1 **TbasCO: Trait-based Comparative 'Omics Identifies Ecosystem-Level and Niche-**
2 **Differentiating Adaptations of an Engineered Microbiome**

3

4 McDaniel, E.A.[1,2#*], van Steenbrugge, J.J.M[3,4,5#+], Noguera, D.R.[6], McMahon, K.D.[1,6],
5 Raaijmakers, J.M.[4,7], Medema, M.H.[3,7], Oyserman, B.O.[3,4+]

6

7

8 [1] Department of Bacteriology, University of Wisconsin – Madison, Madison, WI, USA

9 [2] Microbiology Doctoral Training Program, University of Wisconsin - Madison, Madison, WI, USA

10 [3] Bioinformatics Group, Wageningen University and Research, Wageningen, The Netherlands

11 [4] Microbial Ecology, Netherlands Institute of Ecological Research, Wageningen, The

12 Netherlands

13 [5] Laboratory of Nematology, Wageningen University, Wageningen, The Netherlands

14 [6] Department of Civil and Environmental Engineering, University of Wisconsin – Madison,

15 Madison WI USA

16 [7]Institute of Biology, Leiden University, Leiden, Netherlands

17 [#] = contributed equally

18 [+] = corresponding author

19

20 * Current address:

21 Department of Microbiology and Immunology, University of British Columbia, Vancouver, CA

22

23 Corresponding authors:

24 Elizabeth McDaniel elizabethmcd93@gmail.com

25 Joris van Steenbrugge jorisvansteebrugge@gmail.com

26 Ben Oyserman BenOyserman@gmail.com

27

28

29    **ABSTRACT**

30    A grand challenge in microbial ecology is disentangling the traits of individual

31    populations within complex communities. Various cultivation-independent approaches

32    have been used to infer traits based on the presence of marker genes. However, marker

33    genes are not linked to traits with complete fidelity, nor do they capture important

34    attributes, such as the timing of expression or coordination among traits. To address this,

35    we present an approach for assessing the trait landscape of microbial communities by

36    statistically defining a trait attribute as shared transcriptional pattern across multiple

37    organisms. Leveraging the KEGG pathway database as a trait library and the Enhanced

38    Biological Phosphorus Removal (EBPR) model microbial ecosystem, we demonstrate

39    that a majority (65%) of traits present in 10 or more genomes have niche-differentiating

40    expression attributes. For example, while 14 genomes containing the high-affinity

41    phosphorus transporter *pstABCS* display a canonical attribute (e.g. up-regulation under

42    phosphorus starvation), we identified another attribute shared by 11 genomes where

43    transcription was highest under high phosphorus conditions. Taken together, we provide

44    a novel framework for revealing hidden metabolic versatility when investigating genomic

45    data alone by assigning trait-attributes through genome-resolved time-series

46    metatranscriptomics.

47

48

49

50

51

## INTRODUCTION

A longstanding cornerstone of deterministic ecological theory is that the environment selects for traits. Traits may be defined as any physiological, morphological, or genomic signature that affects the fitness or function of an individual [1]. Trait-based approaches have become indispensable in macroecological systems to describe fitness trade-offs and the effects of biodiversity on ecosystem functioning [2–5]. Recently, trait-based frameworks have been proposed as an alternative to taxonomy-based methods for describing microbial ecosystem processes [6, 7]. Connecting microbial traits and their phylogenetic distributions to ecosystem performance can provide powerful insights into the ecological and evolutionary dynamics underpinning community assembly, microbial biogeography, and organismal responses to changes in the environment [8–10]. Additionally, pinpointing the organismal distribution of traits and the selective pressures that enrich them may be leveraged to reproducibly and rationally engineer stable, functionally redundant ecosystems [11–15]. However, applying trait-based approaches to microbial communities is challenging due to the difficulty in identifying and measuring relevant ecological traits for a given ecosystem [16].

High-throughput sequencing technologies and multi-omics techniques have been used to describe the diversity, activity, and functional potential of uncultivated microbial lineages [17–20]. Improvements in bioinformatics algorithms, and in particular metagenomic binning methods, have allowed for genome-resolved investigations of microbial communities rather than gene-based analyses of assembled contigs [21]. These (meta) genomes are subsequently leveraged to detect the presence of key genes or pathways and predict specific traits of the whole community [22, 23]. Integrating

75  metatranscriptomics data addresses a key limitation, as expression patterns better reflect

76  the actual functional dynamics of a trait compared to gene presence alone. Here, we

77  present TbasCO, a software package and statistical framework for *Trait-bas*ed

78  *C*omparative 'Omics to identify expression attributes. We adopt the terminology *attribute*

79  as a hierarchically structured feature of a trait and assert that statistically similar

80  transcriptional patterns of traits across multiple organisms be treated as *attributes* (Figure

81  1)*.* In this manner, the identification of expression-based *attributes* provides a high-

82  throughput and intuitive framework for extending trait-based methods to time-series

83  expression patterns in microbial communities. We implement this trait-based approach to

84  classify transcriptional attributes in a microbial community performing Enhanced

85  Biological Phosphorus Removal (EBPR), a globally important biotechnological process

86  implemented in numerous wastewater treatment plants (WWTPs).

87       The fundamental feature of the engineered EBPR ecosystem is the decoupled and

88  cyclic availability of an external carbon source and terminal electron acceptor. This cycling

89  is often referred to as "feast-famine" conditions and provides a strong selective pressure

90  for traits such as polymer cycling. Accumulation of intracellular polyphosphate through

91  cyclic anaerobic-aerobic conditions ultimately results in net phosphorus removal and

92  accomplishes the EBPR process [24, 25]. One of the most well-studied polyphosphate

93  accumulating organisms (PAOs) belongs to the uncultivated bacterial lineage

94  '*Candidatus* Accumulibacter phosphatis' (hereby referred to as Accumulibacter) [24, 26].

95  Numerous genome-resolved 'omics methods have been used to investigate the

96  physiology and regulation of this model PAO enriched in engineered lab-scale enrichment

97  bioreactor systems [27–34]. However, novel and putative PAOs have been discovered

98 that remove phosphorus without exhibiting the hallmark traits of Accumulibacter [35–39].

99 Additionally, although these lab-scale systems are designed to specifically enrich for

100 Accumulibacter, a diverse "flanking community" persists in these environments [27], and

101 their ecological roles have largely remained unexplored. As a result, the general

102 adaptations of microbial lineages inhabiting the EBPR community are not well

103 understood. Using genome-resolved metagenomics and metatranscriptomics, we

104 assembled 66 species-representative genomes spanning several significant EBPR

105 lineages and identified the distribution of expression-based attributes. Using our novel

106 trait-based comparative 'omics approach, we show that while some expression attributes

107 are distributed in few genomes, many are redundant and shared across many lineages.

108 Furthermore, we find that a majority of core traits (as defined by the presence of marker

109 genes) have multiple attributes, suggesting that identifying niche-differentiating

110 expression attributes may be used to reveal a large hidden metabolic versatility when

111 investigating genomic data alone.

112

113 **MATERIALS AND METHODS**

114 **Metagenomic Assembly, Annotation, and Metatranscriptomic Mapping**

115 Three metagenomes sampled from an EBPR bioreactor with linked time-series

116 metatranscriptomics data [40] were collected for metagenomic sequencing and

117 assembled into 66 species-representative bins as described in the Supplemental

118 Methods. All bins are greater than 75% complete and contain less than 10%

119 contamination, with a large majority (44/66) >95% complete and <5% redundant as

120 calculated by CheckM [41] (Table 1). Each bin was functionally annotated using the

121     KEGG database through an HMM-based approach under KEGG release 93.0 using the

122     command-line KofamKOALA pipeline [42, 43], selecting annotations that were significant

123     hits above the specific HMM threshold. This resulted in 117,657 total annotations with

124     5,228 unique annotations. We used a metatranscriptomic dataset of six timepoints

125     collected over a single EBPR cycle from Oyserman et al. 2016 [40], with three timepoints

126     from the anaerobic phase and three from the aerobic phase. Raw reads were quality

127     filtered using BBtools suite v38.07 [44]  and ribosomal rRNA was removed from each

128     sample using SortMeRNA [45]. Reads from each sample were mapped against the

129     concatenated set of open reading frames from all 66 bins using kallisto v0.44.0 and

130     parsed using the R package tximport [46, 47].

131     **TbasCO Method Implementation**

132     The TbasCO package identifies expression-based attributes of predefined traits

133     using time-series (meta)transcriptomics data (Figure 1). In general, traits are defined by

134     the presence of a pathway or other collection of genes from an externally provided

135     database. A weighted distance metric between expression patterns for all genes that

136     define a trait is calculated, and statistically significant similarity is determined based on

137     the background distribution of a trait of equal size. Thereby, two or more organisms with

138     a statistically similar expression pattern for a trait share an *attribute.*

139     ***Input and Preprocessing***

140     The input that is accepted by TbasCO is a table of RNAseq counts in csv format.

141     Each row is treated as gene that has columns for the gene/locus name, counts per

142     sample, the genome the gene belongs to, and the KEGG Orthology (KO) identifier. The

143     RNAseq counts table may be provided pre-normalized or can be normalized by the

144    program. The default normalization method is designed to minimize compositional bias in

145    the differential abundance and activity of constituent populations in metatranscriptomics

146    studies. Raw RNA expression counts are therefore normalized by genomic bin and

147    sample [40]. These normalization factors are then applied to each sample for each bin

148    individually. Alternatively, custom normalization methods may be implemented. After

149    normalization, a pruning step is introduced to filter genes that have zero counts or a mean

150    absolute deviation of less than one. To make inter-organismal comparisons of the relative

151    contribution of a gene to total measured organismal RNA, an additional statistic is

152    calculated ranking the expression counts from each sample from highest to lowest. The

153    ranks for each sample are then normalized by dividing them by the maximum rank value

154    in that sample. This normalization is applied to make ranks comparable between

155    organisms with different genome sizes.

156        To assess the statistical significance of the calculated distances between the

157    expression patterns of all genes within a trait, random background distributions are

158    created for 1) individual genes and 2) traits of N genes. For individual genes, three

159    different distributions were calculated, based on the distances between randomly

160    sampled open reading frames, randomly sampled genes with an annotation (but not

161    necessarily the same annotation), and randomly sampled genes with the same

162    annotation. The background distribution for a trait of N genes is based on the distances

163    between randomly composed sets of genes. For each gene pair, two distances metrics

164    are calculated, the Pearson Correlation (PC) and the Normalized Rank Euclidean

165    Distance (NRED). In practice, it is often found that a certain annotation is assigned to

166    multiple genes in the same genome. If this occurs, there is an option to use either a

167     random selection, or the highest scoring pair. In the latter case, a correction for multiple

168     testing is implemented. This process is repeated N-times, where N equals the number of

169     genes in any given trait. The background distribution for traits is determined by first

170     randomly sampling two genomes, identifying the overlap in annotations, and finally

171     artificially defining a trait containing N annotations. For each annotation in the trait, the

172     distances are calculated between genome A and genome B, as described in the previous

173     section. As modules vary in size, this process is repeated for traits of different sizes.

174     ***Identifying Attributes***

175     TbasCO provides both a cluster-based and pair-wise approach to identify

176     attributes. In both methods, the distance between expression patterns of a trait between

177     two genomes is first calculated based on a composite Z score of the PC and NRED for

178     each gene composing the trait. In the cluster-based analysis, the distances are

179     subsequently clustered using the Louvain clustering algorithm to identify trait attributes.

180     To determine if an attribute is significantly similar or not, a one-sided T-test between the

181     attribute and the random background distribution of traits is conducted. This is done for

182     both cluster-based and model-based comparisons. Many traits are complex and

183     represented in databases such as KEGG by numerous alternative routes. To deal with

184     this complexity, each pathway is expanded into the Disjunctive Normative Form (DNF).

185     Due to the extremely high number of DNFs for some traits, attributes are pruned based

186     on a strict requirement of 100% completion.

187     ***Distance Calculations***

188     To determine the similarity in expression patterns between genes, two distance

189     metrics are calculated: the PC between RNAseq counts across samples, and the NRED,

190   where ranks are a measure of relative abundance of RNA in each sample, normalized

191   the abundance of RNA in the corresponding genome. These distance scores are

192   converted to Z scores using a background distribution of distances between randomly

193   sampled genes as previously described. To determine statistically significant similarities

194   between the expression patterns of a trait between two genomes, a composite distance

195   score is calculated based on the distance between genes in two different genomes. For

196   each of these genes the PC and NRED are calculated and transformed to Z scores and

197   combined as (-1*PC + NRED). The distance of the trait between two genomes is defined

198   as the average of these composite distance scores, and then normalized by the Jaccard

199   distance between these genomes.

200   $$(-PC + NRED) * (1 - dJ)$$

201   ***Statistical Assessment of Trait Attributes***

202   In both model-based and pair-wise approaches, the distance is first calculated

203   between expression patterns of a trait between two genomes based on the composite Z

204   score of the PC and NRED for each gene composing the trait. In the clustering-based

205   analysis, the distances are subsequently clustered using the Louvain clustering algorithm

206   to identify trait-attributes. To determine if attributes are significantly similar, a one-sided

207   T-test is conducted between the attribute and a background distribution of randomly

208   sampled traits with the same number of genes. To derive the random background

209   distributions, multiple distributions are calculated ranging in gene numbers from the

210   smallest trait to the largest trait in the dataset as described previously. For each

211   background distribution, N (default: 10,000) traits are randomly composed. The distances

212   between these artificial traits are calculated in the same way as for the actual traits. In

213     addition to a statistical pruning step, the attributes are pruned based on a strict

214     requirement of 100% completion of each DNF module. A benchmarking analysis to

215     examine the effects of different parameters was conducted to determine their influence

216     on the number of attributes identified and may be found in the supplementary materials

217     (Supplementary Table 1, Supplementary Figures 2-4).

218     **Data and Code Availability**

219         All supplementary files and figures including functional annotations and

220     transcriptome count files are available at https://figshare.com/projects/EBPR_Trait-

221     Based_Comparative_Omics/90437. All 64 flanking genomes have been deposited in

222     NCBI at Bioproject PRJNA714686. The remaining two reassembled Accumulibacter

223     genomes have not been deposited in NCBI to not confuse between the original CAPIA

224     and CAPIIA assemblies [27, 28]. These contemporary assemblies are available at the

225     Figshare repository. The three metagenomes and six metatranscriptomes used in this

226     study are available on the JGI/IMG at accession codes 3300026302, 3300026286,

227     3300009517, and 3300002341-46, respectively. All code for performing metagenomic

228     assembly,         binning,         and         annotation         can         be         found         at

229     https://github.com/elizabethmcd/EBPR-MAGs.         The         TbasCO         method         has         been

230     implemented     as     a     reproducible     R     package     and     can     be     accessed     at

231     https://github.com/Jorisvansteenbrugge/TbasCO.

232

233

234

235

**RESULTS AND DISCUSSION**

**Reconstructing a Diverse EBPR SBR Community**

To explore trait-based transcriptional dynamics of a semi-complex microbial community, we applied genome-resolved metagenomics and metatranscriptomics to an EBPR sequencing-batch reactor (SBR) ecosystem (Figure 2). We previously performed a metatranscriptomics time-series experiment over the course of a normally operating EBPR cycle to investigate the regulatory controls of Accumulibacter gene expression [40]. In this experiment, six samples were collected for RNA sequencing: three from the anaerobic phase and three from the aerobic phase (Figure 2A). Additionally, three metagenomes were collected from the same month of the metatranscriptomic experiment, including a sample from the same date of the experiment. We reassembled contemporary Accumulibacter clade IIA and IA genomes that were previously assembled from the same bioreactor system [27, 28]. The genomes of Accumulibacter clades IA and IIA are similar by approximately 85% average-nucleotide identity [28], and although this is well below the common species-resolved cutoff of 95% [48], we refer to the clade nomenclature defined based on polyphosphate kinase (*ppk1*) sequence identity [49, 50]. During the experiment, the bioreactor was highly enriched in Accumulibacter clade IIA, accounting for approximately 50% of the mapped metagenomic reads and the highest transcriptional counts (Figures 2B and 2C) [40]. Whereas Accumulibacter clade IA exhibited low abundance patterns but was within the top 10 genomes with the highest total transcriptional counts (Figure 2C).

Although this bioreactor system was highly enriched in Accumulibacter, a diverse flanking community persisted and was active in this ecosystem (Figure 2B, C). We

11

259   reconstructed representative population genomes of the microbial community of the SBR

260   system, resulting in 64 metagenome-assembled genomes (MAGs) of the flanking

261   community. Interestingly, we recovered genomes of experimentally verified and putative

262   PAOs, including two *Tetrasphaera spp.* (TET1 and TET2) '*Candidatus Obscuribacter*

263   *phosphatis'* (OBS1)*,* and *Gemmatimonadetes* (GEMMA1). Pure cultures of *Tetrasphaera*

264   have been experimentally shown to cycle polyphosphate without incorporating PHA [36],

265   deviating from the hallmark Accumulibacter PAO model. The first cultured representative

266   of the *Gemmatimonadetes* phylum *Gemmatimonas aurantiaca* was isolated from an SBR

267   simulating EBPR and was shown to accumulate polyphosphate through Neisser and

268   DAPI staining [51]. Additionally, *Ca. Obscuribacter phosphatis* has been hypothesized to

269   cycle phosphorus based on the presence of genes for phosphorus transport,

270   polyphosphate incorporation, and potential for both anaerobic and aerobic respiration

271   [37], and has also been enriched in photobioreactor EBPR systems [52]. Both

272   *Tetrasphaera spp.* TET1 and TET2, OBS1, and GEMMA1 groups exhibit higher relative

273   abundance patterns than CAPIA but have similar relative transcriptional levels (Figure 2B

274   and 2C, Table 1).

275       Numerous SBR MAGs among the *Actinobacteria* and *Proteobacteria* contain the

276   metabolic potential for phosphorus cycling based on the presence of the high-affinity

277   phosphorus transporter *pstABCS* system, polyphosphate kinase *ppk1*, and the low-

278   affinity *pit* phosphorus transporter (Supplementary Figure 5). Additionally, select MAGs

279   within the *Alphaproteobacteria*, *Betaproteobacteria*, and *Gammaproteobacteria* contain

280   all required subunits for polyhydroxyalkanoate synthesis (Supplementary Figure 5). Other

281   abundant and transcriptionally active groups in the SBR ecosystem that are not predicted

282    to be PAOs are members of the *Bacteroidetes* such as CHIT1 within the

283    *Chitinophagaceae,* and *Cytophagales* members *Runella* sp*.* RUN1 and *Leadbetterella* sp.

284    LEAD1 (Figure 2B and 2C, Table 1). Interestingly, an uncharacterized group within the

285    *Bacteroidetes* BAC1 contributed the third most to the pool of transcripts (Figure 2C), and

286    did not show phylogenetic similarity to MAGs assembled from Danish full-scale

287    wastewater treatment systems [39] (Supplementary Figure 1). Other groups from which

288    we assembled MAGs for that do not exhibit clear roles in EBPR systems were *Chloroflexi*

289    ANAER1 and HERP1 MAGs, *Armatimonadetes* FIMBRI1, *Firmicutes* FUSI1, and

290    *Patescibacteria* SACCH1. Members of the *Chloroflexi* are filamentous bacteria that have

291    been associated with bulking and foaming events in full-scale WWTPS [53–55], but also

292    aid in forming the scaffolding around floc aggregates and degrade complex polymers [55–

293    57]. The *Patescibacteria* (formerly TM7) are widespread but low abundant members of

294    natural and engineered ecosystems, contain reduced genome sizes, and may contribute

295    to filamentous bulking in activated sludge [21, 58]. To summarize, lab-scale SBRs

296    designed to enrich for Accumulibacter contain diverse flanking community members [27,

297    32], but their ecological functions and putative interactions remain to be fully understood

298    in the context of the EBPR ecosystem.

299    **Identifying Expression-Based Trait Attributes Among the EBPR SBR Community**

300    **with TbasCO**

301        Current metatranscriptomics approaches often employ either a gene-centric [31,

302    59–61] or genome-centric approaches [40, 62–64]. In both approaches, highly,

303    differentially, or co-expressed genes are identified and tested for enrichment of specific

304    functions. Enrichment- or annotation-based approaches are employed in numerous

13

305    metatranscriptomics tools such as MG-RAST, MetaTrans, SAMSA2, COMAN, IMP, and

306    Anvi'o [65–70]. Here, we expand on the use of molecular markers as traits by defining

307    expression attributes by leveraging *a priori* knowledge from predefined trait libraries, such

308    as the KEGG database [71], to statistically assess inter-species expression patterns of

309    genes that together form a trait (Figure 1). First, our results showed that there is

310    statistically significant transcriptional conservation of genes at the community level; genes

311    that share an annotation were significantly more similar than expected using two different

312    distance metrics (NRED: p-value < 2.2e-16, PC: p-value < 2.2e-16). Extending this

313    statistical analysis to the trait level, we identified 1674 attributes distributed across the 66

314    genomes. On average, we identified 9.12 genomes per attribute (SD - 5.22), with a

315    minimum of 3 genomes and a maximum of 35 (Figure 3A). Based on these statistics, we

316    defined redundant attributes as those two standard deviations above the mean (19

317    genomes). With this cutoff applied, we identified 79 redundant trait attributes mostly

318    belonging to pathways among carbohydrate metabolism, purine metabolism, and fatty

319    acid metabolism categories (Table 2). Of 290 traits, we identified 97 traits with two or

320    more attributes identified (33%). Of these, traits in 10 or more genomes were twice as

321    likely to have two or more attributes (65%), suggesting that divergent expression patterns

322    for a trait are common, and may represent a niche-differentiating feature (Figure 3A).

323    Henceforth, when multiple attributes are identified for a trait, we refer to these as niche-

324    differentiating attributes.

325        From the ecosystem perspective, a clear phylogenetic signal is observed in the

326    distribution of attributes, as genomes cluster together by shared trait attributes by phylum

327    with some exceptions, such as genomes belonging to the *Bacteroidetes, Actinobacteria,*

328    and *Proteobacteria* clustering together, respectively (Figure 3B). For simplicity, we filtered

329    the network to only include nodes with more than 5 connections. Highly redundant trait

330    attributes belonged to modules in the lipid metabolism, energy metabolism, and

331    nucleotide metabolism KEGG functional categories. In contrast, more specialized trait

332    attributes on the periphery of the network or amongst group-specific clusters such as

333    within the *Actinobacteria* or subsets of the *Proteobacteria* belonged to amino acid

334    metabolism, biosynthesis of terpenoids and polyketides, metabolism of cofactors and

335    vitamins, and carbohydrate metabolism KEGG modules. Pathways of note that showed

336    a high level of redundancy include the TCA cycle, isoleucine biosynthesis, acyl-CoA

337    synthesis, threonine biosynthesis, and cytochrome c oxidase activity (Table 2). Large

338    pathways with hundreds of possible routes such as glycolysis, the TCA cycle,

339    gluconeogenesis, and the pentose phosphate pathway are not included in the main

340    network and are displayed as individual networks (Supplementary Figure 6).

341        We next explored the distribution of non-redundant attributes (e.g. 3-18 genomes)

342    (Figure 3A). A total of 796 trait attributes with low redundancy were identified belonging

343    to pathways involved in carbohydrate cofactor and vitamin metabolism including

344    glycolysis, gluconeogenesis, parts of the TCA cycle, tetrahydrofolate biosynthesis,

345    tryptophan biosynthesis, and the pentose phosphate pathway (Table 3). Different sets of

346    low redundancy trait attributes were identified within respective phyla (Supplementary

347    Figure 7). Between genomes belonging to the *Actinobacteria*, *Alphaproteobacteria,*

348    *Bacteroidetes, Betaproteobacteria,* and *Gammaproteobacteria,* low redundancy

349    attributes (belonging to less than half of the total genomes within the phylum) include

350    carbohydrate metabolism, amino acid metabolism and metabolism of cofactors and

15

351     vitamins (Supplementary Figure 7). Redundant trait attributes within individual phyla

352     belong to core energy metabolism pathways, fatty acid biosynthesis, and carbohydrate

353     metabolism. However, even within individual phyla,  non-redundant attributes include

354     different amino acids and cofactors (Extended Table 1 - available on Figshare

355     https://figshare.com/articles/dataset/Lineage-

356     Specific_Core_and_Niche_Differentiating_Traits/15001200).

357         As noted previously, one of the most striking findings is that a majority, 65%  of

358     traits present in 10 or more genomes have multiple expression attributes. Thus, it seems

359     that while the presence of marker genes suggests many organisms share a particular

360     trait, the presence of niche-differentiating expression profiles suggest an alternative story,

361     that there is a level of hidden metabolic diversity. For example, central carbon metabolism

362     and energy pathways such as the TCA cycle, glycolysis, gluconeogenesis, and the

363     pentose phosphate pathway are oftentimes considered core traits when only analyzing

364     the presence and/or absence of individual markers belonging to these pathways. Among

365     over 1000 high-quality MAGs assembled from full-scale Danish WWTPs, the TCA cycle

366     and pentose phosphate pathway are highly represented among the abundant

367     microorganisms, with glycolysis less so [39]. Whereas the TCA cycle and pentose

368     phosphate pathway are present among a high number of genomes in the EBPR SBR

369     community, different routes or parts of these pathways have niche-differentiating

370     distributions (Supplementary Figure 4, Tables 2 and 3). These finer-scale differences in

371     expression of "core" traits may explain the persistence of a diverse community when

372     solely fed acetate, as different lineages could employ similar carbon utilization pathways

373     differently or in more versatile ways. Another salient aspect of this analysis is the

374 astonishingly high number of possible routes within individual pathways here represented

375 by their Disjunctive Normal Forms. For example, accounting for all alternative routes and

376 enzymes, the glycolysis pathway has 100s of possible routes. Layering upon this many

377 expression attributes reveals a large hidden metabolic versatility.

378 **Dimensionality of the High-Affinity Phosphorus Transporter System *PstABCS***

379 The EBPR ecosystem is characterized by its highly dynamic phosphorus cycles.

380 To explore how different lineages respond to fluctuating phosphorus concentrations, we

381 explored the expression-based attributes for the KEGG module of the high-affinity

382 phosphorus transporter *pstABCS* (Figure 4). The *pstABCS* system is an ABC-type

383 transporter that strongly binds phosphate under phosphorus-limiting conditions;

384 therefore, it would be expected that the highest expression levels would be at the end of

385 the aerobic cycle [72]. In contrast, we found that expression of the *pstABCS* was

386 characterized by two different trait attributes. In the first attribute shared by 14 community

387 members, all components of *pstABCS* displayed the highest activity towards the end of

388 the aerobic cycle, when phosphorus concentrations were depleted (Figure 4, Attribute 1).

389 Conversely, 11 community members displayed an alternate attribute where the highest

390 activity of *pstABCS* was at the transition from anaerobic to aerobic phases when

391 phosphorus concentrations are highest (Figure 4, Attribute 2).

392 These results are in agreement with previous results showing that Accumulibacter

393 clade IIC has a canonical *pstABCS* expression pattern (as in Figure 4, Attribute 1) ,

394 whereas the Accumulibacter clade IA has a non-canonical expression (as in Figure 4,

395 Attribute 2) [31]. By assigning trait attributes, we are able to extend these findings beyond

396 Accumulibacter to other flanking community members in the SBR ecosystem suggesting

397    that there are conserved ecological pressures driving niche differentiating expression

398    patterns in *pstABCS* within the EBPR community.

399    **Distribution and Expression of Truncated Denitrification Steps Among EPBR**

400    **Community Members**

401    Understanding the induction of denitrification is an important ecosystem property

402    linked to the redox status of an environment. In EBPR communities, there are many

403    diverse and incomplete denitrification pathways, distributed across many lineages

404    denitrification steps expected in denitrifying systems (Figure 5). Among all 66 MAGs, we

405    did not identify any single MAG with a complete denitrification pathway consisting of the

406    genetic repertoire necessary to fully reduce nitrate to nitrogen gas (Supplementary Figure

407    5). Instead, we identified multiple groups of organisms with truncated denitrification

408    pathways, with steps distributed among cohorts of community members (Figure 5).

409    For the first steps of reducing nitrate to nitrite, we explored expression attributes

410    of the *napAB* and *narGH* pathways (Figure 5B, C). For the *narGH* pathway, two attributes

411    were identified (Figure 5B). The first *narGH* attribute was characterized by high

412    expression in the anaerobic phase, with decreasing activity by the second time point of

413    the anaerobic phase. Genomes containing this attribute included the experimentally

414    verified and putative PAOs *Tetrasphaera* (TET1 and TET2) and *Ca.* Obscuribacter

415    (OBS1), respectively. The second attribute was exhibited among members of the

416    *Actinobacteria* (PROP2, PHYC2, PROP3, and NANO1), *Proteobacteria* (BEIJ4), and

417    *Bacteroidetes* (BAC1). The attribute identified for *napAB* was also more highly expressed

418    anaerobically and included CAPIA, CAPIIA, ALIC1, REYR2, RUBRI1, and BEIJ3.

419    Interestingly, this *napAB* attribute had expression patterns that quickly decreased in the

420 first aerobic time point, suggesting a tighter regulation than Attribute 1 for *narGH*.

421 Together, this suggests that the regulation of denitrification within the EBPR ecosystem

422 is a niche-differentiating feature whereby the induction of denitrification pathways occurs

423 either anaerobically or only after anaerobic carbon contact.

424      A smaller cohort contained the genetic repertoire to reduce nitrite to nitrogen gas

425 and exhibited hallmark anaerobic-aerobic expression patterns (Figure 5E) These

426 members within the *Proteobacteria* (OTTO2, BEIJ3, VITREO1, and ZOO1) contained the

427 *nirS* nitrite reductase, the *norBC* nitric oxide reductase, and *nosZ,* and showed highest

428 expression of these subunits towards the beginning of the anaerobic cycle, slowly

429 decreasing over the aerobic period to their lowest in the end of the aerobic cycle. Although

430 BEIJ2 was lacking the *norBC* system, it contained the *nirS* nitrite reductase and *nosZ*

431 subunit, and exhibited similar expression patterns to others in this cohort. Other

432 *Proteobacteria* lineages only contained the *norBC* subunits but were expressed in similar

433 fashions (RHODO2, FLAVO1, RHIZO1, and LEAD1) (Figure 5D). Accumulibacter clades

434 IA and IIA as well as ALIC1 were the only lineages with near-complete denitrification

435 pathways. These lineages contained the *napAB* nitrate reductase system as mentioned

436 above, the *nirS* nitrite reductase, *norB* (missing a confident hit for the *norC* subunit), and

437 *nosZ.* These three lineages also exhibited hallmark upregulation of all steps in the

438 anaerobic phase, with decreased activity after aerobic contact (Figure 5F).

439      Interestingly, Accumulibacter clade IA exhibited a higher magnitude of expression

440 of denitrification steps when activity levels were normalized relative to clade IIA,

441 supporting the hypothesis that denitrification is a niche-differentiating feature among

442 clades [28, 31, 73], and possibly a strain-specific trait since denitrification traits cannot be

443     predicted based on *ppk1* clade designations [32]. For example, independent observations

444     in differences among denitrification activities among strains within Accumulibacter clade

445     IC are inconsistent [34, 74]. Within the same bioreactor environment, coexisting

446     Accumulibacter clades differ between denitrification abilities and expression profiles [31,

447     33, 75]. Truncated denitrification pathways have also been previously shown to be

448     distributed among community members, with the complete denitrification genetic

449     repertoire only present in few members [33, 75], which could be due to extensive

450     horizontal gene transfer of genes comprising denitrification steps [75, 76]. Although this

451     experiment was not conducted under denitrifying conditions, our approach could be

452     applied to denitrifying EBPR systems to further understand the distribution of

453     denitrification traits among community members and how to selectively enrich for diverse

454     DPAOs.

455     **Biosynthetic Potential and Expression Dynamics of Amino Acid and Vitamin**

456     **Synthesis Pathways**

457     Although SBRs are designed to enrich for Accumulibacter by providing acetate as

458     the sole carbon source, a diverse flanking community persists in these setups [27, 75].

459     One hypothesis for the persistence of flanking community members may be cooperative

460     interactions due to underlying auxotrophies of amino acid and vitamin biosynthetic

461     pathways in Accumulibacter. Amino acids and vitamin cofactors are metabolically

462     expensive to synthesize, and widespread auxotrophies have been widely documented

463     among microbial communities [77, 78]. Specifically, auxotrophies of vitamin cofactors

464     have been shown to fuel bacterial and cross-kingdom interactions with *de novo* bacterial

465     and cross-kingdom interactions with *de novo* synthesizers [79, 80]. To explore this

466     hypothesis in the EPBR SBR community, we analyzed the presence of amino acid and

467     vitamin biosynthetic pathways and their expression patterns among the top 15 genomes

468     based on transcript abundance (Figure 6).

469         Within Accumulibacter, there are a few key vitamin cofactor and amino acid

470     auxotrophies that could fuel potential interactions with flanking community members. Both

471     Accumulibacter clade genomes are missing the riboflavin pathway for FAD cofactor

472     synthesis, as well as the pathways for serine and aspartic acid (Figure 6A). The

473     biosynthetic pathway for aspartic acid is distributed among members of the *Bacteroidetes*

474     and *Proteobacteria*, whereas only TET2 contains the pathway for serine synthesis (Figure

475     5A). The lack of serine biosynthesis pathways in Accumulibacter and other flanking

476     genomes seems striking given that serine is one of the least metabolically costly amino

477     acids to synthesize [81]. Interestingly, Accumulibacter clade IIA does not contain the

478     biosynthetic machinery for thiamine and pantothenate synthesis, whereas clade IA does

479     (Figure 6A). Only the CAULO1, HYPHO1, and PSEUDO1 genomes within the

480     *Proteobacteria* can synthesize thiamine, whereas several other members can synthesize

481     pantothenate (Figure 6A). The absence of the pantothenate biosynthetic pathway in

482     Accumulibacter CAP IIA is particularly interesting given that coenzyme A is essential for

483     polyhydroxyalkanoate biosynthesis, which fuels the polymer cycling PAO phenotype of

484     Accumulibacter [24].

485         In addition to flanking community members potentially supporting the growth of

486     Accumulibacter due to underlying auxotrophies, the reciprocal logic may be possible as

487     well. Both Accumulibacter clades contain the pathways for synthesizing tyrosine and

488     phenylalanine, which are missing in a majority of the top 15 active flanking genomes

489   (Figure 6A). Only two other members within the *Proteobacteria* can synthesize tyrosine

490   and phenylalanine, where RAM1 can synthesize both and PSEUDO1 only phenylalanine.

491   Interestingly, phenylalanine and tyrosine are the second and third most metabolically

492   expensive amino acids to synthesize, respectively, with tryptophan the most costly [81].

493   Additionally, a few highly active flanking community members lack the biosynthetic

494   machinery for several vitamin cofactors and amino acids, such as FLAVO1 and BAC3

495   within the *Bacteroidetes* and the putative PAO *Ca.* Obscuribacter phosphatis OBS1

496   (Figure 6A). Particularly, RAM1 within the *Proteobacteria* is missing the biosynthetic

497   machinery for all vitamin cofactors but can synthesize most amino acids including the

498   most metabolically expensive as mentioned above.

499         We next analyzed the distribution of trait-attributes of vitamin and amino acid

500   pathways among these genomes to understand how these biosynthetic pathways are

501   expressed similarly or differently in the EBPR SBR ecosystem (Figure 6B and C).

502   Members of the *Proteobacteria* containing thiamine and cobalamin biosynthetic pathways

503   all express these traits similarly (Figure 6B). However, the pantothenate synthesis

504   pathway contains two trait-attributes and is expressed differently among two cohorts. In

505   the first attribute, RUN1, TET1, CAULO1, CAPIA, and PSEUDO1 express the

506   pantothenate pathway similarly. However, OBS1 and TET2 express the pantothenate

507   pathway differently (Figure 6B). Because tetrahydrofolate can be synthesized through

508   different metabolic routes, we analyzed the differences in trait attribute expression for all

509   routes in genomes that contained sufficient coverage of this trait. Members of the

510   *Bacteroidetes* and *Proteobacteria* mostly cluster together among tetrahydrofolate

511   attributes, whereas the TET1 and TET2 genomes are differentiated (Figure 6B).

512       Expression of various groups of amino acids show more differentiated patterns of

513       expression for genomes with these pathways. Several amino acids also contain different

514       metabolic routes for biosynthesis, and we analyzed all trait attributes for each amino acid

515       for all routes grouped by type (Figure 6C). For the charged amino acids arginine, histidine,

516       and lysine, members of the *Proteobacteria* and *Bacteroidetes* cluster within their

517       phylogenetic groups, respectively, with lysine and histidine expressed differently among

518       these groups (Figure 6C). In contrast, arginine is expressed similarly among all

519       *Proteobacteria* genomes. Among the polar charged amino acids, TET2 is the only

520       genome among the top 15 genomes that contains the metabolic pathway to synthesize

521       serine (Figure 6A). Several groups contain the pathway for threonine synthesis, and

522       expression of different threonine routes are differentiated among the *Proteobacteria,*

523       *Bacteroidetes,* and *Tetrasphaera spp.,* but mostly clusters phylogenetically (Figure 6C)*.*

524       Notably, the expression patterns for the cysteine and proline biosynthetic pathways do

525       not cluster phylogenetically, such as both *Tetrasphaera* genomes expressing the proline

526       pathway more similarly to other *Proteobacteria* and *Bacteroidetes* (Figure 6C)*.* The few

527       lineages that can synthesize tyrosine and phenylalanine (CAPIA, CAPIIA, RAM1,

528       PSEUDO1) show different patterns of expression. These results show that beyond the

529       presence or absence of key vitamin cofactor and amino acid biosynthetic pathways,

530       EBPR SBR organisms also display coherent and differentiated patterns of expression for

531       these traits, of which the functional consequences remain to be further understood.

532

533

534

**CONCLUSIONS AND FUTURE PERSPECTIVES**

In this work, we applied a novel trait-based 'omics pipeline to a semi-complex, engineered bioreactor microbial community to explore ecosystem-level and niche-differentiating traits. Through assembling high-quality MAGs of the EBPR SBR community and using a time-series metatranscriptomics experiment, we were able to extend functional predictions and ecosystem inferences beyond hypotheses made from gene presence/absence data. Using our novel trait-based comparative 'omics pipeline, we identified how similarities and differences in the expression of significant EBPR traits are conferred among community members such as phosphorus cycling, denitrification, and amino acid metabolism. Specifically, we demonstrate that traits with similar expression profiles may be clustered into attributes providing a new layer to trait-based approaches.

We believe that identifying expression-based attributes will be a powerful tool to explore microbial traits in natural, engineered, and host-associated microbiomes. Outside of activated sludge systems, trait-based approaches could illuminate how similar secondary metabolite clusters are expressed among different species in a community [82, 83], how auxotrophies for amino acid and vitamin cofactors govern interactions [84], how rhizosphere microorganisms respond to day-night cycles, and identify putative traits that universally exhibit ecosystem-level or niche-differentiating patterns across ecosystems [19, 23]. Importantly, our trait-based approach can be used to screen for expected expression patterns of a key trait compared to a model organism, and then prioritize specific microbial lineages for downstream experimental verification with techniques such as Raman-FISH [85, 86]. Overall, our trait-based comparative 'omics pipeline is a novel

558    and high-throughput approach to understand how microbial traits connect to ecosystem-

559    level processes in diverse microbiomes.

560

561    **ACKNOWLEDGEMENTS**

573

574

575

576

577

578

579

580

## REFERENCES CITED

581  1.  Violle C, Navas M-L, Vile D, Kazakou E, Fortunel C, Hummel I, et al. Let the

583      concept of trait be functional! *Oikos* 2007; **116**: 882–892.

584  2.  Lavorel S, Garnier E. Predicting changes in community composition and

585      ecosystem functioning from plant traits: revisiting the Holy Grail. *Funct Ecol* 2002;

586      **16**: 545–556.

587  3.  Hooper DU, Chapin FS, Ewel JJ, Hector A, Inchausti P, Lavorel S, et al.

588      EFFECTS OF BIODIVERSITY ON ECOSYSTEM FUNCTIONING: A

589      CONSENSUS OF CURRENT KNOWLEDGE. *Ecol Monogr* 2005; **75**: 3–35.

590  4.  Pianka ER. On r-and K-Selection. *Am Nat* 1970.

591  5.  Wright IJ, Reich PB, Westoby M, Ackerly DD, Baruch Z, Bongers F, et al. The

592      worldwide leaf economics spectrum. *Nature* 2004; **428**: 821–827.

593  6.  Krause S, Le Roux X, Niklaus PA, Van Bodegom PM, Lennon JT, Bertilsson S, et

594      al. Trait-based approaches for understanding microbial biodiversity and

595      ecosystem functioning. *Front Microbiol* 2014; **5**: 251.

596  7.  Malik AA, Martiny JBH, Brodie EL, Martiny AC, Treseder KK, Allison SD, et al.

597      Defining trait-based microbial strategies with consequences for soil carbon cycling

598      under climate change. *bioRxiv* 2019.

599  8.  Guittar J, Shade A, Litchman E. Trait-based community assembly and succession

600      of the infant gut microbiome. *Nat Commun* 2019; **10**: 512.

601  9.  Wolfe BE, Button JE, Santarelli M, Dutton RJ. Cheese Rind Communities Provide

602      Tractable Systems for In Situ and In Vitro Studies of Microbial Diversity. *Cell*

603      2014; **158**: 422–433.

26

604   10.   Enke TN, Datta MS, Schwartzman J, Barrere J, Pascual-García A, Cordero

605         Correspondence OX. Modular Assembly of Polysaccharide-Degrading Marine

606         Microbial Communities. *Curr Biol* 2019; **29**.

607   11.   Herrera Paredes S, Gao T, Law TF, Finkel OM, Mucyn T, Teixeira PJPL, et al.

608         Design of synthetic bacterial communities for predictable plant phenotypes. *PLOS*

609         *Biol* 2018; **16**: e2003962.

610   12.   Lindemann SR, Bernstein HC, Song H-S, Fredrickson JK, Fields MW, Shou W, et

611         al. Engineering microbial consortia for controllable outputs. *ISME J* 2016; **10**:

612         2077–2084.

613   13.   Oyserman BO, Medema MH, Raaijmakers JM. Road MAPs to engineer host

614         microbiomes. *Curr Opin Microbiol* 2018; **43**: 46–54.

615   14.   Lawson CE, Harcombe WR, Hatzenpichler R, Lindemann SR, Löffler FE,

616         O'Malley MA, et al. Common principles and best practices for engineering

617         microbiomes. *Nat Rev Microbiol* . 2019. Nature Publishing Group. , **17**: 725–741

618   15.   Gutierrez CF, Sanabria J, Raaijmakers JM, Oyserman BO. Restoring degraded

619         microbiome function with self-assembled communities. *FEMS Microbiol Ecol*

620         2020; **96**.

621   16.   Allison SD. A trait-based approach for modelling microbial litter decomposition.

622         *Ecol Lett* 2012; **15**: 1058–1070.

623   17.   Tringe SG, von Mering C, Kobayashi A, Salamov AA, Chen K, Chang HW, et al.

624         Comparative metagenomics of microbial communities. *Science* 2005; **308**: 554–7.

625   18.   Tyson GW, Chapman J, Hugenholtz P, Allen EE, Ram RJ, Richardson PM, et al.

626         Community structure and metabolism through reconstruction of microbial

627      genomes from the environment. *Nature* 2004; **428**: 37–43.

628    19.   Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, et al.

629      Thousands of microbial genomes shed light on interconnected biogeochemical

630      processes in an aquifer system. *Nat Commun* 2016; **7**: 13219.

631    20.   Woodcroft BJ, Singleton CM, Boyd JA, Evans PN, Emerson JB, Zayed AAF, et al.

632      Genome-centric view of carbon processing in thawing permafrost. *Nature* 2018;

633      **560**: 49–54.

634    21.   Albertsen M, Hugenholtz P, Skarshewski A, Nielsen KL, Tyson GW, Nielsen PH.

635      Genome sequences of rare, uncultured bacteria obtained by differential coverage

636      binning of multiple metagenomes. *Nat Biotechnol* 2013; **31**: 533–538.

637    22.   Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, et al.

638      Thousands of microbial genomes shed light on interconnected biogeochemical

639      processes in an aquifer system. *Nat Commun* 2016; **7**: 13219.

640    23.   Wrighton KC, Thomas BC, Sharon I, Miller CS, Castelle CJ, VerBerkmoes NC, et

641      al. Fermentation, hydrogen, and sulfur metabolism in multiple uncultivated

642      bacterial phyla. *Science* 2012; **337**: 1661–5.

643    24.   Hesselmann RPX, Werlen C, Hahn D, van der Meer JR, Zehnder AJB.

644      Enrichment, Phylogenetic Analysis and Detection of a Bacterium That Performs

645      Enhanced Biological Phosphate Removal in Activated Sludge. *Syst Appl Microbiol*

646      1999; **22**: 454–465.

647    25.   Seviour RJ, Mino T, Onuki M. The microbiology of biological phosphorus removal

648      in activated sludge systems. *FEMS Microbiol Rev* 2003; **27**: 99–127.

649    26.   Crocetti GR, Hugenholtz P, Bond PL, Schuler A, Ju¨ J, Keller J, et al. Identification

650      of Polyphosphate-Accumulating Organisms and Design of 16S rRNA-Directed

651      Probes for Their Detection and Quantitation. *APPLIED AND ENVIRONMENTAL*

652      *MICROBIOLOGY* . 2000.

653    27.    Martín HG, Ivanova N, Kunin V, Warnecke F, Barry KW, McHardy AC, et al.

654      Metagenomic analysis of two enhanced biological phosphorus removal (EBPR)

655      sludge communities. *Nat Biotechnol* 2006; **24**: 1263–1269.

656    28.    Flowers JJ, He S, Malfatti S, del Rio TG, Tringe SG, Hugenholtz P, et al.

657      Comparative genomics of two 'Candidatus Accumulibacter' clades performing

658      biological phosphorus removal. *ISME J* 2013; **7**: 2301–2314.

659    29.    Oyserman BO, Moya F, Lawson CE, Garcia AL, Vogt M, Heffernen M, et al.

660      Ancestral genome reconstruction identifies the evolutionary basis for trait

661      acquisition in polyphosphate accumulating bacteria. *ISME J* 2016; **10**: 2931–

662      2945.

663    30.    Wilmes P, Andersson AF, Lefsrud MG, Wexler M, Shah M, Zhang B, et al.

664      Community proteogenomics highlights microbial strain-variant protein expression

665      within activated sludge performing enhanced biological phosphorus removal.

666      *ISME J* 2008; **2**: 853–864.

667    31.    McDaniel E, Moya-Flores F, Keene Beach N, Oyserman B, Kizaric M, Hoe Khor

668      E, et al. Metabolic differentiation of co-occurring Accumulibacter clades revealed

669      through genome-resolved metatranscriptomics. *bioRxiv* 2020;

670      2020.11.23.394700.

671    32.    Gao H, Mao Y, Zhao X, Liu WT, Zhang T, Wells G. Genome-centric

672      metagenomics resolves microbial diversity and prevalent truncated denitrification

673    pathways in a denitrifying PAO-enriched bioprocess. *Water Res* 2019; **155**: 275–

674    287.

675  33.  Wang Y, Gao H, Wells G. Integrated Omics Analyses Reveal Differential Gene

676    Expression and Potential for Cooperation Between Denitrifying Polyphosphate

677    and Glycogen Accumulating Organisms. *bioRxiv* . 2020. bioRxiv. ,

678    2020.01.10.901413

679  34.  Camejo PY, Oyserman BO, McMahon KD, Noguera DR. Integrated Omic

680    Analyses Provide Evidence that a 'Candidatus Accumulibacter phosphatis' Strain

681    Performs Denitrification under Microaerobic Conditions. *mSystems* 2019; **4**:

682    e00193-18.

683  35.  Kong Y, Nielsen JL, Nielsen PH. Identity and ecophysiology of uncultured

684    actinobacterial polyphosphate-accumulating organisms in full-scale enhanced

685    biological phosphorus removal plants. *Appl Environ Microbiol* 2005; **71**: 4076–85.

686  36.  Kristiansen R, Nguyen HTT, Saunders AM, Nielsen JL, Wimmer R, Le VQ, et al.

687    A metabolic model for members of the genus Tetrasphaera involved in enhanced

688    biological phosphorus removal. *ISME J* 2013; **7**: 543–54.

689  37.  Soo R, Skennerton CT, Sekiguchi Y, Imelfort M, Paech S, Dennis P, et al. An

690    Expanded Genomic Representation of the Phylum Cyanobacteria. *Genome*

691    *Biology and Evolution*. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4040986/.

692    Accessed 11 Jul 2020.

693  38.  Petriglieri F, Singleton C, Peces M, Petersen JF, Nierychlo M, Nielsen H.

694    'Candidatus Dechloromonas phosphatis' and 'Candidatus Dechloromonas

695    phosphovora', two novel polyphosphate accumulating organisms abundant in

696         wastewater treatment systems. *bioRxiv* 2020.

697   39.   Singleton CM, Petriglieri F, Kristensen JM, Kirkegaard RH, Michaelsen TY,

698         Andersen MH, et al. Connecting structure to function with the recovery of over

699         1000 high-quality metagenome-assembled genomes from activated sludge using

700         long-read sequencing. *Nat Commun* 2021; **12**: 2009.

701   40.   Oyserman BO, Noguera DR, del Rio TG, Tringe SG, McMahon KD.

702         Metatranscriptomic insights on gene expression and regulatory controls in

703         Candidatus Accumulibacter phosphatis. *ISME J* 2016; **10**: 810–822.

704   41.   Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM:

705         assessing the quality of microbial genomes recovered from isolates, single cells,

706         and metagenomes. *Genome Res* 2015; **25**: 1043–55.

707   42.   Kanehisa M, Sato Y, Morishima K. BlastKOALA and GhostKOALA: KEGG Tools

708         for Functional Characterization of Genome and Metagenome Sequences. *J Mol*

709         *Biol* 2016; **428**: 726–731.

710   43.   Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, et al.

711         KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive

712         score threshold. *bioRxiv* 2019; 602110.

713   44.   Bushnell B, Rood J, Singer E. BBMerge – Accurate paired shotgun read merging

714         via overlap. *PLoS One* 2017; **12**: e0185056.

715   45.   Kopylova E, Noé L, Touzet H. SortMeRNA: fast and accurate filtering of ribosomal

716         RNAs in metatranscriptomic data. *Bioinformatics* 2012; **28**: 3211–3217.

717   46.   Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq

718         quantification. *Nat Biotechnol* 2016; **34**: 525–527.

719    47.    Soneson C, Love MI, Robinson MD. Differential analyses for RNA-seq: transcript-

720           level estimates improve gene-level inferences. *F1000Research* 2015; **4**: 1521.

721    48.    Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High

722           throughput ANI analysis of 90K prokaryotic genomes reveals clear species

723           boundaries. *Nat Commun* 2018; **9**: 5114.

724    49.    He S, Gall DL, McMahon KD. 'Candidatus accumulibacter' population structure in

725           enhanced biological phosphorus removal sludges as revealed by polyphosphate

726           kinase genes. *Appl Environ Microbiol* 2007; **73**: 5865–5874.

727    50.    Camejo PY, Owen BR, Martirano J, Ma J, Kapoor V, Santo Domingo J, et al.

728           Candidatus Accumulibacter phosphatis clades enriched under cyclic anaerobic

729           and microaerobic conditions simultaneously use different electron acceptors.

730           *Water Res* 2016; **102**: 125–137.

731    51.    Zhang H, Sekiguchi Y, Hanada S, Hugenholtz P, Kim H, Kamagata Y, et al.

732           Gemmatimonas aurantiaca gen. nov., sp. nov., a Gram-negative, aerobic,

733           polyphosphate-accumulating micro-organism, the first cultured representative of

734           the new bacterial phylum Gemmatimonadetes phyl. nov. *Int J Syst Evol Microbiol*

735           2003; **53**: 1155–1163.

736    52.    McDaniel EA, Wever R, Oyserman BO, Noguera DR, McMahon KD. Genome-

737           Resolved Metagenomics of a Photosynthetic Bioreactor Performing Biological

738           Nutrient Removal. *Microbiol Resour Announc* 2021; **10**.

739    53.    Speirs LBM, Rice DTF, Petrovski S, Seviour RJ. The Phylogeny, Biodiversity, and

740           Ecology of the Chloroflexi in Activated Sludge. *Front Microbiol* . 2019. Frontiers

741           Media S.A. , **10**: 2015

742    54.    Andersen MH, McIlroy SJ, Nierychlo M, Nielsen PH, Albertsen M. Genomic

743            insights into Candidatus Amarolinea aalborgensis gen. nov., sp. nov., associated

744            with settleability problems in wastewater treatment plants. *Syst Appl Microbiol*

745            2019; **42**: 77–84.

746    55.    Nierychlo M, Miłobędzka A, Petriglieri F, McIlroy B, Nielsen PH, McIlroy SJ. The

747            morphology and metabolic potential of the Chloroflexi in full-scale activated

748            sludge wastewater treatment plants. *FEMS Microbiol Ecol* 2019; **95**.

749    56.    McIlroy SJ, Karst SM, Nierychlo M, Dueholm MS, Albertsen M, Kirkegaard RH, et

750            al. Genomic and in situ investigations of the novel uncultured Chloroflexi

751            associated with 0092 morphotype filamentous bulking in activated sludge. *ISME J*

752            2016; **10**: 2223–2234.

753    57.    Kragelund C, Levantesi C, Borger A, Thelen K, Eikelboom D, Tandoi V, et al.

754            Identity, abundance and ecophysiology of filamentous Chloroflexi species present

755            in activated sludge treatment plants. *FEMS Microbiol Ecol* 2007; **59**: 671–682.

756    58.    Kindaichi T, Yamaoka S, Uehara R, Ozaki N, Ohashi A, Albertsen M, et al.

757            Phylogenetic diversity and ecophysiology of Candidate phylum Saccharibacteria

758            in activated sludge. *FEMS Microbiol Ecol* 2016; **92**: 1–11.

759    59.    Mann E, Wetzels SU, Wagner M, Zebeli Q, Schmitz-Esser S. Metatranscriptome

760            Sequencing Reveals Insights into the Gene Expression and Functional Potential

761            of Rumen Wall Bacteria. *Front Microbiol* 2018; **9**: 43.

762    60.    Jiang Y, Xiong X, Danska J, Parkinson J. Metatranscriptomic analysis of diverse

763            microbial communities reveals core metabolic pathways and microbiome-specific

764            functionality. *Microbiome* 2016; **4**: 2.

765  61.  Linz AM, Aylward FO, Bertilsson S, McMahon KD. Time-series

766       metatranscriptomes reveal conserved patterns between phototrophic and

767       heterotrophic microbes in diverse freshwater systems. *Limnol Oceanogr* 2019.

768  62.  Lawson CE, Wu S, Bhattacharjee AS, Hamilton JJ, McMahon KD, Goel R, et al.

769       Metabolic network analysis reveals microbial community interactions in anammox

770       granules. *Nat Commun* 2017; **8**: 15416.

771  63.  Aylward FO, Eppley JM, Smith JM, Chavez FP, Scholin CA, DeLong EF.

772       Microbial community transcriptional networks are conserved in three domains at

773       ocean basin scales. *Proc Natl Acad Sci* 2015; **112**: 5443–5448.

774  64.  Hao L, Michaelsen TY, Singleton CM, Dottorini G, Kirkegaard RH, Albertsen M, et

775       al. Novel syntrophic bacteria in full-scale anaerobic digesters revealed by

776       genome-centric metatranscriptomics. *ISME J* 2020; **14**: 906–918.

777  65.  Glass EM, Meyer F. The Metagenomics RAST Server: A Public Resource for the

778       Automatic Phylogenetic and Functional Analysis of Metagenomes. *Handb Mol*

779       *Microb Ecol I Metagenomics Complement Approaches* 2011; **8**: 325–331.

780  66.  Martinez X, Pozuelo M, Pascal V, Campos D, Gut I, Gut M, et al. MetaTrans: an

781       open-source pipeline for metatranscriptomics. *Sci Rep* 2016; **6**: 26447.

782  67.  Westreich ST, Treiber ML, Mills DA, Korf I, Lemay DG. SAMSA2: a standalone

783       metatranscriptome analysis pipeline. *BMC Bioinformatics* 2018; **19**: 175.

784  68.  Ni Y, Li J, Panagiotou G. COMAN: a web server for comprehensive

785       metatranscriptomics analysis. *BMC Genomics* 2016; **17**: 622.

786  69.  Narayanasamy S, Jarosz Y, Muller EEL, Heintz-Buschart A, Herold M, Kaysen A,

787       et al. IMP: a pipeline for reproducible reference-independent integrated

788    metagenomic and metatranscriptomic analyses. *Genome Biol* 2016; **17**: 260.

789  70.  Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, et al. Anvi'o:

790    an advanced analysis and visualization platform for 'omics data. *PeerJ* 2015; **3**:

791    e1319.

792  71.  Kanehisa M, Sato Y, Kawashima M, Furumichi M, Tanabe M. KEGG as a

793    reference resource for gene and protein annotation. *Nucleic Acids Res* 2016; **44**:

794    D457–D462.

795  72.  Wanner BL. Gene regulation by phosphate in enteric bacteria. *J Cell Biochem*

796    1993; **51**: 47–54.

797  73.  Flowers JJ, He S, Yilmaz S, Noguera DR, McMahon KD. Denitrification

798    capabilities of two biological phosphorus removal sludges dominated by different

799    'Candidatus Accumulibacter' clades. *Environ Microbiol Rep* 2009; **1**: 583–588.

800  74.  Rubio-Rincón FJ, Weissbrodt DG, Lopez-Vazquez CM, Welles L, Abbas B,

801    Albertsen M, et al. 'Candidatus Accumulibacter delftensis': A clade IC novel

802    polyphosphate-accumulating organism without denitrifying activity on nitrate.

803    *Water Res* 2019; **161**: 136–151.

804  75.  Gao H, Mao Y, Zhao X, Liu WT, Zhang T, Wells G. Genome-centric

805    metagenomics resolves microbial diversity and prevalent truncated denitrification

806    pathways in a denitrifying PAO-enriched bioprocess. *Water Res* 2019; **155**: 275–

807    287.

808  76.  Parsons C, Stüeken EE, Rosen CJ, Mateos K, Anderson RE. Radiation of

809    nitrogen-metabolizing enzymes across the tree of life tracks environmental

810    transitions in Earth history. *Geobiology* 2020.

811   77.   Gómez-Consarnau L, Sachdeva R, Gifford SM, Cutter LS, Fuhrman JA, Sañudo-

812         Wilhelmy SA, et al. Mosaic patterns of B-vitamin synthesis and utilization in a

813         natural marine microbial community. *Environ Microbiol* 2018; **20**: 2809–2823.

814   78.   Hamilton JJ, Garcia SL, Brown BS, Oyserman BO, Moya-Flores F, Bertilsson S,

815         et al. Metabolic Network Analysis and Metatranscriptomics Reveal Auxotrophies

816         and Nutrient Sources of the Cosmopolitan Freshwater Microbial Lineage acI.

817         *mSystems* 2017; **2**: e00091-17.

818   79.   McClure RS, Overall CC, Hill EA, Song H-S, Charania M, Bernstein HC, et al.

819         Species-specific transcriptomic network inference of interspecies interactions.

820         *ISME J* 2018; 1.

821   80.   Croft MT, Lawrence AD, Raux-Deery E, Warren MJ, Smith AG. Algae acquire

822         vitamin B12 through a symbiotic relationship with bacteria. *Nature* 2005; **438**: 90–

823         93.

824   81.   Akashi H, Gojobori T. Metabolic efficiency and amino acid composition in the

825         proteomes of Escherichia coli and Bacillus subtilis. *Proc Natl Acad Sci U S A*

826         2002; **99**: 3695–3700.

827   82.   Lozano GL, Bravo JI, Diago MFG, Park HB, Hurley A, Peterson SB, et al.

828         Introducing THOR, a Model Microbiome for Genetic Dissection of Community

829         Behavior. *MBio* 2019; **10**: e02846-18.

830   83.   Crits-Christoph A, Diamond S, Butterfield CN, Thomas BC, Banfield JF. Novel soil

831         bacteria possess diverse genes for secondary metabolite biosynthesis. *Nature*

832         2018; **558**: 440–444.

833   84.   Zengler K, Zaramela LS. The social network of microorganisms — how

834    auxotrophies shape complex communities. *Nat Rev Microbiol* 2018; **16**: 383–390.

835    85.    Fernando EY, McIlroy SJ, Nierychlo M, Herbst FA, Petriglieri F, Schmid MC, et al.

836          Resolving the individual contribution of key microbial populations to enhanced

837          biological phosphorus removal with Raman–FISH. *ISME J* 2019; **13**: 1933–1946.

838    86.    Petriglieri F, Petersen JF, Peces M, Nierychlo M, Hansen K, Baastrand CE, et al.

839          Quantification of biologically and chemically bound phosphorus in activated

840          sludge from full-scale plants with biological P-removal. *bioRxiv* 2021.

841    87.    Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH. GTDB-Tk: a toolkit to classify

842          genomes with the Genome Taxonomy Database. *Bioinformatics* 2019.

843    88.    Seemann T. Prokka: Rapid prokaryotic genome annotation. *Bioinformatics* 2014;

844          **30**.

845

846

847

848

849

850

851

852

853

854

855

856

857 **FIGURE AND TABLE LEGENDS**

858 **Figure 1. Overview of Trait-based Comparative Transcriptomics Approach**

859 In genome-resolved metagenomics approaches, representative MAGs are assembled

860 from a microbial community of interest, and the presence and/or absence of key metabolic

861 pathways are used to make inferences of metabolic potential and ecosystem processes.

862 However, metagenomic data alone can only assess the metabolic potential of a given

863 pathway, and do not provide other biologically relevant information such as the timing or

864 induction of these traits. Using time-series metatranscriptomics, we developed a trait-

865 based comparative 'omics (TbasCO) pipeline that statistically assesses the inter-

866 organismal differences in gene expression pattern of a given trait to cluster into trait

867 attributes.

868

869 **Figure 2. Genome-Resolved Metatranscriptomics Approach of an EBPR System**

870 Application of a genome-resolved metatranscriptomics approach to a lab-scale

871 sequencing batch reactor (SBR) designed to enrich for Accumulibacter. **1A)** Schematic

872 of the main cycle parameters and analyte dynamics of an SBR simulating EBPR. Six

873 samples were taken for RNA sequencing within the cycle at time-points denoted by

874 arrows. **1B)** Phylogenetic identity and abundance patterns of 66 assembled MAGs from

875 the EBPR system. The phylogenetic tree was constructed from concatenated markers

876 contained in the GTDB-tk with muscle, calculated with RAxML, and visualized in iTOL. A

877 phylogenetic tree of all 66 MAGs with reference genomes and high-quality genomes from

878 Singleton et al. constructed with concatenated markers from GTDB-tk are provided in

879 Supplementary Figure 1. Sizes of circles represent abundance patterns of metagenomic

880  reads mapping back to genomes from the same day as the metatranscriptomic

881  experiment and are not to scale. **1C)** Transcriptional patterns of each MAG in the

882  anaerobic and aerobic phases of the EBPR cycle. RNA-seq reads from each time-point

883  were competitively mapped to all 66 assembled MAGs and counts normalized by

884  transcripts per million (TPM). Total counts in the anaerobic and aerobic phases for each

885  genome were averaged separately and plotted on a log scale. Order of MAGs from left to

886  right mirrors the order of MAGs in the phylogenetic tree in 1B from the top of the circle

887  going clockwise.

888

889  **Figure 3. Clustering and Distribution of Trait Attributes Across EBPR SBR**

890  **Community Members.** Using the TbasCO method, we identified expression-based trait

891  attributes from predefined trait modules in the KEGG library and explored the distribution

892  of these trait attributes across community members. **A)** Distribution of trait-attributes

893  among sets of genomes. Bars represent the number of trait-attributes present in a set

894  number of genomes and colored by KEGG module category. Among a total of 35

895  genomes, trait attributes present between 3-18 genomes are designated as niche

896  differentiating, whereas trait attributes present in 19 or greater genomes are designated

897  as core trait attributes. Inset figure demonstrates the maximum number of attributes for

898  the maximum number of genomes. **B)** Cytoscape network showing the connectedness of

899  genomes to trait attributes. The network was filtered to only include nodes with more than

900  5 connections, therefore filtering out both genomes with few trait attributes and trait

901  attributes connected to less than 5 genomes. Genomes are represented as squares

902  colored by phylum, and trait attributes are represented as circles colored by KEGG

903    category. The size of both the squares and circles represents the number of connections

904    to that genome or trait attribute, respectively.

905

906    **Figure 4. Trait Attributes of the High-Affinity Phosphorus Transporter System**

907    ***pstABCS***

908    Using the TbasCO method, two trait attributes of the high-affinity phosphorus transporter

909    system *pstABCS* were identified. The *pstABCS* system consists of a phosphate-binding

910    protein and ABC-type transporter, and the corresponding KEGG orthologs for each

911    subunit are shown. Timepoints 1-3 refer to the three anaerobic phase timepoints, and

912    timepoints 4-6 refer to the three anaerobic phase timepoints (Figure 1). Expression values

913    are log-transformed based on setting the lowest expression value within each genome

914    across the time-series to 0 for each subunit. Specific subunits for some genomes in both

915    attributes are missing to the high cutoff thresholds for annotations. However we kept

916    genomes with 2/4 subunits to show similarities in expression profiles. The first *pstABCS*

917    trait-attribute includes microbial lineages that exhibited the highest expression of all

918    subunits towards the end of the aerobic cycle, when phosphate concentrations are

919    expected to be lowest. This includes microbial lineages within the *Actinobacteria,*

920    *Proteobacteria, Gemmatimonadetes,* and *Chloroflexi.* The second *pstABCS* trait-attribute

921    includes lineages that exhibited highest expression of all subunits upon the switch from

922    anaerobic to aerobic phases, or when phosphate concentrations are expected to be the

923    highest. This includes lineages within the *Actinobacteria* and *Proteobacteria.*

924

925    **Figure 5. Expression Dynamics of Distributed Denitrification Routes**

926   Expression of denitrification traits distributed among community members in the EBPR

927   SBR ecosystem. Timepoints 1-3 correspond to the anaerobic phase and timepoints 4-6

928   correspond to the aerobic phase as referenced in Figure 1. **A)** Complete denitrification

929   pathway and associated genetic repertoire with each sequential step. **B)** Trait attributes

930   of expression dynamics for community members with the *narGH* nitrate reductase

931   system. This trait was the only denitrification trait identified with more than one attribute.

932   **C)** Expression dynamics of the *napAB* nitrate reductase system. **D)** Expression dynamics

933   of the *norBC* nitrous oxide reductase system. **E)** Expression of all steps of denitrification

934   starting at nitrite reduction. **F)** Expression of the most complete denitrification route

935   among three community members, with the *norC* subunit for nitrous oxide reduction

936   missing. Note that OTTO1 only contains *nirS* but is included in this trait attribute because

937   the expression dynamics are similar to that of the other three genomes for this subunit.

938

939   **Figure 6. Biosynthetic Potential Compared to Expression of Amino Acid and**

940   **Vitamin Synthesis Pathways for Top 15 Expressed MAGs**

941   Biosynthetic potential and expression patterns of amino acid and vitamin pathways were

942   analyzed for the top 15 genomes with the highest transcriptional counts (Table 1). **A)** For

943   a pathway to be considered present for downstream analysis in the TbasCO pipeline,

944   80% of the pathway had to be present in a genome. Thus, we used this cutoff criterion to

945   discern whether a specific pathway was present or absent in a genome (with the

946   expectation of methionine, as all genomes did not contain at least 80% of the subunits in

947   the KEGG methionine synthase pathway, we inferred the presence of the methionine

948   synthase as presence of this pathway). Orange colored boxes for cofactor biosynthesis

949    pathways represents the presence of that pathway, whereas grey infers absence. For

950    amino acid biosynthetic pathways, amino acids are listed by their side chain groups –

951    charged, polar, hydrophobic, and other. **B)** Mini-networks of vitamin co-factors. Squares

952    are genomes with the colors matching the color bar in A. Nodes are attributes, where the

953    colored nodes for the tetrahydrofolate attributes represent the different routes. **C)** Mini-

954    networks of amino acid biosynthesis pathways split by type. Colors of nodes for each

955    amino acid represent the different routes for that pathway. Squares represent genomes

956    with colors matching the color bar in A.

957

958    **Table 1.** Genome quality statistics and relative abundance calculations for all 66 EBPR

959    SBR MAGs. Genome code names match names used in all figures and within the text.

960    Classifications were assigned using the GTDB-tk [87] and confirmed by comparing

961    against select publicly available references and a subset of HQ MAGs from Singleton et

962    al. 2021 [39]. Completeness and redundancy estimates and GC content were calculated

963    by CheckM [41]. tRNA and rRNA predictions were performed with Barrnap as part of the

964    Prokka software [88]. Relative abundance estimates reflect the proportion of reads

965    mapped to the genome in that sample divided by the total number of reads mapped to all

966    genomes        as        performed        with        SingleM.        Table        available        at

967    https://figshare.com/articles/dataset/EBPR_SBR_MAGs_Metadata/13063874.

968

969    **Table 2.** KEGG Pathways for core trait-attributes present in greater than 19 genomes.

970    **Table 3.** KEGG Pathways for differentiating trait-attributes present between 3 and 18

971    genomes.

972    **Figure 1.**



973

974

975

976

977

978

979

980

981

982

**Figure 2.**

994     **Figure 3.**



995

**Figure 4.**



TP = Timepoint    TP1, TP2, TP3 = Anaerobic Phase    TP4, TP5, TP6 = Aerobic Phase

1006    **Figure 5.**
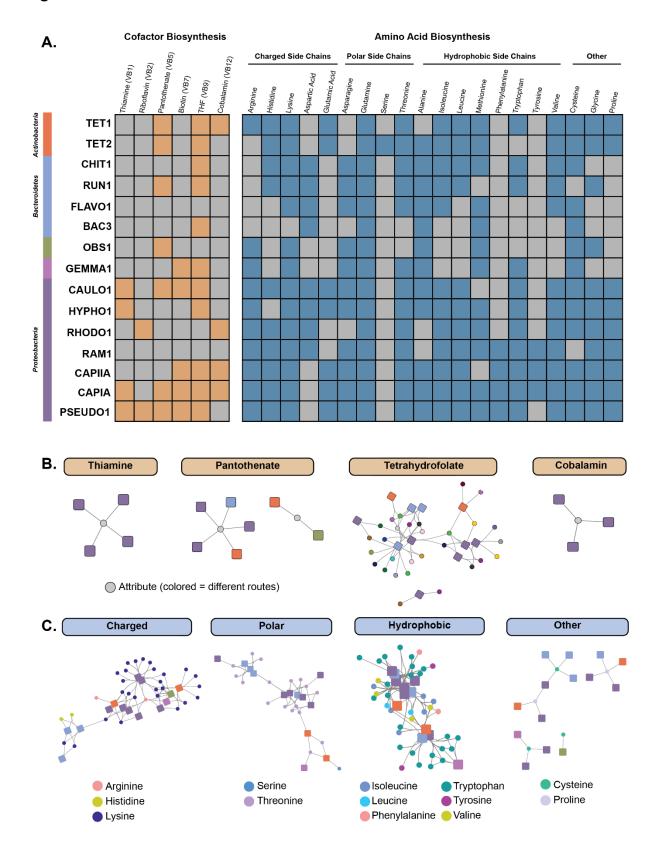


1007

1008

1009

1010 **Figure 6.**

# Table 1.

| Code | Genbank Accession | Classification | Completeness | Contamination | Size (Mbp) | Contigs | GC | Abundance 2013-5-13 | Abundance 2013-5-23 | Abundance 2013-5-28 | Total Transcriptional Reads Mapped | Total rRNAs | Total tRNAs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AUS1 | GCA_020161845.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Actinomycetales;f__Dermatophilaceae;g__Austwickia;s__ | 99.45 | 5.01 | 4.39 | 82 | 71.2 | 0.261 | 0.720 | 0.124 | 255331 | 3 | 61 |
| PHYC1 | GCA_020161815.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Actinomycetales;f__Dermatophilaceae;g__Phycicoccus;s__ | 98.02 | 0.54 | 3.06 | 34 | 71 | 1.355 | 3.007 | 0.341 | 332509 | 1 | 49 |
| PHYC2 | GCA_020161155.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Actinomycetales;f__Dermatophilaceae;g__Phycicoccus;s__Phycicoccus | 95.82 | 1.89 | 3.20 | 111 | 69.2 | 0.047 | 0.174 | 0.112 | 152031 | 1 | 52 |
| TET1 | GCA_020160805.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Actinomycetales;f__Dermatophilaceae;g__Tetrasphaera_A;s__ | 98.42 | 0.54 | 3.75 | 57 | 67.9 | 0.446 | 0.436 | 0.507 | 1378316 | 2 | 47 |
| TET2 | GCA_020160795.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Actinomycetales;f__Dermatophilaceae;g__Tetrasphaera_A;s__Tetrasphaera_A | 98.92 | 0.05 | 3.96 | 66 | 69.3 | 0.803 | 0.236 | 1.244 | 2538782 | 1 | 76 |
| LEU1 | GCA_020161315.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Actinomycetales;f__Microbacteriaceae;g__Leucobacter;s__ | 96.06 | 2.05 | 3.01 | 74 | 63.5 | 0.272 | 0.083 | 0.093 | 99061 | 3 | 47 |
| LEU2 | GCA_020161175.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Actinomycetales;f__Microbacteriaceae;g__Leucobacter;s__Leucobacter | 83.22 | 1.48 | 2.31 | 140 | 64.8 | 0.065 | 0.101 | 0.092 | 22050 | 2 | 44 |
| SAL1 | GCA_020160915.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Actinomycetales;f__Microbacteriaceae;g__Salinibacterium;s__ | 97.81 | 0 | 2.93 | 8 | 67.2 | 0.335 | 0.142 | 0.559 | 178111 | 2 | 45 |
| NANO1 | GCA_020161245.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Nanopelagicales;f__;g__;s__ | 99.14 | 3.68 | 4.29 | 95 | 72.7 | 0.106 | 0.047 | 0.172 | 64510 | 1 | 52 |
| PROP1 | GCA_020161795.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Propionibacteriales;f__Propionibacteriaceae;g__;s__ | 91.04 | 0.91 | 3.47 | 67 | 69.3 | 0.063 | 0.108 | 0.206 | 100351 | 0 | 60 |
| PROP2 | GCA_020161755.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Propionibacteriales;f__Propionibacteriaceae;g__Propionicimonas;s__ | 93.63 | 3.02 | 4.08 | 61 | 70.7 | 0.094 | 0.046 | 0.413 | 130384 | 3 | 52 |
| PROP3 | GCA_020161015.1 | d__Bacteria;p__Actinobacteriota;c__Actinobacteria;o__Propionibacteriales;f__Propionibacteriaceae;g__Propionicimonas;s__ | 94.14 | 3.15 | 3.67 | 65 | 71.6 | 0.074 | 0.176 | 0.249 | 96105 | 0 | 51 |
| FIMBRI1 | GCA_020161505.1 | d__Bacteria;p__Armatimonadota;c__Fimbriimonadia;o__Fimbriimonadales;f__Fimbriimonadaceae;g__Uphvl-Ar1;s__ | 96.55 | 0 | 3.14 | 38 | 58.8 | 0.068 | 0.234 | 0.009 | 27830 | 1 | 48 |
| BAC1 | GCA_020161835.1 | d__Bacteria;p__Bacteroidota;o__;f__;g__;s__ | 94.52 | 0 | 4.40 | 36 | 41.6 | 0.345 | 0.024 | 0.003 | 32140 | 4 | 42 |
| BAC2 | GCA_020162035.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__AKYH767;f__b-17BO;g__;s__ | 99.05 | 0.48 | 3.17 | 31 | 29.6 | 0.757 | 0.010 | 0.015 | 46346 | 3 | 52 |
| CHIT1 | GCA_020161435.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__Chitinophagales;f__Chitinophagaceae;g__;s__ | 99.01 | 0 | 4.19 | 10 | 46.3 | 0.183 | 0.174 | 3.613 | 3141341 | 0 | 34 |
| CHIT2 | GCA_020161535.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__Chitinophagales;f__Chitinophagaceae;g__Flavihumibacter;s__ | 100 | 1.23 | 4.03 | 23 | 48.2 | 0.195 | 0.383 | 0.033 | 24003 | 3 | 40 |
| SAP1 | GCA_020160935.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__Chitinophagales;f__Saprospiraceae;g__;s__ | 96.53 | 0.99 | 5.84 | 51 | 50.3 | 0.226 | 0.007 | 0.128 | 702648 | 3 | 36 |
| SAP2 | GCA_020160855.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__Chitinophagales;f__Saprospiraceae;g__OLB8;s__ | 97.52 | 0.5 | 3.73 | 65 | 37.2 | 0.290 | 0.167 | 0.016 | 10636 | 3 | 34 |
| LEAD1 | GCA_020161355.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__Cytophagales;f__Spirosomaceae;g__Leadbetterella;s__ | 99.11 | 0.6 | 4.81 | 17 | 37.7 | 0.136 | 0.002 | 0.836 | 1017458 | 2 | 36 |
| RUN1 | GCA_020161055.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__Cytophagales;f__Spirosomaceae;g__Runella;s__Runella | 100 | 0 | 7.44 | 60 | 44.4 | 0.124 | 1.088 | 1.749 | 10725342 | 2 | 40 |
| FLAVO1 | GCA_020161455.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__Flavobacteriales;f__Flavobacteriaceae;g__Flavobacterium;s__ | 99.29 | 0.35 | 3.08 | 18 | 32.5 | 0.030 | 0.002 | 0.742 | 3002991 | 3 | 36 |
| CHRYS1 | GCA_020161485.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__Flavobacteriales;f__Weeksellaceae;g__Chryseobacterium_A;s__Chryseobacterium_A | 100 | 0.25 | 2.57 | 11 | 36.7 | 0.107 | 3.917 | 0.358 | 209940 | 2 | 35 |
| BAC3 | GCA_020162015.1 | d__Bacteria;p__Bacteroidota;c__Bacteroidia;o__NS11-12g;f__UKL13-3;g__B1;s__ | 100 | 0 | 3.74 | 45 | 41.1 | 0.445 | 0.892 | 0.693 | 9991372 | 0 | 44 |
| IGNAVII | GCA_020161395.1 | d__Bacteria;p__Bacteroidota;c__Ignavibacteria;o__Ignavibacteriales;f__Ignavibacteriaceae_A;g__UTCHB3;s__ | 97.27 | 0.55 | 4.07 | 21 | 42.2 | 0.163 | 0.635 | 0.025 | 58496 | 3 | 44 |
| RTHERM1 | GCA_020160835.1 | d__Bacteria;p__Bacteroidota;c__Rhodothermia;o__Rhodothermales;f__;g__;s__ | 98.36 | 1.38 | 3.25 | 36 | 67 | 0.328 | 0.050 | 0.060 | 116472 | 3 | 52 |
| ANAER1 | GCA_020161935.1 | d__Bacteria;p__Chloroflexota;c__Anaerolineae;o__SBR1031;f__A4b;g__;s__ | 98.17 | 0 | 7.64 | 32 | 54.2 | 0.375 | 0.190 | 0.153 | 910673 | 4 | 48 |
| HERP1 | GCA_020161265.1 | d__Bacteria;p__Chloroflexota;c__Chloroflexia;o__Chloroflexales;f__Herpetosiphonaceae;g__Herpetosiphon;s__ | 99.09 | 0.91 | 6.04 | 13 | 50.2 | 0.774 | 0.025 | 0.008 | 7917 | 0 | 55 |
| OBS1 | GCA_020161235.1 | d__Bacteria;p__Cyanobacteria;c__Vampirovibrionia;o__Obscuribacterales;f__Obscuribacteraceae;g__Obscuribacter;s__Obscuribacter | 98.28 | 0.94 | 5.09 | 17 | 49.2 | 6.272 | 0.681 | 0.197 | 1713299 | 6 | 42 |
| FUSI1 | GCA_020161505.1 | d__Bacteria;p__Firmicutes_A;c__Clostridia;o__Peptostreptococcales;f__Fusibacteraceae;g__UBA5201;s__ | 96.5 | 1.75 | 3.08 | 41 | 42.8 | 0.001 | 0.580 | 0.001 | 11649 | 3 | 57 |
| GEMMA1 | GCA_020161135.1 | d__Bacteria;p__Gemmatimonadota;c__Gemmatimonadetes;o__Gemmatimonadales;f__Gemmatimonadaceae;g__;s__ | 98.35 | 3.3 | 4.55 | 8 | 70.1 | 0.004 | 0.031 | 0.494 | 2624259 | 3 | 55 |
| SACCH1 | GCA_020160975.1 | d__Bacteria;p__Patescibacteria;c__Saccharimonadia;o__Saccharimonadales;f__Saccharimonadaceae;g__Saccharimonas;s__ | 84.48 | 0 | 0.97 | 1 | 49.6 | 0.637 | 1.437 | 0.035 | 65079 | 3 | 43 |
| ALPHA1 | GCA_020161965.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__;f__;g__;s__ | 82.43 | 2.65 | 3.94 | 581 | 64.6 | 0.015 | 0.165 | 0.007 | 1283274 | 3 | 39 |
| CAED1 | GCA_020161545.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Caedimonadales;f__UBA1908;g__;s__ | 86.36 | 1.1 | 1.88 | 96 | 52.8 | 0.034 | 0.201 | 0.002 | 41264 | 3 | 45 |
| BREV1 | GCA_020161595.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Caulobacterales;f__Caulobacteraceae;g__Brevundimonas;s__Brevundimonas | 97.51 | 3.41 | 3.07 | 155 | 67.2 | 0.011 | 0.254 | 0.004 | 27852 | 2 | 45 |
| CAULO1 | GCA_020161365.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Caulobacterales;f__Caulobacteraceae;g__Caulobacter;s__ | 100 | 0 | 4.43 | 25 | 66.9 | 0.048 | 0.093 | 0.589 | 4627825 | 3 | 55 |
| HYPHO1 | GCA_020161405.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Caulobacterales;f__Hyphomonadaceae;g__UBA1942;s__ | 98.43 | 0.32 | 2.98 | 6 | 39.4 | 0.844 | 0.006 | 2.208 | 4138107 | 0 | 33 |
| REYR1 | GCA_020160955.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Reyranellales;f__Reyranellaceae;g__Reyranella;s__ | 89.96 | 7.34 | 5.08 | 210 | 70 | 0.057 | 0.090 | 0.238 | 224063 | 3 | 53 |
| REYR2 | GCA_020160995.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Reyranellales;f__Reyranellaceae;g__Reyranella;s__Reyranella | 91.04 | 6.01 | 5.71 | 258 | 65.3 | 0.074 | 0.102 | 0.134 | 62018 | 1 | 53 |
| ANDERS1 | GCA_020161855.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhizobiales;f__Anderseniellaceae;g__PALSA-927;s__ | 97.64 | 0.4 | 3.36 | 19 | 61.6 | 0.187 | 0.175 | 0.029 | 25238 | 2 | 44 |
| BEIJ1 | GCA_020161915.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhizobiales;f__Beijerinckiaceae;g__Bosea;s__ | 81.6 | 8.48 | 4.44 | 777 | 66.3 | 0.156 | 0.319 | 0.423 | 338238 | 0 | 43 |
| BEIJ2 | GCA_020161975.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhizobiales;f__Beijerinckiaceae_A;g__;s__ | 81.18 | 5.25 | 3.99 | 465 | 62.5 | 0.042 | 0.157 | 0.018 | 28432 | 0 | 41 |
| BEIJ3 | GCA_020161475.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhizobiales;f__Beijerinckiaceae_A;g__PAR1;s__ | 76.21 | 1.72 | 3.08 | 320 | 63.3 | 0.017 | 1.744 | 0.099 | 77102 | 0 | 41 |
| BEIJ4 | GCA_020161575.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhizobiales;f__Beijerinckiaceae_A;g__PAR1;s__ | 97.89 | 0 | 3.19 | 17 | 63.2 | 0.176 | 0.538 | 0.014 | 26820 | 0 | 47 |
| PHREA1 | GCA_020161695.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhizobiales;f__Phreatobacteraceae;g__Phreatobacter;s__ | 98.35 | 3.96 | 4.69 | 38 | 67.7 | 0.022 | 0.273 | 0.103 | 133243 | 1 | 50 |
| RHIZO1 | GCA_020161035.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhizobiales;f__Rhizobiaceae;g__Aminobacter;s__Aminobacter | 94.26 | 5.5 | 5.50 | 80 | 63.8 | 0.136 | 0.095 | 0.095 | 219213 | 3 | 48 |
| RHIZO2 | GCA_020161665.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhizobiales;f__Rhizobiaceae;g__QFOR01;s__ | 88.41 | 2.12 | 3.39 | 43 | 60.6 | 0.035 | 0.335 | 0.003 | 24536 | 0 | 47 |
| RHIZO3 | GCA_020161615.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhizobiales;f__Rhizobiaceae;g__Shinella;s__Shinella | 78.53 | 6.03 | 6.98 | 935 | 63.6 | 0.010 | 0.169 | 0.037 | 149921 | 0 | 43 |
| RHODO1 | GCA_020161655.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhodobacterales;f__Rhodobacteraceae;g__Defluviimonas;s__ | 100 | 0.35 | 4.08 | 24 | 65.5 | 0.321 | 0.141 | 0.848 | 3645270 | 0 | 44 |
| RHODO2 | GCA_020161615.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhodobacterales;f__Rhodobacteraceae;g__Pararhodobacter;s__ | 99.09 | 1.19 | 4.87 | 26 | 67.9 | 0.084 | 0.534 | 0.009 | 25807 | 1 | 45 |
| RHODO3 | GCA_020160875.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rhodospirillales_C;f__Rhodospirillaceae_A;g__;s__ | 91.2 | 2.27 | 3.76 | 236 | 62.2 | 0.153 | 0.046 | 0.185 | 153017 | 1 | 39 |
| RICK1 | GCA_020160775.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Rickettsiales;f__Rickettsiaceae;g__GCA-2402195;s__ | 75.59 | 1.58 | 1.18 | 82 | 34.5 | 0.085 | 0.075 | 0.052 | 17671 | 2 | 45 |
| SPHING1 | GCA_020160755.1 | d__Bacteria;p__Proteobacteria;c__Alphaproteobacteria;o__Sphingomonadales;f__Sphingomonadaceae;g__Sphingopyxis;s__ | 99.98 | 1.56 | 4.31 | 20 | 65.1 | 0.026 | 0.014 | 0.607 | 600695 | 3 | 47 |
| ALIC1 | GCA_020161945.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Burkholderiales;f__Burkholderiaceae;g__Alicycliphilus;s__ | 99.64 | 1.04 | 3.83 | 33 | 66.3 | 0.166 | 2.959 | 0.738 | 770970 | 1 | 46 |
| OTTO1 | GCA_020161215.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Burkholderiales;f__Burkholderiaceae;g__Ottowia;s__ | 93.66 | 5.56 | 4.52 | 250 | 67.1 | 0.011 | 0.276 | 0.004 | 26717 | 1 | 46 |
| OTTO2 | GCA_020161715.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Burkholderiales;f__Burkholderiaceae;g__Ottowia;s__Ottowia | 99.26 | 0.62 | 3.40 | 35 | 69.1 | 0.372 | 4.140 | 0.424 | 121379 | 1 | 45 |
| RAM1 | GCA_020161775.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Burkholderiales;f__Burkholderiaceae;g__Ramlibacter;s__ | 99.84 | 0.06 | 4.36 | 32 | 66.1 | 0.778 | 0.536 | 1.814 | 1832037 | 1 | 45 |
| RUBRI1 | GCA_020161065.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Burkholderiales;f__Burkholderiaceae;g__Rubrivivax;s__ | 99.52 | 0.05 | 6.29 | 41 | 71.2 | 0.236 | 0.347 | 0.306 | 1259737 | 1 | 73 |
| VITREO1 | GCA_020161145.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Burkholderiales;f__Burkholderiaceae;g__Vitreoscilla_A;s__ | 100 | 0.7 | 3.51 | 13 | 68.9 | 0.397 | 4.498 | 0.536 | 382529 | 1 | 46 |
| CAPIA | NA | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Burkholderiales;f__Rhodocyclaceae;g__Accumulibacter;s__Accumulibacter | 100 | 0.03 | 4.59 | 61 | 63.8 | 18.797 | 10.533 | 0.106 | 2411395 | 0 | 46 |
| CAPIIA | NA | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Burkholderiales;f__Rhodocyclaceae;g__Accumulibacter;s__Accumulibacter | 99.84 | 0.24 | 4.64 | 81 | 64.3 | 33.479 | 26.824 | 49.334 | 102762132 | 0 | 44 |
| ZOO1 | GCA_020161115.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Burkholderiales;f__Rhodocyclaceae;g__Zoogloea;s__ | 91.62 | 3.51 | 4.99 | 501 | 65.7 | 0.090 | 0.026 | 0.106 | 913411 | 4 | 59 |
| LEG1 | GCA_020161725.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Legionellales;f__Legionellaceae;g__Legionella;s__ | 92.74 | 1.07 | 2.58 | 182 | 36.1 | 0.094 | 0.126 | 0.006 | 19591 | 1 | 27 |
| LUTEI1 | GCA_020161335.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Xanthomonadales;f__Xanthomonadaceae;g__Luteimonas;s__ | 96.89 | 0.71 | 3.56 | 252 | 69.9 | 0.002 | 0.309 | 0.011 | 49418 | 1 | 39 |
| PSEUDO1 | GCA_020160895.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Xanthomonadales;f__Xanthomonadaceae;g__Pseudoxanthomonas_A;s__ | 99.95 | 0.89 | 3.67 | 28 | 67.8 | 0.416 | 0.730 | 3.125 | 3964795 | 2 | 50 |
| PSEUDO2 | GCA_020161075.1 | d__Bacteria;p__Proteobacteria;c__Gammaproteobacteria;o__Xanthomonadales;f__Xanthomonadaceae;g__Pseudoxanthomonas;s__ | 99.02 | 0 | 2.99 | 6 | 69.6 | 1.750 | 6.111 | 1.228 | 515369 | 3 | 52 |

File at https://figshare.com/account/projects/90437/articles/13063874

**Table 2.**

| Module Description | Number of Attributes |
|---|---|
| Citrate cycle, second carbon oxidation, 2-oxoglutarate => oxaloacetate [PATH:map00020 map01200 map01100] | 13 |
| Citrate cycle (TCA cycle, Krebs cycle) [PATH:map00020 map01200 map01100] | 10 |
| Shikimate pathway, phosphoenolpyruvate + erythrose-4P => chorismate [PATH:map00400 map01230 map01100 map01110] | 8 |
| Fatty acid biosynthesis, initiation [PATH:map00061 map01212 map01100] | 7 |
| Glycolysis, core module involving three-carbon compounds [PATH:map00010 map01200 map01230 map01100] | 7 |
| Adenine ribonucleotide biosynthesis, IMP => ADP,ATP [PATH:map00230 map01100] | 4 |
| Guanine ribonucleotide biosynthesis IMP => GDP,GTP [PATH:map00230 map01100] | 4 |
| Inosine monophosphate biosynthesis, PRPP + glutamine => IMP [PATH:map00230 map01100] | 4 |
| Isoleucine biosynthesis, threonine => 2-oxobutanoate => isoleucine [PATH:map00290 map01230 map01100] | 3 |
| NADH:quinone oxidoreductase, prokaryotes [PATH:map00190] | 3 |
| beta-Oxidation, acyl-CoA synthesis [PATH:map00061 map00071 map01212 map01100] | 2 |
| F-type ATPase, prokaryotes and chloroplasts [PATH:map00190 map00195] | 2 |
| Valine/isoleucine biosynthesis, pyruvate => valine / 2-oxobutanoate => isoleucine [PATH:map00290 map00770 map01210 map01230 map01100 map01110] | 2 |
| CAM (Crassulacean acid metabolism), dark [PATH:map00620 map00710 map01200 map01100 map01120] | 1 |
| Cytochrome c oxidase, cbb3-type [PATH:map00190] | 1 |
| Cytochrome c oxidase, prokaryotes [PATH:map00190] | 1 |
| dTDP-L-rhamnose biosynthesis [PATH:map00521 map00523 map01100 map01130] | 1 |
| Leucine biosynthesis, 2-oxoisovalerate => 2-oxoisocaproate [PATH:map00290 map01210 map01230 map01100 map01110] | 1 |
| Phosphatidylethanolamine (PE) biosynthesis, PA => PS => PE [PATH:map00564 map01100] | 1 |
| PRPP biosynthesis, ribose 5P => PRPP [PATH:map00030 map00230 map01200 map01230 map01100] | 1 |
| Pyruvate oxidation, pyruvate => acetyl-CoA [PATH:map00010 map00020 map00620 map01200 map01100] | 1 |
| Semi-phosphorylative Entner-Doudoroff pathway, gluconate => glycerate-3P [PATH:map00030 map01200 map01100 map01120] | 1 |
| Threonine biosynthesis, aspartate => homoserine => threonine [PATH:map00260 map01230 map01100 map01110] | 1 |

**Table 3.**

| Module_description | Number of Attributes |
|---|---:|
| Glycolysis (Embden-Meyerhof pathway), glucose => pyruvate [PATH:map00010 map01200 map01100] | 279 |
| Citrate cycle (TCA cycle, Krebs cycle) [PATH:map00020 map01200 map01100] | 208 |
| Gluconeogenesis, oxaloacetate => fructose-6P [PATH:map00010 map00020 map01100] | 76 |
| Inosine monophosphate biosynthesis, PRPP + glutamine => IMP [PATH:map00230 map01100] | 45 |
| Citrate cycle, second carbon oxidation, 2-oxoglutarate => oxaloacetate [PATH:map00020 map01200 map01100] | 31 |
| Heme biosynthesis, plants and bacteria, glutamate => heme [PATH:map00860 map01100 map01110] | 27 |
| Tetrahydrofolate biosynthesis, GTP => THF [PATH:map00790 map00670 map01100] | 25 |
| Tryptophan biosynthesis, chorismate => tryptophan [PATH:map00400 map01230 map01100 map01110] | 25 |
| Ornithine biosynthesis, glutamate => ornithine [PATH:map00220 map01210 map01230 map01100] | 24 |
| Histidine biosynthesis, PRPP => histidine [PATH:map00340 map01230 map01100 map01110] | 17 |
| Pentose phosphate pathway (Pentose phosphate cycle) [PATH:map00030 map01200 map01100 map01120] | 16 |
| Lysine biosynthesis, succinyl-DAP pathway, aspartate => lysine [PATH:map00300 map01230 map01100] | 12 |
| Uridine monophosphate biosynthesis, glutamine (+ PRPP) => UMP [PATH:map00240 map01100] | 11 |