# METABOLIC: A scalable high-throughput metabolic and biogeochemical functional trait profiler based on microbial genomes

Zhichao Zhou[1], Patricia Tran[1], Yang Liu[2], Kristopher Kieft[1], Karthik Anantharaman[1,*]

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

[2] Institute for Advanced Study, Shenzhen University, Shenzhen, Guangdong Province, 518060, China

*To whom correspondence should be addressed.

## Abstract

**Summary:** Microbial metabolism mediates fundamental transformations of chemistry and energy that drive biogeochemical cycles on our planet. Increasingly**,** we can read genomic blueprints of microorganisms, decipher their functional capacities and activities, and reconstruct their roles in biogeochemical processes using omic-based techniques such as metagenomics. Currently available tools for analyses of genomic data can annotate and depict metabolic functions to some extent, but they are not comprehensive. No standardized approaches are currently available for bioinformatic validation of metabolic predictions and identifying contributions of microorganisms and genes to biogeochemical cycles. Here we present METABOLIC (**MET**abolic **A**nd **B**ioge**O**chemistry ana**L**yses **I**n mi**C**robes), a scalable metabolic and biogeochemical functional trait profiler to comprehensively study microbial metabolism using genome data. METABOLIC uses metagenome-assembled (MAG), single-cell (SAG), or isolate genomes as input, annotates and processes genomes for identification and characterization of metabolism markers using KEGG and curated custom protein HMM databases, and applies motif confirmation of biochemically validated conserved residues in proteins. The output report includes functionally important HMM hit tables, protein collections for downstream analysis, tables (KEGG modules) and diagrams representing metabolic pathways for individual genomes, and a summary figure representing selected biogeochemical cycling processes on a community scale. We expect that METABOLIC will facilitate the study of genome-informed microbial metabolism and biogeochemistry and transform our understanding of environmental microbiomes.

**Availability and implementation:** METABOLIC is available on github: https://github.com/AnantharamanLab/METABOLIC.

**Contact:** karthik@bact.wisc.edu

## 1. Introduction

43

44 Microbially-mediated biogeochemical processes serve as important driving forces for

45 transformation and cycling of elements, energy, and matter among the lithosphere,

46 atmosphere, hydrosphere and biosphere (Madsen, 2011). Metagenomics and single-cell

47 genomics have transformed the field of microbial ecology by revealing a rich diversity

48 of microorganisms from diverse settings, including terrestrial and marine environments

49 and human body (Anantharaman, et al., 2016; Dombrowski, et al., 2018; Parks, et al.,

50 2017; Pasolli, et al., 2019). These approaches can provide an unbiased and insightful

51 view into microorganisms mediating and contributing to the biogeochemical activities

52 at a number of scales ranging from individual organisms to communities

53 (Anantharaman, et al., 2016; Bowers, et al., 2017; Hug, et al., 2016; Parks, et al., 2017).

54 Prediction of microbial metabolism relies on annotation of protein function for

55 microorganisms using a number of established databases, e.g., KEGG (Kanehisa and

56 Goto, 2000), MetaCyc (Caspi, et al., 2006), Pfam (Finn, et al., 2014), TIGRfam

57 (Selengut, et al., 2007), SEED (Overbeek, et al., 2013), and eggNOG (Huerta-Cepas,

58 et al., 2016). However, these results are often highly detailed. Obtaining a functional

59 profile and identifying metabolic pathways in a microbial genome can involve manual

60 inspection of thousands of genes. Interpreting, organizing and visualizing such datasets

61 remains a challenge and is often untenable, and there is a critical need for a tool to

62 identify and validate the presence of metabolic pathways and genes of biogeochemical

63 function in a user-friendly manner. Such a tool would also allow standardization and

64 easy integration of genome-informed metabolism into biogeochemical models which

65 currently rely primarily on physico-chemical data and treats microorganisms as black

66 boxes.

67

68 Here we present the software METABOLIC, a tool to profile metabolic and

69 biogeochemical functional traits based on microbial genomes. It integrates annotation

70 of proteins using KEGG (Kanehisa and Goto, 2000), TIGRfam (Selengut, et al., 2007),

71 Pfam (Finn, et al., 2014), and custom HMM databases (Anantharaman, et al., 2016),

72    incorporates a motif validation step to accurately identify proteins based on prior

73    biochemical validation, determines presence or absence of metabolic pathways based

74    on KEGG modules, and produces user-friendly outputs in the form of tables and figures

75    including a summary of biogeochemically-relevant pathways and their abundance for

76    individual genomes and at the community scale.

77

## 2. Methods

79    METABOLIC is written in Perl and R and is expected to run in Unix/Linux and MacOS.

80    The prerequisites are described on METABOLIC's GitHub page

81    (https://github.com/AnantharamanLab/METABOLIC). The input folder requires

82    microbial genome sequences in FASTA format and an optional set of metagenomic

83    reads in which were used to reconstruct those genomes (Supplementary Figure S1).

84    Genomic sequences are annotated by Prodigal (Hyatt, et al., 2010), or a user can provide

85    self-annotated proteins (with extensions of ".faa") in order to facilitate incorporation

86    into existing pipelines. Proteins will be queried against hidden Markov model (HMM)

87    databases using hmmsearch implemented within HMMER (Finn, et al., 2011) which

88    implements methods to detect remote homologs as sensitively and efficiently as

89    possible. The HMM databases include Kofam prokaryotic (KEGG) (Aramaki, et al.,

90    2019), TIGRfam (Selengut, et al., 2007), Pfam (Finn, et al., 2014) and custom metabolic

91    HMM files (Anantharaman, et al., 2016). The cutoff threshold values for HMM

92    databases were used as follows: Kofam - Kofam suggested values;

93    TIGRfam/Pfam/Custom databases - Manually curated by adjusting noise cutoffs (NC)

94    and trusted cutoffs (TC) to avoid potential false positive hits; detailed curation methods

95    are described previously (Anantharaman, et al., 2016).

96

97    To computationally validate protein hits and avoid false positives, we have introduced

98    a motif validation step that including comparison of protein motifs against a manually

99    curated set of highly conserved residues in important proteins. As an example, DsrC

100    (sulfite reductase subunit C) and TusE (tRNA 2-thiouridine synthesizing protein E) are

101 similar proteins that are routinely misannotated. Both are assigned to the family

102 KO:K11179 in the KEGG database. To avoid assigning TusE as a sulfite reductase, we

103 identified a specific motif for DsrC but not TusE (GPXKXXCXXXGXPXPXXCX"

104 where "X" stands for any amino acid) (Venceslau, et al., 2014). We use these specific

105 motifs to filter for proteins which have high sequence similarity but functionally

106 divergent homologs.

107

108 The software output integrates the presence and absence of genes from the outputs of

109 individual HMM runs and relates them to microbial functional traits. Individual KEGG

110 annotations are inferred in the context of KEGG modules for better interpretation of

111 metabolic pathways. A KEGG module is a collection of manually defined functional

112 units. A module is comprised of multiple steps with each step representing a distinct

113 metabolic function. Since genomes can often have incomplete metabolic pathways, we

114 determine the completeness of specific metabolic pathways by parsing KEGG module

115 IDs. A user-defined cutoff is used to estimate the completeness of a given module (the

116 default value is 75%), which is then used to produce KEGG module presence/absence

117 table. All modules exceeding the cutoff are determined to be complete in the given

118 genome. Outputs consist of four different results that are reported in an Excel

119 spreadsheet (Supplementary Figure S2). These contain details of HMM hits

120 (Supplementary Figure S2A), presence/absence of functional traits (Supplementary

121 Figure S2B), presence/absence of KEGG modules (Supplementary Figure S2C), and

122 presence/absence of KEGG module steps (Supplementary Figure S2D). Each collection

123 of HMM hits can be extracted from input genomes for the downstream phylogenetic

124 analysis. A detailed workflow of METABOLIC is available in Supplementary Figure

125 S1.

126

127 To visualize pathways of biogeochemical importance, the software draws schematic

128 profiles for nitrogen, carbon, sulfur and other element cycles for each genome. A

129 summary schematic diagram at the scale of a microbial community integrates results

130 from all genomes from a given dataset (Figure 1) and includes computed abundances

5

131      for each step in a biogeochemical cycle if the metagenomic reads datasets are provided.

## 3. Results

METABOLIC has been successfully applied on a metagenomic dataset which includes 98 MAGs from a deep-sea hydrothermal plume at Guaymas Basin in the Pacific Ocean, and two sets of metagenomic reads (that are subsets of original reads with 10 million read numbers for each pair comprising ~10% of the total reads). The total run time was ~8 hours using 25 CPU threads in a Linux version 4.15.0-48-generic server (Ubuntu v5.4.0). The resulting summary scheme on various biogeochemical cycling processes reflects the pattern on a community scale (Figure 1) (Supplementary Data S1 contains tables and figures from the METABOLIC output).

In order to test the accuracy of the results predicted by METABOLIC, we picked 15 bacterial and archaeal genomes from *Chloroflexi*, *Thaumarchaeota*, and *Crenarchaeota* which are reported to have 3 Hydroxypropionate cycle (3HP) or 3-hydroxypropionate/4-hydroxybutyrate cycle (3HP/4HB) for carbon fixation. METABOLIC predicts results in line with KEGG genome database annotations and can also be visualized with the KEGG Mapper (Supplementary Table S1). Our predictions are also in accord with biochemical evidence of the existence of corresponding carbon fixation pathways in each microbial group: only organisms from the phylum *Chloroflexi* are known to possess the 3HP pathway and 3HP/4HB pathway could only be detected in *Crenarchaeota* and *Thaumarchaeota* (Supplementary Table S1 and references therein). These results suggest that METABOLIC can provide accurate annotations and genomic profiles of metabolism and serve as a good functional predictor for microbial genomes at the individual and community scales.

## Funding

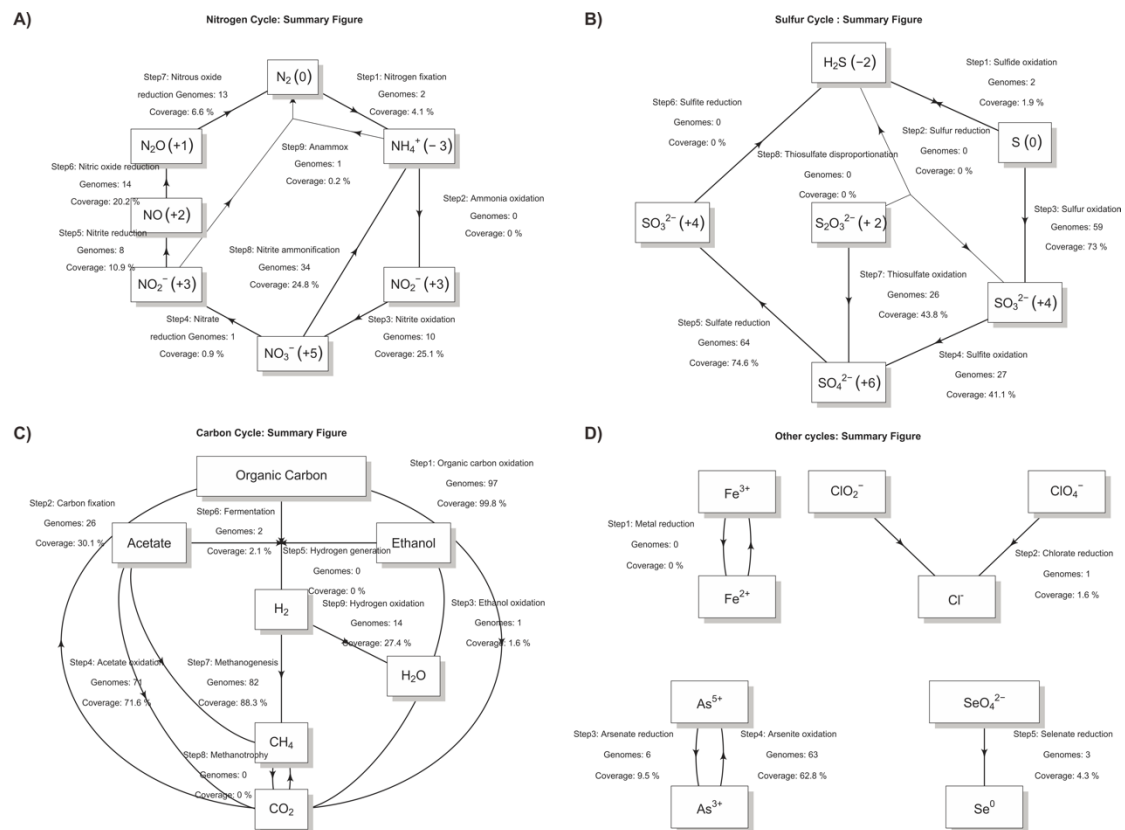161     *Conflict of interest*: none declared.

162

163

164

165

**Figure 1. The summary scheme of biogeochemical cycling processes on a community scale.** Above each arrow (which represent each step within a cycle) there are three lines. The first one indicates the step name and the reaction, the second one indicates the number of genomes that acquire these reactions, the third one indicates the percentage of metagenomic coverage on each step.

171

172

# References

Anantharaman, K.*, et al.* (2016) Thousands of microbial genomes shed light on interconnected biogeochemical processes in an aquifer system, *Nat. Commun.*, **7**, 13219.

Aramaki, T.*, et al.* (2019) KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive score threshold, *bioRxiv*, 602110.

Bowers, R.M.*, et al.* (2017) Minimum information about a single amplified genome (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and archaea, *Nat. Biotechnol.*, **35**, 725.

Caspi, R.*, et al.* (2006) MetaCyc: a multiorganism database of metabolic pathways and enzymes, *Nucleic Acids Res.*, **34**, D511-516.

Dombrowski, N., Teske, A.P. and Baker, B.J. (2018) Expansive microbial metabolic versatility and biodiversity in dynamic Guaymas Basin hydrothermal sediments, *Nat. Commun.*, **9**, 4999.

Finn, R.D.*, et al.* (2014) Pfam: the protein families database, *Nucleic Acids Res.*, **42**, D222-230.

Finn, R.D., Clements, J. and Eddy, S.R. (2011) HMMER web server: interactive sequence similarity searching, *Nucleic Acids Res.*, **39**, W29-37.

Huerta-Cepas, J.*, et al.* (2016) eggNOG 4.5: a hierarchical orthology framework with improved functional annotations for eukaryotic, prokaryotic and viral sequences, *Nucleic Acids Res.*, **44**, D286-D293.

Hug, L.A.*, et al.* (2016) A new view of the tree of life, *Nat. Microbiol.*, **1**, 16048.

Hyatt, D.*, et al.* (2010) Prodigal: prokaryotic gene recognition and translation initiation site identification, *BMC Bioinformatics*, **11**, 119.

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

Madsen, E.L. (2011) Microorganisms and their roles in fundamental biogeochemical cycles, *Curr. Opin. Biotechnol.*, **22**, 456-464.

Overbeek, R.*, et al.* (2013) The SEED and the Rapid Annotation of microbial genomes using Subsystems Technology (RAST), *Nucleic Acids Res.*, **42**, D206-D214.

Parks, D.H.*, et al.* (2017) Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life, *Nat. Microbiol.*, **2**, 1533-1542.

Pasolli, E.*, et al.* (2019) Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle, *Cell*, **176**, 649-662 e620.

Selengut, J.D.*, et al.* (2007) TIGRFAMs and Genome Properties: tools for the assignment of molecular function and biological process in prokaryotic genomes, *Nucleic Acids Res.*, **35**, D260-264.

Venceslau, S.S.*, et al.* (2014) The "bacterial heterodisulfide" DsrC is a key protein in dissimilatory sulfur metabolism, *Biochim. Biophys. Acta*, **1837**, 1148-1164.