genus” to calculate the corresponding contribution percentage
278 of each microbial group. Two equations are provided as follows to calculate each function’s
279 MW-score (1) and the percentage of contribution of each microbial group to the MW-score (2):

280
$$MW_{f_i} = \frac{\sum_{g=g_1}^{g_n} C_{g_n} \cdot S_{f_i}}{\sum_{g=g_1, f=f_1}^{g_n, f_n} C_{g_n} \cdot S_{f_n}} \quad (1)$$

281
$$C_{prec_{f_i, p_j}} = \left(\frac{\sum_{g=g_1}^{g_k} C_{g_n} \cdot S_{f_i}}{\sum_{g=g_1, f=f_1}^{g_n, f_n} C_{g_n} \cdot S_{f_n}} \Big/ \frac{\sum_{g=g_1}^{g_n} C_{g_n} \cdot S_{f_i}}{\sum_{g=g_1, f=f_1}^{g_n, f_n} C_{g_n} \cdot S_{f_n}} \right) \times 100\% \quad (2)$$

282 within which $g_k \dots g_1 \in p_j$

283

284 In equation (1), MW refers to MW-score. f_i refers to the studied function (f) which ranks in the
285 (i) position amongst all functions. g_1 and g_n indicate the first and the last genome amongst all
286 genomes. f_1 and f_n indicate the first and the last function amongst all functions. C_g means the
287 coverage of a genome and S_f means the presence (denoted as 1) or absence (denoted as 0) state
288 of a function within that genome. In equation (2), C_{prec} refers to the contribution percentage
289 of a microbial group to the MW-score. p_j means the studied group (p) which ranks in the (j)
290 position amongst all groups. g_k and g_1 indicate the genomes which rank in the (k) position and
291 the (l) position amongst all genomes; the additional note $g_k \dots g_1 \in p_j$ indicates all the
292 genomes between these two belong to the studied group p_j .

293

294 **Example of METABOLIC analysis**

295 An example of community-scale analyses including element biogeochemical cycling and
296 sequential reaction analyses, functional network and metabolic Sankey visualization, and MW-
297 score calculation were conducted using a metagenomic dataset of microbial community
298 inhabiting deep-sea hydrothermal vent environment of Guaymas Basin in the Pacific Ocean
299 [53]. It contains 98 MAGs and 1 set of metagenomic reads (genomes were available at NCBI
300 BioProject PRJNA522654 and metagenomic reads were deposited to NCBI SRA with
301 accession as SRR3577362).

302

303 A metagenomic-based study of the microbial community from an aquifer adjacent to Colorado
304 River, located near Rifle, has provided an accurate reconstruction of the metabolism and
305 ecological roles of the microbial majority [2]. From underground water and sediments of the
306 terrestrial subsurface at Rifle, 2545 reconstructed MAGs were obtained (genomes are under
307 NCBI BioProject PRJNA288027). They were used as the *in silico* dataset to test
308 METABOLIC’s performance. First, all the microbial genomes were dereplicated by dRep
309 v2.0.5 [54] to pick the representative genomes for downstream analysis using the setting of ‘-
310 comp 85’. Then, METABOLIC-G was applied to profile the functional traits of these
311 representative genomes using default settings. Finally, the metabolic profile chart was depicted
312 by assigning functional traits to GTDB taxonomy-clustered genome groups.

313

314 **Test on software performance for different environments**
315 To benchmark and test the performance of METABOLIC in different environments, eight
316 datasets of metagenomes and metagenomic reads from marine, terrestrial, and human
317 environments were used. These included marine subsurface sediments [55] (Deep biosphere
318 beneath Hydrate Ridge offshore Oregon), freshwater lake [56] (Lake Tanganyika, eastern
319 Africa), colorectal cancer (CRC) patient gut [57], healthy human gut [57], deep-sea
320 hydrothermal vent [53] (Guaymas Basin, Gulf of California), terrestrial subsurface sediments
321 and water [2] (Rifle, CO, USA), meadow soils [58] (Angelo Coastal Range Reserve, CA, USA),
322 and advanced water treatment facility [59] (Groundwater Replenishment System, Orange
323 County, CA, USA). Default settings were used for running METABOLIC-C.
324

325 **Comparison of community-scale metabolism**
326 To compare the metabolic profile of two environments at the community scale, MW-score was
327 used as the benchmarker. Two sets of environment pairs were compared, including the pair of
328 marine subsurface sediments [55] and terrestrial subsurface sediments [2] and the pair of
329 freshwater lake [56] and deep-sea hydrothermal vent [53]. To demonstrate differences between
330 these environments in specific biogeochemical processes, we focused on the biogeochemical
331 cycling of sulfur. The sulfur biogeochemical cycling diagrams were depicted with the
332 annotation of the number and the coverage of genomes that contain each biogeochemical
333 cycling step.
334

335 **Metabolism in human microbiomes**
336 To inspect the metabolism of microorganisms in the human microbiome (associated with skin,
337 oral mucosa, conjunctiva, gastrointestinal tracts, etc.), a subset of KOfam HMMs (139 HMM
338 profiles) were used as markers to depict the human microbiome metabolism (parsed by
339 HuMiChip targeted functional gene families [60]). They included 10 function categories as
340 follows: amino acid metabolism, carbohydrate metabolism, energy metabolism, glycan
341 biosynthesis and metabolism, lipid metabolism, metabolism of cofactors and vitamins,
342 metabolism of other amino acids, metabolism of terpenoids and polyketides, nucleotide
343 metabolism, and translation. The CRC and healthy human gut (healthy control) sample datasets
344 were used as the input (Accession IDs: BioProject PRJEB7774 Sample 31874 and Sample
345 532796). Heatmap of presence/absence of these functions were depicted by R package
346 “pheatmap” [61] with 189 horizontal entries (there are duplications of HMM profiles among
347 function categories; for detailed human microbiome metabolism markers, refer to [Additional
348 file 9: Dataset S2](#)).
349

350 **Representation of microbial cell metabolism**
351 To provide a schematic representation of the metabolism of microbial cells, two microbial
352 genomes were used as examples, Hadesarchaea archaeon 1244-C3-H4-B1 and Nitrospirae
353 bacteria M_DeepCast_50m_m2_151. METABOLIC-G results of these two genomes, including
354 functional traits and KEGG modules, were used to draw the cell metabolism diagrams.
355

356 **Metatranscriptome analysis by METABOLIC**

357 METABOLIC-C can take metatranscriptomic reads as input into transcript coverage
358 calculation and integrate the result into downstream community analyses. METABOLIC-C
359 uses a similar method to that of gene coverage calculation, including mapping transcriptomic
360 reads to the gene collection from input genomes, converting BAM files to sorted BAM files,
361 and calculating the transcript coverage. The raw transcript coverage was further normalized by
362 the gene length and metatranscriptomic read number in Reads Per Kilobase of transcript, per
363 Million mapped reads (RPKM). Hydrothermal vent and background seawater transcriptomic
364 reads from Guaymas Basin (NCBI SRA accessions: SRR452448 and SRR453184) were used
365 to test the outcome of metatranscriptome analysis.

366 **RESULTS**

367 Given the ever-increasing number of microbial genomes from microbiome studies, we
368 developed METABOLIC to enable the metabolic pathway analysis and the visualization of
369 biogeochemical cycles and community-scale functional networks. METABOLIC has an
370 improved methodology to get fast, accurate, and robust annotation results, and it integrates a
371 variety of visualization functions for better interpreting the community-level functional
372 interactions and microbial contributions. While METABOLIC relies on microbial genomes and
373 metagenomic reads for underpinning its analyses for community-level functional interactions,
374 it can easily integrate transcriptomic datasets to provide an activity-based measure of
375 community networks. The scalable capacity, wide utility, and compatibility for analyzing
376 datasets from various environments make it a well-tailored tool for metabolic profiling of large
377 sets of genomes. In the following sections, the microbial community consisting of 98 MAGs
378 from a deep-sea hydrothermal vent was used as the input dataset if not mentioned otherwise.

380 **Workflow to determine the presence of metabolic pathways**

381 METABOLIC is written in Perl and R and is expected to run on Unix, Linux, or macOS. The
382 prerequisites are described on METABOLIC's GitHub wiki pages
383 (<https://github.com/AnantharamanLab/METABOLIC/wiki>). The input folder requires
384 microbial genome sequences in FASTA format and an optional set of genomic/metagenomic
385 reads which were used to reconstruct those genomes (Figure 1). The annotated proteins from
386 input genomic sequences are queried against HMM databases (KEGG KOfam, Pfam,
387 TIGRFam, and custom HMMs) using hmmsearch implemented within HMMER [42] which
388 applies methods to detect remote homologs as sensitively and efficiently as possible. After the
389 hmmsearch step, METABOLIC subsequently validates the primary outputs by a motif-
390 checking step for a subset of protein families; only those protein hits which successfully pass
391 this step are regarded as positive hits.

392
393 METABOLIC relies on matches to the above databases to infer the presence of specific
394 metabolic pathways in microbial genomes. Individual KEGG annotations are inferred in the
395 context of KEGG modules for a better interpretation of metabolic pathways. A KEGG module
396 is comprised of multiple steps with each step representing a distinct metabolic function. We

398 parsed the KEGG module database [62] to link the existing relationship of KO identifiers to
399 KEGG module identifiers to project our KEGG annotation result into the interactive network
400 which was constructed by individual building blocks – modules – for better representation of
401 metabolic blueprints of input genomes. In most cases, we used KOfam HMM profiles for
402 KEGG module assignments. For a specific set of important metabolic marker proteins and
403 commonly misannotated proteins, we also applied the TIGRfam/Pfam/custom HMM profiles
404 and motif-validation steps. The software has customizable settings for increasing or decreasing
405 the priority of specific databases, primarily meant to increase annotation confidence by
406 preferentially using custom HMM databases over KEGG KOfam when both targeting the same
407 set of proteins.

408
409 Since individual genomes from metagenomes and single-cell genomes can often have
410 incomplete metabolic pathways due to their low completeness compared to isolate genomes,
411 we provide an option to determine the completeness of a metabolic pathway (or a module here).
412 A user-defined cutoff is used to set the threshold of completeness for a given module to be
413 assigned as present (the default cutoff is the presence of 75% of metabolic steps/genes within
414 a given module), which is then used to produce a KEGG module presence/absence table. All
415 modules exceeding the cutoff value are determined to be present. Meanwhile, the
416 presence/absence information for each module step is also summarized in an overall output
417 table to facilitate further detailed investigations.

418
419 Outputs consist of six different results that are reported in an Excel spreadsheet ([Additional file
420 1: Figure S1](#)). These contain details of protein hits ([Additional file 1: Figure S1A](#)) which include
421 both presence/absence and protein names, presence/absence of functional traits ([Additional file
422 1: Figure S1B](#)), presence/absence of KEGG modules ([Additional file 1: Figure S1C](#)),
423 presence/absence of KEGG module steps ([Additional file 1: Figure S1D](#)), carbohydrate-active
424 enzyme (CAZyme) hits ([Additional file 1: Figure S1E](#)) and peptidase/inhibitor hits ([Additional
425 file 1: Figure S1F](#)). For each HMM profile, the protein hits from all input genomes can be used
426 to construct phylogenetic trees or further be combined with reference protein collections for
427 detailed evolutionary analyses.

428
429 **Quantitative visualization of biogeochemical cycles and sequential reactions**
430 After METABOLIC generates protein and pathway annotation results, the software further
431 identifies and highlights specific pathways of importance in microbiomes associated with
432 energy metabolism and biogeochemistry. To visualize pathways of biogeochemical
433 importance, it generates schematic profiles for nitrogen, carbon, sulfur, and other elemental
434 cycles for each genome. The set of genomes used as input is considered the “community”, and
435 each genome within is considered an “organism”. A summary schematic diagram at the
436 community level integrates results from all individual genomes within a given dataset ([Figure
437 2](#)) and includes computed abundances for each step in a biogeochemical cycle if the
438 genomic/metagenomic read datasets are provided. The genome number labeled in the figure
439 indicates the number/quantity of genomes that contain the specific gene components of a

440 biogeochemical cycling step (Figure 2) [2]. In other words, it represents the number of
441 organisms within a given community inferred to be able to perform a given metabolic or
442 biogeochemical transformation. The abundance percentage indicates the relative abundance of
443 microbial genomes that contain the specific gene components of a biogeochemical cycling step
444 among all microbial genomes in a given community (Figure 2) [2].

445

446 Microorganisms in nature often do not encode pathways for the complete transformation of
447 compounds. For example, microorganisms possess partial pathways for denitrification that can
448 release intermediate compounds like nitrite, nitric oxide, and nitrous oxide in lieu of nitrogen
449 gas which is produced by complete denitrification [63]. A greater energy yield could be
450 achieved if one microorganism conducts all steps associated with a pathway (such as
451 denitrification) [2] since it could fully use all available energy from the reaction. However, in
452 reality, few organisms in microbial communities carry out multiple steps in complex pathways;
453 organisms commonly rely on other members of microbial communities to conduct sequential
454 reactions in pathways [2, 64, 65]. Thus, to study this metabolic scenario in microbial
455 communities, METABOLIC summarizes and enables visualization of the genome number and
456 coverage (relative abundance) of microorganisms that are putatively involved in the sequential
457 transformation of both important inorganic and organic compounds (Figure 3). This provides a
458 quantitative calculation of microbial interactions and connections using shared metabolites
459 associated with inorganic and organic transformations. Additionally, it shows the intuitive
460 pattern of quantity and abundance of microorganisms that are able to conduct partial or all steps
461 for a given pathway, which potentially reflects the degree of resilience of a microbial
462 community.

463

464 **Calculation and visualization of functional networks, metabolic weight scores (MW- 465 scores), and microbial contribution to metabolic reactions**

466 Given the microbial pathway abundance information generated by METABOLIC, we identified
467 co-existing metabolisms in microbial genomes as a measure of connections between different
468 metabolic functions and biogeochemical steps. In the context of biogeochemistry, this approach
469 allows the evaluation of relatedness among biogeochemical steps and the connection
470 contribution by microorganisms. This is enabled at the resolution of individual microbial
471 groups based on the phylogenetic classification (Figure 4) assigned by GTDB-Tk [51]. As an
472 example, we have demonstrated this approach on a microbial community inhabiting deep-sea
473 hydrothermal vents. We divided the microbial community of deep-sea hydrothermal vents into
474 18 phylum-level groups (except for Proteobacteria which were divided into their subordinate
475 classes). The functional network diagrams were depicted at the resolution of both individual
476 phyla and the entire community level (Additional file 10: Dataset S3). Figure 4 demonstrates
477 metabolic connections that were represented with individual metabolic/biogeochemical cycling
478 steps depicted as nodes, and the connections between two given nodes depicted as edges. The
479 size of a given node is proportional to the degree (number of connections to each node). The
480 thickness of a given edge was depicted based on the average of gene coverage values of two
481 biogeochemical cycling steps (the connected nodes). More edges connecting two nodes

482 represent more connections between these two steps. The color of the edge corresponds to the
483 taxonomic group. At the whole community level, more abundant microbial groups were more
484 represented in the diagram (Figure 4). Overall, METABOLIC provides a comprehensive
485 approach to construct and visualize functional networks associated with important pathways of
486 energy metabolism and biogeochemical cycles in microbial communities and ecosystems.

487

488 To address the lack of quantitative and reproducible measures to represent potential metabolic
489 interactions in microbial communities, we developed a new metric that we termed MW-score
490 (metabolic weight scores) (Equations 1 and 2). MW-scores quantitatively measure “function
491 weights” within a microbial community as reflected by the metabolic profile and gene coverage.
492 As metabolic potential for the whole community was profiled into individual functions that
493 either mediated specific pathways or transformed certain substrates into products, a function
494 weight that reflects the abundance fraction for each function can be used to represent the overall
495 metabolic potential of the community. MW-scores resolved the functional capacity and
496 abundance in the co-sharing functional networks as studied and visualized in the above section.
497 More frequently shared functions and their higher abundances lead to higher MW-scores, which
498 quantitatively reflects the function weights in functional networks (Figure 5). MW-score
499 reflects the same functional networking pattern as the above description on the edges
500 (networking lines) connecting the nodes (metabolic steps) that – more edges connecting two
501 nodes indicates two steps are more co-shared, thicker edges indicate higher gene abundance for
502 the metabolic steps. The MW-scores can integratively represent these two networking patterns
503 and serve as metrics to measure these function weights. At the same time, we also calculated
504 each microbial group’s (phylum in this case) contribution to the MW-score of a specific
505 function within the community (Figure 5). A higher microbial group contribution percentage
506 value indicates that one function is more represented by the microbial group (for both gene
507 presence and abundance) in the functional networks. MW-scores provide a quantitative
508 measure of comparing function weights and microbial group contributions within functional
509 networks.

510

511 To understand the contributions of microbial groups associated with specific metabolic and
512 biogeochemical transformations, we developed an approach to visualize the connections among
513 specific taxonomic groups, metabolic reactions, and entire biogeochemical cycles such as
514 carbon, nitrogen, and sulfur cycles. Our approach involves the use of Sankey diagrams (also
515 called ‘Alluvial’ plots) to represent the fractions of metabolic functions that are contributed by
516 various microbial groups in a given community (Figure 6). It allows visualization of metabolic
517 reactions as the link between microbial contributors clustered as taxonomic groups and
518 biogeochemical cycles at a community level (Figure 6 and Additional file 10: Dataset S3). The
519 function fraction was calculated by accumulating the genome coverage values of genomes from
520 a specific microbial group that possesses a given functional trait. The width of curved lines
521 from a specific microbial group to a given functional trait indicates their corresponding
522 proportional contribution to a specific metabolism (Figure 6). Alternatively, the
523 genomic/metagenomic datasets which are used in constructing the above two diagrams:

524 functional network diagram (Figure 4) and metabolic Sankey diagram (Figure 6), can be
525 replaced by transcriptomic/metatranscriptomic datasets, and correspondingly, the gene
526 coverage values will be replaced by gene expression values, and therefore, diagrams will
527 represent the transcriptional activity patterns of functional network and microbial contribution
528 to metabolic reactions (Additional file 2, 3, 4, and 5: Figure S2, S3, S4, and S5).

529

530 To demonstrate this part of the workflow in reality, the microbial community consisting of 98
531 MAGs from a deep-sea hydrothermal vent was used as a test dataset. After running the
532 bioinformatic analyses described above, resulting tables and diagrams were compiled and
533 visualized accordingly (Figure 4, 5, 6 and Additional file 10: Dataset S3). Results for functional
534 networks and MW-scores of the deep-sea hydrothermal vent environment indicate that the
535 microbial community depends on mixotrophy and sulfur oxidation for energy conservation and
536 involves arsenate reduction potentially responsible for detoxification/arsenate resistance [66].
537 MW-scores indicate that amino acid utilization, complex carbon degradation, acetate oxidation,
538 and fermentation are the major heterotrophic metabolisms for this environment; CO₂-fixation
539 and sulfur oxidation also occupy a considerable functional fraction, which indicates
540 heterotrophy and autotrophy both contribute to energy conservation (Figure 5). As represented
541 by both MW-scores and metabolic Sankey diagram, Gammaproteobacteria are the most
542 numerically abundant group in the community and they occupy significant functional fractions
543 amongst both heterotrophic and autotrophic metabolisms (MW-score contribution ranging from
544 59-100%) (Figure 5, 6), which is consistent with previous findings in the Guaymas Basin
545 hydrothermal environment [53, 67]. Meanwhile, MW-scores also explicitly reflect the
546 involvement of other minor electron donors in energy conservation which are mainly
547 contributed by Gammaproteobacteria, such as hydrogen and methane (Figure 5). This is also
548 consistent with previous findings [53, 67] and indicates the accuracy and sensitivity of MW-
549 scores to reflect metabolic potentials.

550

551 **METABOLIC performance demonstration**

552 To test METABOLIC's performance on speed, we applied the software (METABOLIC-C
553 mode) to analyze the metagenomic dataset which includes 98 MAGs from a deep-sea
554 hydrothermal vent, and two sets of metagenomic reads (that are subsets of original reads with
555 10 million reads for each pair comprising ~10% of the total reads). The total running time was
556 ~3 hours using 40 CPU threads in a Linux version 4.15.0-48-generic server (Ubuntu v5.4.0).
557 The most compute-demanding step is hmmsearch, which took ~45 mins. When tested on
558 another dataset comprising ~3600 microbial genomes (data not shown), METABOLIC could
559 complete hmmsearch in ~5 hours by using 40 CPU threads, indicating its scalable capability on
560 analyzing thousands of genomes.

561

562 In order to test the accuracy of the results predicted by METABOLIC, we picked 15 bacterial
563 and archaeal genomes from Chloroflexi, Thaumarchaeota, and Crenarchaeota which are
564 reported to have 3 hydroxypropionate cycle (3HP) and/or 3-hydroxypropionate/4-
565 hydroxybutyrate cycle (3HP/4HB) for carbon fixation. METABOLIC predicted results in line

566 with annotations from the KEGG genome database which can be visualized in KEGG Mapper
567 ([Table 1](#)). Our predictions are also in accord with biochemical evidence of the existence of
568 corresponding carbon fixation pathways in each microbial group: 1) 3 out of 5 *Chloroflexi*
569 genomes are predicted by both METABOLIC and KEGG to possess the 3HP pathway and none
570 of all these *Chloroflexi* genomes are predicted to possess the 3HP/4HB pathway. This is
571 consistent with current reports based on biochemical and molecular experiments that only
572 organisms from the phylum *Chloroflexi* are known to possess the 3HP pathway [68] ([Table 1](#)).
573 2) All 5 *Thaumarchaeota* genomes and 2 out of 5 *Crenarchaeota* genomes are predicted by
574 both METABOLIC and KEGG to possess the 3HP/4HB pathway and none of these
575 *Thaumarchaeota* and *Crenarchaeota* genomes are predicted to possess the 3HP pathway. This
576 is consistent with current reports that only the 3HP/4HB pathway could be detected in
577 *Crenarchaeota* and *Thaumarchaeota* [69, 70] ([Table 1](#)). We have also applied METABOLIC
578 on a large well-studied dataset comprising 2545 metagenome-assembled genomes from
579 terrestrial subsurface sediments and groundwater [2]. The annotation results of METABOLIC
580 are consistent with previously described reports ([Additional file 6, 10: Figure S6, Dataset S3](#)).
581 These results suggest that METABOLIC can provide accurate annotations and perform well as
582 a functional predictor for microbial genomes and communities.
583

584 Currently, several software packages and online servers are available for genome annotation
585 and metabolic profiling. Comparing to other software/online servers including GhostKOALA
586 [71], BlastKOALA [71], KAAS [72], RAST/SEED [34], and eggNOG-mapper [73],
587 METABOLIC is unique in its ability to integrate multi-omic information towards elucidating
588 and visualizing community-level functional connections and the contribution of
589 microorganisms to biogeochemical cycles ([Figure 7A](#)). Additionally, in order to compare the
590 prediction performance of METABOLIC to others, we conducted parallel *in silico* experiments
591 ([Figure 7B](#)). We used two representative bacterial genomes as the test datasets. We randomly
592 picked 100 protein sequences from individual genomes and submitted them to annotation by
593 these six software/online servers. Predicted protein annotations by individual software and
594 online servers were compared to their original annotations that were provided by the NCBI
595 database ([Additional file 11, 12: Dataset S4, S5](#)). According to statistical methods of evaluating
596 binary classification [74], the following parameters were used to make the comparison: 1) recall
597 (also referred to as the sensitivity) as the true positive rate, 2) precision (also referred to as the
598 positive predictive value) which indicates the reproducibility and repeatability of a
599 measurement system, 3) accuracy which indicates the closeness of measurements to their true
600 values, and 4) F_1 value which is the harmonic mean of precision and recall, and reflects both
601 these two parameters. Among the tested software/online servers, the performance parameters
602 of METABOLIC consistently placed it as the top 3 and top 2 software for recall and F_1 and the
603 top 1 and top 2 software for precision and accuracy. These results demonstrate that
604 METABOLIC ([Figure 7B](#)) provides robust performance and consistent metabolic prediction
605 that facilitate accurate and reliable applicability for downstream data visualization and
606 community-level analyses.
607

608 To demonstrate the application and performance of METABOLIC in different samples, we
609 tested eight distinct environments (marine subsurface, terrestrial subsurface, deep-sea
610 hydrothermal vent, freshwater lake, gut microbiome from patients with colorectal cancer, gut
611 microbiome from healthy control, meadow soil, wastewater treatment facility). Overall, we
612 found METABOLIC to perform well across all the environments to profile microbial genomes
613 with functional traits and biogeochemical cycles ([Additional file 10: Dataset S3](#)). Among these
614 tested environments, we also performed community-scale metabolic comparisons based on the
615 MW-score ([Figure 8](#)). MW-score fraction at the community scale reflects the overall metabolic
616 profile distribution pattern. Specifically, we compared samples from terrestrial and marine
617 subsurface and samples from hydrothermal vent and freshwater lake. We observed that
618 terrestrial subsurface contains more abundant metabolic functions related to nitrogen cycling
619 compared to the marine subsurface ([Figure 8A](#)), consistent with the previous characterization
620 of these two environments [2, 75]. Deep-sea hydrothermal vent samples had a considerably
621 high concentration of methane and hydrogen [53] as compared to Lake Tanganyika (freshwater
622 lake). Consistent with this phenomenon, the deep-sea hydrothermal vent microbial community
623 has more abundant metabolic functions associated with methanotrophy and hydrogen oxidation
624 ([Figure 8B](#)). In order to focus on a specific biogeochemical cycle, we applied METABOLIC to
625 compare sulfur-related metabolisms at the community scale for these two environment pairs
626 ([Additional file 7: Figure S7](#)). Terrestrial subsurface contains genomes covering more sulfur
627 cycling steps compared to marine subsurface (7 steps vs 3 steps) ([Additional file 7: Figure
628 S7A](#)). Freshwater lake contains genomes involving almost all the sulfur cycling steps except
629 for sulfur reduction, while deep-sea hydrothermal vent contains less sulfur cycling steps (8
630 steps vs 6 steps) ([Additional file 7: Figure S7B](#)). Nevertheless, deep-sea hydrothermal vent has
631 a higher fraction of genomes (59/98) and a higher relative abundance (73%) of these genomes
632 involving sulfur oxidation compared to the freshwater lake ([Additional file 7: Figure S7B](#)). This
633 indicates that the deep-sea hydrothermal vent microbial community has a more biased sulfur
634 metabolism towards sulfur oxidation, which is consistent with previous metabolic
635 characterization on the dependency of elemental sulfur in this environment [53, 76-78].
636 Collectively, by characterizing community-scale metabolism, METABOLIC can facilitate the
637 comparison of overall functional profiles as well as for a particular elemental cycle.
638

639 **METABOLIC enables accurate reconstruction of cell metabolism**

640 To demonstrate applications of reconstructing and depicting cell metabolism based on
641 METABOLIC results, two microbial genomes were used as an example ([Figure 9](#)). As
642 illustrated in [Figure 9A](#), Hadesarchaea archaeon 1244-C3-H4-B1 has no TCA cycling gene
643 components, which is consistent with previous findings in archaea within this class [79].
644 Gluconeogenesis/glycolysis pathways are also lacking in the genome; since gluconeogenesis is
645 the central carbon metabolism responsible for generating sugar monomers which will be further
646 biosynthesized to polysaccharides as important cell structural components [80], the lack of this
647 pathway could be due to genome incompleteness. As an enigmatic archaeal class newly
648 discovered in the recent decade, Hadesarchaea have distinctive metabolisms that separate them
649 from conventional euryarchaeotal groups. They almost lost all TCA cycle gene components for

650 the production of acetyl-CoA; while they could metabolize amino acids in a heterotrophic
651 lifestyle [79]. It is posited that the Hadesarchaea genome has been subjected to a streamlining
652 process possibly due to nutrient limitations in their surrounding environments [79]. Due to their
653 metabolic novelty and limited available genomes at the current time, there are still uncertainties
654 on unknown/hypothetical genes and pathways and unclassified metabolic potential across the
655 whole class. The previous metabolic characterization on four Hadesarchaea genomes indicates
656 Hadesarchaea members could anaerobically oxidize CO, and H₂ was produced as the side
657 product [79]. In the Hadesarchaea archaeon 1244-C3-H4-B1 genome, METABOLIC results
658 indicate the loss of all anaerobic carbon-monoxide dehydrogenase gene components, which
659 suggests the distinctive metabolism of this Hadesarchaea archaeon from others and highlights
660 the accuracy of METABOLIC in reflecting functional details.

661

662 We also reconstructed the metabolism for Nitrospirae bacteria M_DeepCast_50m_m2_151, a
663 member of the Nitrospirae phylum reconstructed from Lake Tanganyika [56] ([Figure 9B](#)). It
664 contains the full pathway for the TCA cycle and gluconeogenesis/glycolysis. Furthermore, it
665 also has the full set of oxidative phosphorylation complexes for energy conservation and
666 functional genes for nitrite oxidation to nitrate. Other nitrogen cycling metabolisms identified
667 in this genome include ammonium oxidation, urea utilization, and nitrite reduction to nitric
668 oxide. The reverse TCA cycle pathway was identified for carbon fixation. The metabolic
669 profiling result is in accord with the fact that Nitrospirae is a well-known nitrifying bacterial
670 class capable of nitrite oxidation and living an autotrophic lifestyle [80]. Additionally, their
671 more abundant distribution in nature compared to other nitrite-oxidizing bacteria such as
672 *Nitrobacter* indicates their significant contribution to nitrogen cycling in the environment [80].
673 This highlights the ability of METABOLIC in reflecting functional details of more common
674 and prevalent microorganisms compared to the Hadesarchaea archaeon. Notably as discovered
675 from METABOLIC analyses, this bacterial genome also contains a wide range of transporter
676 enzymes on the cell membrane, including mineral and organic ion transporters, sugar and lipid
677 transporters, phosphate and amino acid transporters, heme and urea transporters,
678 lipopolysaccharide and lipoprotein releasing system, bacterial secretion system, etc., which
679 indicates its metabolic versatility and potential interactive activities with other organisms and
680 the ambient environment. Collectively, METABOLIC result of functional profiling provides
681 an intuitively-represented summary of a single microbial genome which enables depicting cell
682 metabolism for better visualizing the functional capacity.

683

684 **METABOLIC accurately represents metabolism in the human microbiome**

685 In addition to resolving microbial metabolism and biogeochemistry in environmental
686 microbiomes, METABOLIC also accurately identifies metabolic traits associated with human
687 microbiomes. The implications of microbial metabolism on human health largely remain a
688 black box, much like microbial contributions to biogeochemical cycling. We demonstrate the
689 utility of METABOLIC in human microbiomes using publicly available data from stool
690 samples collected from patients with colorectal cancer and healthy individuals. From this study,
691 we selected stool metagenomes from one colorectal cancer (CRC) and an age and sex-matched

692 healthy control to conduct the comparison. The heatmap indicates the human microbiome
693 functional profiles of both samples based on the marker gene presence/absence patterns (Figure
694 10). As an example of METABOLIC's application, we demonstrate that there were 28 markers
695 with variations > 10% in terms of the marker-containing genome fractions between these two
696 samples (Figure 10, Additional file 13: Dataset S6). These 28 markers involved all the ten
697 metabolic categories except for lipid metabolism and translation, suggesting the broad
698 functional differences between these two samples. In addition to analyzing human microbiome
699 specific functional markers, METABOLIC can be used to visualize elemental nutrient cycling
700 and analyze metabolic interactions in human microbiomes. Overall it enables systematic
701 characterization of the composition, structure, function, and interaction of microbial
702 metabolisms in the human microbiome and facilitates omics-based studies of microbial
703 community on human health [60].

705 DISCUSSION

706 The rapid increase in the availability of sequenced microbial genomes, metagenome-assembled
707 genomes, and single-cell genomes has significantly benefited ecogenomic research on
708 unraveling microbial functional roles and their metabolic contribution to biogeochemical
709 cycles. Tools that enable to conduct accurate and reproducible functional profiling on genomic
710 blueprints at the scale of both individual microorganisms and the whole microbial community
711 offered significant applications and advances. They are fundamental to facilitate understanding
712 of community-level functions, activities, interactions, and functional contributions in the era of
713 multi-omics. An ideal tool for microbial biogeochemical profiling needs consideration on better
714 organizing, interpreting, and visualizing the functional profile information; this is especially
715 important for dealing with thousands of genomes reconstructed from metagenomes and
716 studying the community-scale interactive metabolisms. Meanwhile, fast, accurate, robust
717 performance and wide usage of the tool will allow for providing reliability and efficiency.

718 Here we developed METABOLIC for profiling metabolisms, biogeochemical pathways, and
719 community-scale functional networks. Instead of solely depending on widely adopted protein
720 annotation databases, in METABOLIC two additional steps were added in order to accurately
721 predict protein functions and reconstruct metabolic pathways. First, for
722 TIGRfam/Pfam/Custom HMM profile databases, default NC/TC thresholds are often set too
723 low to avoid noisy signals especially for annotating proteins from large sets of metagenomes
724 wherein similar protein families often co-exist. This frequently leads to misannotations. To
725 avoid this, we collected hmmsearch scores of previous annotation results and plotted these
726 scores as a function of all annotations, and manually curated NC/TC by specifically picking the
727 sharpest decreasing interval as the adjusted cutoff. Second, the motif validation step involves
728 comparing potential hits to a set of manually curated highly conserved amino acid residues.
729 This helps to distinguish two protein families with high sequence identity but different
730 functions which are often difficult to separate by HMM profile-based annotations. These two
731 steps help to filter out non-specific and cross-talking hits of important functional proteins for
732 downstream bioinformatic analyses. After obtaining predicted metabolic pathways, many other

734 software/online servers mostly provide raw annotation results with overwhelming yet
735 unorganized details on characterizing protein functions. For microbial ecologists it is
736 fundamental to provide organized and intuitive results to facilitate understanding on the whole
737 landscape of biogeochemical cycling capacities. In METABOLIC, such a function was
738 developed to enable visualizing the presence/absence state of each step of biogeochemical
739 cycles for individual genomes and the whole microbial community. Combined with gene
740 abundance information calculated by metagenomic read mapping, we can identify the relative
741 abundance for each step of biogeochemical cycles. Furthermore, METABOLIC can also
742 visualize sequential reaction patterns for important organic and inorganic compound
743 transformations. This visualization function of METABOLIC is practical for representing the
744 “metabolic handoff” scenario of within-community interactions [2]. METABOLIC can be
745 implemented in human microbiome with the same performance. Recently, METABOLIC was
746 applied to stool metagenomic samples from 667 individuals who either were healthy or had
747 adenomas or carcinomas of the colon, to profile organic/inorganic sulfate reduction and sulfide
748 production [81]. This has considerably enlarged the utility of METABOLIC in community-
749 scale investigation on human microbiomes for purposes of systematic microbiota-disease
750 studies

751
752 Previously, the community networks reflected by microbial genomes mostly focused on
753 modeling reactions that are linked by metabolizing substrates and generating products [15, 19,
754 26]. On the contrary, METABOLIC was developed for a different purpose to study microbially-
755 mediated biogeochemical processes. In METABOLIC the community-scale functional network
756 provides an intuitive perspective on the metabolic connectivity among
757 biogeochemical/metabolic steps and microbial contributions to these functions. MW-score, a
758 metric that was built based on the same notion and methodology, offers quantitative
759 measurement for these connected functions. Combined together they represent which functions
760 are more centralized (connected with others) and important (weighted with higher relative
761 abundance) in the co-sharing functional networks and which groups of microbial players
762 contribute to these functions. Additionally, metabolic Sankey diagrams can be drawn to further
763 visualize the microbial group contributions to different functions and biogeochemical cycles.
764 As gene coverages generated by metagenomic read mapping can be replaced by transcript
765 coverages generated by transcriptomic reads mapping, we broaden the usage in reflecting active
766 function connections and weights. In practical applications, functional networks and MW-
767 scores can be made in a standardized, reproducible, and normalized manner, so parallel
768 comparisons between communities (or samples) are applicable. The visualized network and
769 Sankey diagram can also offer intuitive representations of functional connections and microbial
770 contribution at both individual function and community-scale levels by using customized color
771 schemes. There are other read-based metagenomic profiling tools, e.g., MetaPhlAn [28] and
772 MEGAN [82], that can study the taxonomical and functional composition of microbiome at the
773 community-scale level. Compared to read-based approaches which largely depend on the
774 comprehensiveness of reference databases to capture microbial organisms, METABOLIC
775 depends on the annotation of MAGs that is free from the limitation of reference databases on

776 novel and rare organism characterization. METABOLIC specifically provides additional
777 functionalities on annotation validation, result organization, and visualization which are
778 meaningful to give reliable and easily accessible functional profiling results for microbial
779 ecologists and biogeochemists to have a comprehensive understanding on the whole landscape
780 of biogeochemical cycling capacities.

781

782 CONCLUSIONS

783 Metabolic functional profile of microbial genomes at the scale of individual organisms and
784 communities is essential to have a comprehensive understanding of ecosystem processes, and
785 as a conduit for enabling functional trait-based modeling of biogeochemistry. We have
786 developed METABOLIC as a metabolic functional profiler that goes above and beyond current
787 frameworks of genome/protein annotation platforms in providing protein annotations and
788 metabolic pathway analyses that are used for inferring the contribution of microorganisms,
789 metabolism, interactions, activity, and biogeochemistry at the community-scale. METABOLIC
790 facilitates standardization and integration of genome-informed metabolism into metabolic and
791 biogeochemical models. We anticipate that METABOLIC will enable easier interpretation of
792 microbial metabolism and biogeochemistry from metagenomes and genomes and enable
793 microbiome research in diverse fields.

794

795 Additional files

796 **Additional file 1: Figure S1.** METABOLIC result table report

797 **Additional file 2: Figure S2.** Functional network diagram based on the transcriptomic dataset
798 from a hydrothermal vent sample

799 **Additional file 3: Figure S3.** Functional network diagram based on the transcriptomic dataset
800 from a deep-sea background sample

801 **Additional file 4: Figure S4.** Microbial Sankey diagram based on the transcriptomic dataset
802 from a hydrothermal vent sample

803 **Additional file 5: Figure S5.** Microbial Sankey diagram based on the transcriptomic dataset
804 from a deep-sea background sample

805 **Additional file 6: Figure S6.** Metabolic profile diagram of a subsurface microbial community
806 from Rifle, Colorado, USA

807 **Additional file 7: Figure S7.** Comparison of sulfur metabolism at the community scale level

808 **Additional file 8: Dataset S1.** Conserved motif sequences and motif pairs used in our analyses

809 **Additional file 9: Dataset S2.** Summary table of Human Microbiome Marker genes

810 **Additional file 10: Dataset S3.** METABOLIC result of metagenomes from eight different
811 environments

812 **Additional file 11: Dataset S4.** Comparison of protein prediction performance among five
813 software packages/online servers on the genome of *Escherichia coli* O157H7 str. Sakai

814 **Additional file 12: Dataset S5.** Comparison of protein prediction performance among five
815 software packages/online servers on the genome of *Pseudomonas aeruginosa* PAO1

816 **Additional file 13: Dataset S6.** Human microbiome functional profiling results for a CRC
817 subject and a healthy control

818

819 **Declarations**

820

821 **Acknowledgments**

822 We thank the comments and suggestions from the users of METABOLIC, which helped to
823 improve and expand the functions of this software.

824

825 **Funding**

826 We thank the University of Wisconsin - Office of the Vice-Chancellor for Research and
827 Graduate Education, University of Wisconsin – Department of Bacteriology, Wisconsin
828 Alumni Research Foundation, and the University of Wisconsin – College of Agriculture and
829 Life Sciences for their support. KA acknowledges support from the National Science
830 Foundation under Grant No. NSF-DBI 2047598. PQT was supported by the Natural Sciences
831 and Engineering Research Council of Canada (NSERC). KK is supported by a Wisconsin
832 Distinguished Graduate Fellowship Award from the University of Wisconsin-Madison, and a
833 William H. Peterson Fellowship Award from the Department of Bacteriology, University of
834 Wisconsin-Madison. ESC is supported by an NLM training grant to the Computation and
835 Informatics in Biology and in part by the Medicine Scientist Training Program (NLM
836 5T15LM007359).

837

838 **Authors' contributions**

839 ZZ and KA conceptualized and designed the study. ZZ and PQT wrote the Perl and R scripts.
840 ZZ ran the test data and improved the software. YL provided a part of the databases. PQT,
841 AMB, KK, ESC, and UK provided ideas and comments, helped to set up the GitHub page, and
842 contributed to improving the overall performance of the software. ZZ and KA wrote the
843 manuscript, and all authors contributed and approved the final edition of the manuscript.

844

845 **Corresponding author**

846 Correspondence to Karthik Anantharaman.

847

848 **Ethics declarations**

849 **Ethics approval and consent to participate**

850 Not applicable.

851

852 **Consent for publication**

853 Not applicable.

854

855 **Competing interests**

856 The authors declare that they have no competing interests.

857 **References:**

858 1. Wu X, Holmfeldt K, Hubalek V, Lundin D, Astrom M, Bertilsson S, Dopson M:
859 **Microbial metagenomes from three aquifers in the Fennoscandian shield**
860 **terrestrial deep biosphere reveal metabolic partitioning among populations.** *ISME*
861 *J* 2016, **10**:1192-1203.

862 2. Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, Thomas BC,
863 Singh A, Wilkins MJ, Karaoz U, et al: **Thousands of microbial genomes shed light**
864 **on interconnected biogeochemical processes in an aquifer system.** *Nat Commun*
865 2016, **7**:13219.

866 3. Probst AJ, Ladd B, Jarett JK, Geller-McGrath DE, Sieber CMK, Emerson JB,
867 Anantharaman K, Thomas BC, Malmstrom RR, Stieglmeier M, et al: **Differential**
868 **depth distribution of microbial function and putative symbionts through**
869 **sediment-hosted aquifers in the deep terrestrial subsurface.** *Nat Microbiol* 2018,
870 **3**:328-336.

871 4. Siegl A, Kamke J, Hochmuth T, Piel J, Richter M, Liang C, Dandekar T, Hentschel U:
872 **Single-cell genomics reveals the lifestyle of Poribacteria, a candidate phylum**
873 **symbiotically associated with marine sponges.** *ISME J* 2011, **5**:61-70.

874 5. Iverson V, Morris RM, Frazar CD, Berthiaume CT, Morales RL, Armbrust EV:
875 **Untangling Genomes from Metagenomes: Revealing an Uncultured Class of**
876 **Marine Euryarchaeota.** *Science* 2012, **335**:587-590.

877 6. Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, Beghini F, Manghi P,
878 Tett A, Ghensi P, et al: **Extensive Unexplored Human Microbiome Diversity**
879 **Revealed by Over 150,000 Genomes from Metagenomes Spanning Age,**
880 **Geography, and Lifestyle.** *Cell* 2019, **176**:649-662 e620.

881 7. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy T,
882 Schulz F, Jarett J, Rivers AR, Eloe-Fadrosh EA: **Minimum information about a**
883 **single amplified genome (MISAG) and a metagenome-assembled genome**
884 **(MIMAG) of bacteria and archaea.** *Nat Biotechnol* 2017, **35**:725.

885 8. Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN,
886 Hugenholtz P, Tyson GW: **Recovery of nearly 8,000 metagenome-assembled**
887 **genomes substantially expands the tree of life.** *Nat Microbiol* 2017, **2**:1533-1542.

888 9. Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, Butterfield
889 CN, Hernsdorf AW, Amano Y, Ise K: **A new view of the tree of life.** *Nat Microbiol*
890 2016, **1**:16048.

891 10. Kraemer S, Ramachandran A, Colatriano D, Lovejoy C, Walsh DA: **Diversity and**
892 **biogeography of SAR11 bacteria from the Arctic Ocean.** *ISME J* 2020, **14**:79-90.

893 11. Ruuskanen MO, Colby G, St Pierre KA, St Louis VL, Aris-Brosou S, Poulain AJ:
894 **Microbial genomes retrieved from High Arctic lake sediments encode for**
895 **adaptation to cold and oligotrophic environments.** *Limnol Oceanogr* 2020,
896 **65**:S233-S247.

897 12. Madsen EL: **Microorganisms and their roles in fundamental biogeochemical cycles.**
898 *Curr Opin Biotechnol* 2011, **22**:456-464.

899 13. Abreu NA, Taga ME: **Decoding molecular interactions in microbial communities.**
900 *FEMS Microbiol Rev* 2016, **40**:648-663.

901 14. Morris BEL, Henneberger R, Huber H, Moissl-Eichinger C: **Microbial syntrophy: 902 interaction for the common good.** *FEMS Microbiol Rev* 2013, **37**:384-406.

903 15. Zelezniak A, Andrejev S, Ponomarova O, Mende DR, Bork P, Patil KR: **Metabolic 904 dependencies drive species co-occurrence in diverse microbial communities.** *Proc 905 Natl Acad Sci U S A* 2015, **112**:6449.

906 16. Baker BJ, Lazar CS, Teske AP, Dick GJ: **Genomic resolution of linkages in carbon, 907 nitrogen, and sulfur cycling among widespread estuary sediment bacteria.** *Microbiome* 908 2015, **3**.

909 17. Morris BE, Henneberger R, Huber H, Moissl-Eichinger C: **Microbial syntrophy: 910 interaction for the common good.** *FEMS Microbiol Rev* 2013, **37**:384-406.

911 18. Graf DR, Jones CM, Hallin S: **Intergenomic comparisons highlight modularity of 912 the denitrification pathway and underpin the importance of community structure 913 for N₂O emissions.** *PLoS One* 2014, **9**:e114118.

914 19. Machado D, Andrejev S, Tramontano M, Patil KR: **Fast automated reconstruction of 915 genome-scale metabolic models for microbial species and communities.** *Nucleic 916 Acids Res* 2018, **46**:7542-7553.

917 20. Karp PD, Billington R, Caspi R, Fulcher CA, Latendresse M, Kothari A, Keseler IM, 918 Krummenacker M, Midford PE, Ong Q, et al: **The BioCyc collection of microbial 919 genomes and metabolic pathways.** *Briefings in Bioinformatics* 2019, **20**:1085-1093.

920 21. Diener C, Gibbons SM, Resendis-Antonio O: **MICOM: Metagenome-Scale 921 Modeling To Infer Metabolic Interactions in the Gut Microbiota.** *mSystems* 2020, 922 **5**:e00606-00619.

923 22. Zimmermann J, Kaleta C, Waschyna S: **gapseq: informed prediction of bacterial 924 metabolic pathways and reconstruction of accurate metabolic models.** *Genome 925 Biol* 2021, **22**:81-81.

926 23. Karp PD, Midford PE, Billington R, Kothari A, Krummenacker M, Latendresse M, 927 Ong WK, Subhraveti P, Caspi R, Fulcher C, et al: **Pathway Tools version 23.0 update: 928 software for pathway/genome informatics and systems biology.** *Briefings in 929 Bioinformatics* 2021, **22**:109-126.

930 24. Zorrilla F, Buric F, Patil KR, Zelezniak A: **metaGEM: reconstruction of genome 931 scale metabolic models directly from metagenomes.** *Nucleic Acids Res* 2021, 932 **49**:e126-e126.

933 25. Belcour A, Frioux C, Aite M, Bretaudeau A, Hildebrand F, Siegel A: **Metage2Metabo, 934 microbiota-scale metabolic complementarity for the identification of key species.** 935 *eLife* 2020, **9**:e61968.

936 26. Gu C, Kim GB, Kim WJ, Kim HU, Lee SY: **Current status and applications of 937 genome-scale metabolic models.** *Genome Biol* 2019, **20**:121.

938 27. Shaffer M, Borton MA, McGivern BB, Zayed AA, La Rosa SL, Soden LM, Liu P, 939 Narrowe AB, Rodríguez-Ramos J, Bolduc B, et al: **DRAM for distilling microbial 940 metabolism to automate the curation of microbiome function.** *Nucleic Acids Res* 941 2020, **48**:8883-8900.

942 28. Beghini F, McIver LJ, Blanco-Míguez A, Dubois L, Asnicar F, Maharjan S, Mailyan A, 943 Manghi P, Scholz M, Thomas AM, et al: **Integrating taxonomic, functional, and 944 strain-level profiling of diverse microbial communities with bioBakery 3.** *eLife*

945 2021, **10**:e65088.

946 29. Hug Laura A, Co R: **It Takes a Village: Microbial Communities Thrive through**
947 **Interactions and Metabolic Handoffs.** *mSystems*, **3**:e00152-00117.

948 30. Kanehisa M, Goto S: **KEGG: kyoto encyclopedia of genes and genomes.** *Nucleic*
949 *Acids Res* 2000, **28**:27-30.

950 31. Caspi R, Foerster H, Fulcher CA, Hopkinson R, Ingraham J, Kaipa P, Krummenacker
951 M, Paley S, Pick J, Rhee SY, et al: **MetaCyc: a multiorganism database of metabolic**
952 **pathways and enzymes.** *Nucleic Acids Res* 2006, **34**:D511-516.

953 32. Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, Heger A,
954 Hetherington K, Holm L, Mistry J, et al: **Pfam: the protein families database.** *Nucleic*
955 *Acids Res* 2014, **42**:D222-230.

956 33. Selengut JD, Haft DH, Davidsen T, Ganapathy A, Gwinn-Giglio M, Nelson WC,
957 Richter AR, White O: **TIGRFAMs and Genome Properties: tools for the assignment**
958 **of molecular function and biological process in prokaryotic genomes.** *Nucleic Acids*
959 *Res* 2007, **35**:D260-264.

960 34. Overbeek R, Olson R, Pusch GD, Olsen GJ, Davis JJ, Disz T, Edwards RA, Gerdes S,
961 Parrello B, Shukla M: **The SEED and the Rapid Annotation of microbial genomes**
962 **using Subsystems Technology (RAST).** *Nucleic Acids Res* 2013, **42**:D206-D214.

963 35. Huerta-Cepas J, Szklarczyk D, Forsslund K, Cook H, Heller D, Walter MC, Rattei T,
964 Mende DR, Sunagawa S, Kuhn M, et al: **eggNOG 4.5: a hierarchical orthology**
965 **framework with improved functional annotations for eukaryotic, prokaryotic and**
966 **viral sequences.** *Nucleic Acids Res* 2016, **44**:D286-D293.

967 36. Kanehisa M, Araki M, Goto S, Hattori M, Hirakawa M, Itoh M, Katayama T,
968 Kawashima S, Okuda S, Tokimatsu T, Yamanishi Y: **KEGG for linking genomes to**
969 **life and the environment.** *Nucleic Acids Res* 2008, **36**:D480-D484.

970 37. Schimel J: **1.13 - Biogeochemical Models: Implicit versus Explicit Microbiology.** In
971 *Global Biogeochemical Cycles in the Climate System.* Edited by Schulze E-D,
972 Heimann M, Harrison S, Holland E, Lloyd J, Prentice IC, Schimel D. San Diego:
973 Academic Press; 2001: 177-183

974 38. Graham EB, Knelman JE, Schindlbacher A, Siciliano S, Breulmann M, Yannarell A,
975 Beman JM, Abell G, Philippot L, Prosser J, et al: **Microbes as Engines of Ecosystem**
976 **Function: When Does Community Structure Enhance Predictions of Ecosystem**
977 **Processes?** *Front Microbio* 2016, **7**.

978 39. Zhang H, Yohe T, Huang L, Entwistle S, Wu P, Yang Z, Busk PK, Xu Y, Yin Y:
979 **dbCAN2: a meta server for automated carbohydrate-active enzyme annotation.**
980 *Nucleic Acids Res* 2018, **46**:W95-W101.

981 40. Rawlings ND, Barrett AJ, Finn R: **Twenty years of the MEROPS database of**
982 **proteolytic enzymes, their substrates and inhibitors.** *Nucleic Acids Res* 2016,
983 **44**:D343-D350.

984 41. Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, Ogata H:
985 **KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive**
986 **score threshold.** *Bioinformatics* 2019, **36**:2251-2252.

987 42. Finn RD, Clements J, Eddy SR: **HMMER web server: interactive sequence**
988 **similarity searching.** *Nucleic Acids Res* 2011, **39**:W29-37.

989 43. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ: **Prodigal: 990 prokaryotic gene recognition and translation initiation site identification.** *BMC 991 Bioinformatics* 2010, **11**:119.

992 44. UniProt C: **UniProt: a worldwide hub of protein knowledge.** *Nucleic Acids Res* 2019, 993 **47**:D506-D515.

994 45. Sondergaard D, Pedersen CN, Greening C: **HydDB: A web tool for hydrogenase 995 classification and analysis.** *Sci Rep* 2016, **6**:34212.

996 46. Venceslau SS, Stockdreher Y, Dahl C, Pereira IA: **The "bacterial heterodisulfide" 997 DsrC is a key protein in dissimilatory sulfur metabolism.** *Biochim Biophys Acta 998* 2014, **1837**:1148-1164.

999 47. Buchfink B, Xie C, Huson DH: **Fast and sensitive protein alignment using 1000 DIAMOND.** *Nat Methods* 2015, **12**:59-60.

1001 48. Langmead B, Salzberg SL: **Fast gapped-read alignment with Bowtie 2.** *Nat Methods 1002* 2012, **9**:357.

1003 49. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, 1004 Durbin R, Genome Project Data Processing S: **The Sequence Alignment/Map format 1005 and SAMtools.** *Bioinformatics* 2009, **25**:2078-2079.

1006 50. Barnett DW, Garrison EK, Quinlan AR, Stromberg MP, Marth GT: **BamTools: a C++ 1007 API and toolkit for analyzing and managing BAM files.** *Bioinformatics* 2011, 1008 **27**:1691-1692.

1009 51. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarszewski A, Chaumeil P-A, 1010 Hugenholtz P: **A standardized bacterial taxonomy based on genome phylogeny 1011 substantially revises the tree of life.** *Nat Biotechnol* 2018, **36**:996-1004.

1012 52. Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH: **GTDB-Tk: a toolkit to classify 1013 genomes with the Genome Taxonomy Database.** *Bioinformatics* 2020, **36**:1925-1927.

1014 53. Anantharaman K, Breier JA, Sheik CS, Dick GJ: **Evidence for hydrogen oxidation 1015 and metabolic plasticity in widespread deep-sea sulfur-oxidizing bacteria.** *Proc 1016 Natl Acad Sci U S A* 2013, **110**:330.

1017 54. Olm MR, Brown CT, Brooks B, Banfield JF: **dRep: a tool for fast and accurate 1018 genomic comparisons that enables improved genome recovery from metagenomes 1019 through de-replication.** *ISME J* 2017, **11**:2864.

1020 55. Glass JB, Ranjan P, Kretz CB, Nunn BL, Johnson AM, Xu M, McManus J, Stewart FJ: 1021 **Microbial metabolism and adaptations in Atribacteria-dominated methane 1022 hydrate sediments.** *Environ Microbiol* 2021, **23**:4646-4660.

1023 56. Tran PQ, Bachand SC, McIntyre PB, Kraemer BM, Vadeboncoeur Y, Kimirei IA, 1024 Tamatamah R, McMahon KD, Anantharaman K: **Depth-discrete metagenomics 1025 reveals the roles of microbes in biogeochemical cycling in the tropical freshwater 1026 Lake Tanganyika.** *ISME J* 2019, **15**:1971-1986.

1027 57. Feng Q, Liang S, Jia H, Stadlmayr A, Tang L, Lan Z, Zhang D, Xia H, Xu X, Jie Z, et 1028 al: **Gut microbiome development along the colorectal adenoma–carcinoma 1029 sequence.** *Nat Commun* 2015, **6**:6528.

1030 58. Diamond S, Andeer PF, Li Z, Crits-Christoph A, Burstein D, Anantharaman K, Lane 1031 KR, Thomas BC, Pan C, Northen TR, Banfield JF: **Mediterranean grassland soil C– 1032 N compound turnover is dependent on rainfall and depth, and is mediated by**

1033 59. **genomically divergent microorganisms.** *Nat Microbiol* 2019, **4**:1356-1367.

1034 59. Stamps BW, Leddy MB, Plumlee MH, Hasan NA, Colwell RR, Spear JR: **Characterization of the Microbiome at the World's Largest Potable Water Reuse**

1035 **Facility.** *Front Microbiol* 2018, **9**.

1036 60. Tu Q, He Z, Li Y, Chen Y, Deng Y, Lin L, Hemme CL, Yuan T, Van Nostrand JD, Wu L, et al: **Development of HuMiChip for Functional Profiling of Human Microbiomes.** *PLoS One* 2014, **9**:e90546.

1037 61. Kolde R, Kolde MR: **Package 'pheatmap'.** *R Package* 2015, **1**:790.

1038 62. Muto A, Kotera M, Tokimatsu T, Nakagawa Z, Goto S, Kanehisa M: **Modular**

1039 **architecture of metabolic pathways revealed by conserved sequences of reactions.**

1040 *Journal of Chemical Information and Modeling* 2013, **53**:613-622.

1041 63. Kuypers MMM, Marchant HK, Kartal B: **The microbial nitrogen-cycling network.**

1042 *Nat Rev Microbiol* 2018, **16**:263-276.

1043 64. Hug LA, Co R: **It Takes a Village: Microbial Communities Thrive through**

1044 **Interactions and Metabolic Handoffs.** *mSystems* 2018, **3**:e00152-00117.

1045 65. Graf DRH, Jones CM, Hallin S: **Intergenomic Comparisons Highlight Modularity**

1046 **of the Denitrification Pathway and Underpin the Importance of Community**

1047 **Structure for N₂O Emissions.** *PLoS One* 2014, **9**:e114118.

1048 66. Mukhopadhyay R, Rosen BP, Phung LT, Silver S: **Microbial arsenic: from geocycles**

1049 **to genes and enzymes.** *FEMS Microbiol Rev* 2002, **26**:311-325.

1050 67. Zhou Z, Liu Y, Pan J, Cron BR, Toner BM, Anantharaman K, Breier JA, Dick GJ, Li M: **Gammaproteobacteria mediating utilization of methyl-, sulfur- and petroleum**

1051 **organic compounds in deep ocean hydrothermal plumes.** *ISME J* 2020, **14**:3136-

1052 3148.

1053 68. Shih PM, Ward LM, Fischer WW: **Evolution of the 3-hydroxypropionate bicycle and**

1054 **recent transfer of anoxygenic photosynthesis into the Chloroflexi.** *Proc Natl Acad*

1055 *Sci U S A* 2017, **114**:10749-10754.

1056 69. Berg IA, Kockelkorn D, Buckel W, Fuchs G: **A 3-hydroxypropionate/4-**

1057 **hydroxybutyrate autotrophic carbon dioxide assimilation pathway in Archaea.**

1058 *Science* 2007, **318**:1782-1786.

1059 70. Pester M, Schleper C, Wagner M: **The Thaumarchaeota: an emerging view of their**

1060 **phylogeny and ecophysiology.** *Curr Opin Microbiol* 2011, **14**:300-306.

1061 71. Kanehisa M, Sato Y, Morishima K: **BlastKOALA and GhostKOALA: KEGG tools**

1062 **for functional characterization of genome and metagenome sequences.** *J Mol Biol*

1063 2016, **428**:726-731.

1064 72. Moriya Y, Itoh M, Okuda S, Yoshizawa AC, Kanehisa M: **KAAS: an automatic**

1065 **genome annotation and pathway reconstruction server.** *Nucleic Acids Res* 2007,

1066 **35**:W182-W185.

1067 73. Huerta-Cepas J, Szklarczyk D, Heller D, Hernandez-Plaza A, Forslund SK, Cook H, Mende DR, Letunic I, Rattei T, Jensen LJ, et al: **eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090**

1068 **organisms and 2502 viruses.** *Nucleic Acids Res* 2019, **47**:D309-D314.

1069 74. Olson DL, Delen D: *Advanced data mining techniques.* Berlin, Heidelberg: Springer-

1070 Verlag Berlin Heidelberg 2008.

1071

1072

1073

1074

1075

1076

1077 75. Glass JB, Ranjan P, Kretz CB, Nunn BL, Johnson AM, McManus J, Stewart FJ: 1078 **Adaptations of Atribacteria to life in methane hydrates: hot traits for cold life.** 1079 *bioRxiv* 2019, **1**:536078.

1080 76. Anantharaman K, Breier JA, Dick GJ: **Metagenomic resolution of microbial 1081 functions in deep-sea hydrothermal plumes across the Eastern Lau Spreading 1082 Center.** *ISME J* 2015, **10**:225.

1083 77. Anantharaman K, Duhaime MB, Breier JA, Wendt K, Toner BM, Dick GJ: **Sulfur 1084 Oxidation Genes in Diverse Deep-Sea Viruses.** *Science* 2014, **344**:757-760.

1085 78. Zhou Z, Tran PQ, Kieft K, Anantharaman K: **Genome diversification in globally 1086 distributed novel marine Proteobacteria is linked to environmental adaptation.** 1087 *ISME J* 2020, **14**:2060-2077.

1088 79. Baker BJ, Saw JH, Lind AE, Lazar CS, Hinrichs K-U, Teske AP, Ettema TJG: **Genomic 1089 inference of the metabolism of cosmopolitan subsurface Archaea, Hadesarchaea.** 1090 *Nat Microbiol* 2016, **1**:16002.

1091 80. Madigan MT, John M. Martinko, Kelly S. Bender, Daniel H. Buckley, and David Allan 1092 Stahl: *Brock Biology of Microorganisms*. Fourteenth edition edn. Boston: Pearson; 1093 2015.

1094 81. Wolf PG, Cowley ES, Breister A, Matatov S, Lucio L, Polak P, Ridlon JM, Gaskins 1095 HR, Anantharaman K: **Diversity and distribution of sulfur metabolism in the 1096 human gut microbiome and its association with colorectal cancer.** *bioRxiv* 1097 2021:2021.2007.2001.450790.

1098 82. Huson DH, Beier S, Flade I, Górska A, El-Hadidi M, Mitra S, Ruscheweyh H-J, Tappu 1099 R: **MEGAN Community Edition - Interactive Exploration and Analysis of Large- 1100 Scale Microbiome Sequencing Data.** *PLoS Comput Biol* 2016, **12**:e1004957.

1101

1102 **FIGURE AND TABLE LEGENDS**

1103 **Figure 1. An outline of the workflow of METABOLIC.** Detailed instructions are available
1104 at <https://github.com/AnantharamanLab/METABOLIC/wiki>. METABOLIC-G workflow is
1105 specifically shown in the blue box and METABOLIC-C workflow is shown in the green square.

1106
1107 **Figure 2. Summary scheme of biogeochemical cycling processes at the community scale.**
1108 Each arrow represents a single transformation/step within a cycle. Labels above each arrow are
1109 (from top to bottom): step number and reaction, number of genomes that can conduct these
1110 reactions, metagenomic coverage of genomes (represented as a percentage within the
1111 community) that can conduct these reactions. The numbers in brackets next to the nitrogen or
1112 sulfur-containing compounds are chemical states of the nitrogen or sulfur atoms in these
1113 compounds.

1114
1115 **Figure 3. Schematic figure of sequential metabolic transformations. (A) the sequential**
1116 **transformation of inorganic compounds; (B) the sequential transformation of organic**
1117 **compounds.** X-axes describe individual sequential transformations indicated by letters. The
1118 two panels describe the number of genomes and genome coverage (represented as a percentage
1119 within the community) of organisms that are involved in certain sequential metabolic
1120 transformations. The deep-sea hydrothermal vent dataset was used for these analyses.

1121
1122 **Figure 4. Functional network showing connections between different functions in the**
1123 **microbial community.** Nodes represent individual steps in biogeochemical cycles; edges
1124 connecting two given nodes represent the functional connections between nodes, which are
1125 enabled by organisms that can conduct both biogeochemical processes/steps. The size of the
1126 node was depicted according to the degree (number of connections to each node). The thickness
1127 of the edge was depicted according to the average gene coverage values of the two connected
1128 biogeochemical cycling steps – for example, thiosulfate oxidation and organic carbon
1129 oxidation. The color of the edges was assigned based on the taxonomy of the represented
1130 genome. The deep-sea hydrothermal vent dataset was used for these analyses.

1131
1132 **Figure 5. Description, calculation, and result table of MW-scores. (A) The calculation**
1133 **method for the MW-score within a community based on a given metagenomic dataset.**
1134 Each circle stands for a genome within the community, and the adjacent bar stands for its
1135 genome coverage within the community. The coverage values of encoded genes for all
1136 functions were summed up as the denominator, and the coverage value of encoded genes for
1137 each function was used as the numerator, and the MW-score was calculated accordingly for
1138 each function. **(B) The resulting table of MW-score for the deep-sea hydrothermal vent**
1139 **metagenomic dataset.** MW-score for each function was given in a separated column, and the
1140 rest of the table indicates the contribution percentage to each MW-score of the genomes
1141 grouped in each phylum. The MW-score of “N-S-07:Nitrous oxide reduction” was not exactly
1142 0 but rounded to 0 due to the original number being less than 0.05. Additionally, contribution
1143 percentages were also rounded to only retain one digit after the decimal points; consequently,
1144 the sum contribution percentages for some functions slightly deviate from 100%.

1146 **Figure 6. Metabolic Sankey diagram representing the contributions of microbial genomes**
1147 **to individual metabolic and biogeochemical processes and entire elemental cycles.**

1148 Microbial genomes are represented at the phylum-level resolution. The three columns from left
1149 to right represent taxonomic groups scaled by the number of genomes, the contribution to each
1150 metabolic function by microbial groups calculated based on genome coverage, and the
1151 contribution to each functional category/biogeochemical cycle. The colors were assigned based
1152 on the taxonomy of the microbial groups. The deep-sea hydrothermal vent dataset was used for
1153 these analyses.

1154
1155 **Figure 7. Comparison of METABOLIC with other software packages and online servers.**
1156 **(A) Comparison of the workflows and services, (B) Comparison of performance of protein**
1157 **prediction for two representative genomes, *Pseudomonas aeruginosa* PAO1, and**
1158 ***Escherichia coli* O157H7 str. sakai.**

1159
1160 **Figure 8. Community metabolism comparison based on MW-scores. (A) Comparison**
1161 **between terrestrial subsurface (left red bars) and marine subsurface (right blue bars); (B)**
1162 **Comparison between deep-sea hydrothermal vent (left red bars) and freshwater lake**
1163 **(right blue bars).** MW-scores were calculated as gene coverage fractions for individual
1164 metabolic functions. Functions with MW-scores in both environments as zero were removed
1165 from each panel, e.g., N-S-02:Ammonia oxidation, N-S-09:Anammox, S-S-02:Sulfur
1166 reduction, and S-S-06:Sulfite reduction in Panel (A), and C-S-07:Methanogenesis, N-S-01:N₂
1167 fixation, N-S-09:Anammox, S-S-02:Sulfur reduction, and S-S-06:Sulfite reduction in Panel
1168 (B). Details for MW-score and each microbial group contribution refer to [Supplementary](#)
1169 [Dataset S3](#).

1170
1171 **Figure 9. Cell metabolism diagrams of two microbial genomes. (A) cell metabolism**
1172 **diagram of Hadesarchaea archaeon 1244-C3-H4-B1 (B) cell metabolism diagram of**
1173 **Nitrospirae bacteria M_DeepCast_50m_m2_151.** The absent functional
1174 pathways/complexes were labeled with dash lines.

1175
1176 **Figure 10. Presence/Absence map of human microbiome metabolisms of a colorectal**
1177 **cancer (CRC) patient and a healthy control gut sample.** The heatmap has summarized 189
1178 horizontal entries (189 lines) based on 139 key functional gene families that covered 10
1179 function categories. Purple cells indicate presence and gray cells indicate absence. Detailed
1180 KEGG KO identifier IDs and protein information for each function category were described in
1181 [Supplementary Dataset S2](#).

1182 **Table 1.** The carbon fixation metabolic traits of 15 tested bacterial and archaeal genomes
 1183 predicted by both METABOLIC and KEGG genome database
 1184

Accession ID	Organism	KEGG Organism Code	Group	METABOLIC result		KEGG genome pathway	
				3HP cycle	3HP/4HB cycle	3HP cycle	3HP/4HB cycle
GCA_000011905.1	<i>Dehalococcoides mccartyi</i> 195	det	Chloroflexi	Absent	Absent	Absent	Absent
GCA_000017805.1	<i>Roseiflexus castenholzii</i> DSM 13941	rca	Chloroflexi	Present	Absent	Present	Absent
GCA_000018865.1	<i>Chloroflexus aurantiacus</i> J-10-fl	cau	Chloroflexi	Present	Absent	Present	Absent
GCA_000021685.1	<i>Thermomicrobium roseum</i> DSM 5159	tro	Chloroflexi	Absent	Absent	Absent	Absent
GCA_000021945.1	<i>Chloroflexus aggregans</i> DSM 9485	cag	Chloroflexi	Present	Absent	Present	Absent
GCA_000299395.1	<i>Nitrosopumilus sediminis</i> AR2	nir	Thaumarchaeota	Absent	Present	Absent	Present
GCA_000698785.1	<i>Nitrososphaera viennensis</i> EN76	nvn	Thaumarchaeota	Absent	Present	Absent	Present
GCA_000875775.1	<i>Nitrosopumilus piranensis</i> D3C	nid	Thaumarchaeota	Absent	Present	Absent	Present
GCA_000812185.1	<i>Nitrosopelagicus brevis</i> CN25	nbv	Thaumarchaeota	Absent	Present	Absent	Present
GCA_900696045.1	<i>Nitrosocosmicus franklandus</i> NFRAN1	nfn	Thaumarchaeota	Absent	Present	Absent	Present
GCA_000015145.1	<i>Hyperthermus butylicus</i> DSM 5456	hbu	Crenarchaeota	Absent	Absent	Absent	Absent
GCA_000017945.1	<i>Caldisphaera lagunensis</i> DSM 15908	clg	Crenarchaeota	Absent	Present	Absent	Present
GCA_000148385.1	<i>Vulcanisaeta distributa</i> DSM 14429	vdi	Crenarchaeota	Absent	Absent	Absent	Absent
GCA_000193375.1	<i>Thermoproteus uzoniensis</i> 768-20	tuz	Crenarchaeota	Absent	Present	Absent	Present
GCA_003431325.1	<i>Acidilobus</i> sp. 7A	acia	Crenarchaeota	Absent	Absent	Absent	Absent

1185
 1186