

001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046

microbetag: simplifying microbial network interpretation through annotation, enrichment and metabolic complementarity analysis

Haris Zafeiropoulos¹, Ermis Ioannis Michail Delopoulos¹,
Andi Erega², Annelies Geirnaert², John Morris³, Karoline Faust^{1*}

^{1*} Department of Microbiology, Immunology and Transplantation, Rega Institute for Medical Research , KU Leuven, Herestraat, Leuven, 3000, , Belgium .

² Institute of Food, Nutrition and Health, ETH Zurich, Street, Zurich, 8092, , Switzerland .

³ Department of Pharmaceutical Chemistry, University of California San Francisco, Street, San Francisco, 94143, California, USA .

*Corresponding author(s). E-mail(s): karoline.faust@kuleuven.be;

Contributing authors: haris.zafeiropoulos@kuleuven.be;

ermisioannis.michaildelopoulos@student.kuleuven.be;

andi.errega@hest.ethz.ch; annelies.geirnaert@hest.ethz.ch;

scooter@cgl.ucsf.edu;

Abstract

*

Up to 350 words.

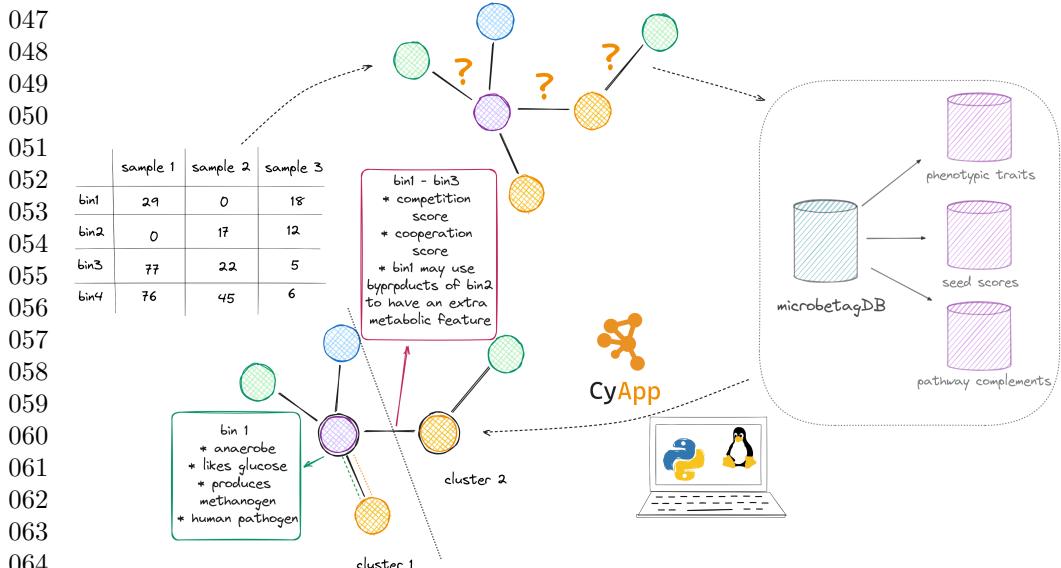
The abstract must include the following separate sections:

Background: the context and purpose of the study

Results: the main findings

Conclusions: a brief summary and potential implications

* Looks like Chris Quince is our editor.



057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092

Figure abstract.

Keywords: microbial associations, enrichment analysis, data integration, pathway complementarity, seed set

Background ¹ ²

Microbial ecology plays a fundamental role in the stability and resilience of ecosystems and their processes; from soils, aquatic environments and biogeochemical cycles [1] to host-associated environments and the human health [2, 3]. Most microbial species live only in communities [4] and most natural microbial communities consist of hundreds or even thousands of species [5]. Each species exhibits a unique repertoire of reactions and adapts to various niches, each with specific nutrient and environmental requirements. Understanding the dynamics governing interactions among microbial species and their relationships with the surrounding environment would shed light in several aspects microbial ecology [6].

Metabolic conditions within an environment influence the composition and distribution of its functions (metabolic niche effect). At the same time, processes like genetic diversity and stochastic events, induce variability within different functional groups of the community. Recent studies suggest that the relative importance of those two processes drive the extent to which functional and taxonomic composition at the

¹We are to submit in the Microbiome journal as a "Software" manuscript, thus we follow [these rules](#)

²The introduction should not include subheadings. The Background section should explain the relevant context and the specific issue that the software described is intended to address.

community level are correlated [7]. +++++++	093
including cooperation, competition, parasitism, commensalism and ammensalism [3].	094
High-throughput sequencing (HTS) has provided great insight into the diversity	095
and composition of microbial communities [8]. Uncultivated species can now be	096
detected and their feautues can be inferred through their genomic information [9]. More-	097
over, the composition of thousands of microbiome samples is now accessible allowing	098
for the inference of patterns among sets of samples. A widely used approach to extraxt	099
such patterns, is the creation of co-occurrence networks based on metagenomic read	100
data (amplicon and/or shotgun) [10]. A great number of approaches is available for	101
co-occurrence network inference based on a range of statistical concepts such as: cor-	102
relation (e.g., CoNet [11], SparCC [12]), linear regression (e.g., SpiecEasi [13]) and	103
causal inference (FlashWeave [14]). Nevertheless, microbial co-occurrence networks	104
continue to encounter various challenges [15], encompassing issues associated with data	105
analysis and network construction, leading to tool-dependent analysis [4, 16, 17]. But	106
also, challenges regarding the interpretation of the networks. Addressing the question	107
of <i>What can we learn from the hairballs</i> posed by Röttjers et al. [4] could provide	108
essential insight on the mechanisms of the interactions.	109
The use of microbial network inference as a means for predicting interactions	110
has underscored its limited accuracy, and the fact that the biological implications	111
of network properties remain unclear [18]. Theoretical principles derived from net-	112
work studies might provide indications of emergent biological characteristics [4, 19].	113
For example, modules (highly interconnected nodes) within microbial co-occurrence	114
networks could serve as indicators of ecological processes that govern community	115
structure, including niche filtering and habitat preference[20]. Data integration and	116
clustering have been suggested to address this challenge [15]. Clusters identified in	117
microbial association networks have demonstrated their ability to mirror key drivers	118
of community composition [21] and sevaral algorithms and implementations are avail-	119
able [22]. However, data integration approaches in microbial co-occurrence networks	120
are so-far limited. Here, we present microbetag , a microbial co-occurrence network	121
annotator that exploits several channels of information to enhance/diminish the con-	122
fidence of the associations suggested by the network and generate hypotheseses for	123
further investigation both at the taxon pair and the community level.	124
microbetag serves as a comprehensive platform that provides information on taxa	125
along with their potential metabolic interactions from multiple channels (see Imple-	126
mentation 3). The key concept here is the reverse ecology approach <i>reverse ecology</i> [23].	127
Reverse ecology leverages genomics to explore community ecology with no <i>a priori</i>	128
assumptions about the taxa involved. Making the most of advancements in systems	129
biology and genomic metabolic modeling, as well as system-level analysis of intricate	130
biological networks, the reverse ecology framework enables the prediction of ecologi-	131
cal traits for less-understood microorganisms, their interactions with others, and the	132
overall ecology of microbial communities [24].	133
A metabolic network's "seed set" is the set of compounds that, based on the net-	134
work topology, need to be acquired exogenously [25] (see Figure 1). Such nodes might	135
be independent, i.e. they cannot be activated by any other node in the network, or	136
they can be interdependent forming groups of seed nodes. Seeds are a useful proxy	137
	138

139 for the habitat of the organism and an essential tool in the frameowrk of reverse ecol-
140 ogy [25, 26]. Based on the seed concept, several graph theory-based metrics (indices)
141 have been described to predict species interactions directly from their networks' topolo-
142 gies [27–30]. Over the last years, the seed approach has been implemented at the
143 Genome-scale metabolic network reconstructions (GENREs) level. GENREs encapsu-
144 late mathematical representations capturing the biochemical reactions that could take
145 place within an organism [31–33].

146 Metabolic complementarity among species, serving as a reflection of potential
147 cooperation within communities, assesses the capacity for collaboration; cross-feeding
148 or syntrophy interactions are typical examples of such a collaboration. In contrast,
149 metabolic competition refers to the metabolic overlap between two species leading to
150 exploitative competition, e.g. for nutrient resources. Seed and non-seed sets can be
151 used to compute such indices. Thorough examination of such complements can reveal
152 metabolic interactions leading to patterns observed on the co-occurrence network.

153 Considering complementarity as a range of alternatives and pairwise microbial
154 interactions in the context of the community as a whole, microbial species may also
155 exchange metabolic compounds that may be not seeds at a certain time but may
156 allow them to perform functions that they are currently not capable of [34, 35]. Such
157 by-products may be even metabolites not even necessary for themeselves but for the
158 community as a whole [36]. To explore the potenial of a species metabolism given
159 they benefit from a partner of theirs, genome annotations combined with collections
160 of functional units to highlight can provide a valid proxy. We present here a naive
161 approach exporting all possible complements between a pair of species based on their
162 KEGG ORTHOLOGY (KOs) annotations and the KEGG MODULES database [37].

163 **microbetag** annotates a user's co-occurrence network by integrating phenotypic
164 traits on the taxa present on the network (nodes) and potential metabolic interactions
165 to their suggested associations (edges). A Graphical User Interface (GUI) is supported
166 as a CytoscapeApp providing a user-friendly environment to investigate annotations
167 in a straightforward way. All annotations present in microbetagDB are also available
168 through an Application Programming Interface (API). **microbetag** 's source code is
169 distributed under a GNU GPL v3 license and available on GitHub. Documentation
170 and further support on how to use **microbetag** is available at [documentation web-site](#).
171 To the best of our knowledge there is not a software with which **microbetag** could
172 be compared with directly. To validate our annotations we used a recently published
173 network with partially known interactions between some pairs of species found associ-
174 ated [38] (see Results section, paragraph 3). To demonstrate **microbetag** 's potential,
175 we present the main features of its interface and we discuss a real-world use-case (see
176 Discussion section, paragraph 3).

177

178

179

180

181

182

183

184

Implementation	³	185
		186
		187
		188
		189
		190
		191
		192
		193
		194
		195
		196
		197
		198
		199
		200
		201
		202
		203
		204
		205
		206
		207
		208
		209
		210
		211
		212
		213
		214
		215
		216
		217
		218
		219
		220
		221
		222
		223
		224
		225
		226
		227
		228
		229
		230
Genomes included		
Using the Genome Taxonomy Database (GTDB) v207 metadata files, we retrieved the NCBI genome accessions of the high quality representative genomes, i.e. completeness $\geq 95\%$ and contamination $\leq 5\%$. A set of 26,778 genomes was obtained, representing 22,009 unique NCBI Taxonomy IDs. Using these accession numbers, we were able to download their corresponding .faa files when available leading to a set of 16,900 amino acid sequence files. The latter were annotated and used to obtain potential pathway complementarities between pairs of genomes (see paragraph 3). Last, when available, their corresponding annotations on PATRIC database [39] were retrieved to reconstruct GENREs (see paragraph 3).		
Taxonomy schemes		
microbetag maps the taxonomy of each entry in the abundance table to their corresponding NCBI Taxonomy ID and, if available, their closest GTDB representative genome(s), since several GTDB representative genomes may map to the same NCBI Taxonomy ID. Two well established taxonomy schemes are supported: the GTDB [40] that is being broadly used for bins and/or MAGs taxonomical classification and the Silva database [41] that is widely used in amplicon studies. Both taxonomy schemes link their taxonomies to NCBI Taxonomy IDs [42]. In case none of those two taxonomies was used and the abundance table contains less than 1,000 taxa, microbetag maps the user provided taxonomies to NCBI Taxonomy. To this end, microbetag makes use of the fuzzywuzzy library that implements the Levenshtein Distance Metric to get the closest NCBI taxon name and thus its corresponding NCBI Taxonomy ID; a relatively high similarity score is used (90) to avoid false positives. Also, using the nodes dump file of NCBI Taxonomy, microbetag may retrieve the child taxa of a taxon in user's data, along with their corresponding NCBI Taxonomy IDs, if requested by the user. If the user provides their abundance table with taxonomies already mapped to the GTDB taxonomy, microbetag will report the best possible annotations in a time efficient manner.		
Network inference		
When a co-occurrence network is not provided by the user, microbetag exploits FlashWeave [14] to build one on the fly. Yet, microbetag supports the annotation of networks built from any algorithm/software, in any format Cytoscape can load.		
microbetag pre-processing		
In order to aid the user to map their sequences to the GTDB taxonomy, DADA2-formatted 16S rRNA gene sequences for both bacteria and archaea [43] were used to train the TAXID classifier of the DECIPHER package [44] and are available through		

³This should include a description of the overall architecture of the software implementation, along with details of any critical issues and how they were addressed.

231 the [microbetag preprocess Docker image](#). Likewise, when the abundance table consists
232 of more than 1,000 taxa, providing a network as an input is mandatory. Again, to help
233 the user, [microbetag](#) preprocess Docker image supports the inference of a network
234 using FlashWeave.

235

236 **Literature based nodes annotation**

237 Using a set of Tara Ocean samples [45] FAPROTAX [46] estimates the functional
238 potential of the bacterial and archaeal communities, by classifying each taxonomic unit
239 into functional group(s) based on current literature, descriptions of cultured represen-
240 tatives and/or manuals of systematic microbiology. In this manually curated approach,
241 a taxon is associated with a function if and only if all the cultured species within the
242 taxon have been shown to exhibit that function. In its current version, FAPROTAX
243 includes more than 80 functions based on 7600 functional annotations and covering
244 more than 4600 taxa. Contrary to gene content based approaches, e.g. PICRUSt2 [47],
245 FAPROTAX estimates metabolic phenotypes based on experimental evidence.
246

247 [microbetag](#) invokes the accompanying script of FAPROTAX and converts the
248 taxonomic microbial community profile of the samples included in the user's abun-
249 dance table or of the taxa present in the provided network, into putative functional
250 profiles. Then, it parses FAPROTAX's subtables to annotate each taxonomic unit
251 present in the user's data with all the functions for which they had a hit. FAPROTAX
252 annotations are not part of the microbetagDB but are computed on the fly.
253

254 **Genomic based nodes annotation**

255 phenDB [48] is a publicly available resource that supports the analysis of bacterial
256 (meta)genomes to identify 47 distinct functional traits, e.g. whether a species is pro-
257 ducing butanol or has an halophilic lifestyle. It relies on support vector machines
258 (SVM) trained with manually curated datasets based on gene presence/absence pat-
259 terns for trait prediction. More specifically, the model for a particular trait is trained
260 using a collection of EggNOG annotated genomes where the knowledge of whether
261 that trait is present or absent among its members is available. These models (classi-
262 fiers) are used to predict presence/absence of their corresponding traits in non-studied
263 species.

264 In the framework of microbetagDB, classifiers were re-trained using the genomes
265 provided by phenDB for each trait to sync with the latest version of eggNOG [49]
266 and the [phenotrex](#) [48] software tool. Genomes were downloaded from NCBI using
267 the [Batch Entrez](#) program. Then, *genotype* files were produced for all the high quality
268 GTDB representative genomes. Each model was then used against all the GTDB
269 *genotype* files to annotate each with the presence or the absence of the trait. A list of all
270 the phenotypic traits available for the genomes present in microbetagDB is available
271 on [microbetag](#)'s [documentation site](#). The updated models are also available
272

273 **Pathway complementarity**

274

275 To infer potential pathway complementarities we consider the modules described in
276 KEGG MODULES database [37]. A KEGG module is defined as a functional unit

within the KEGG framework that represents a set of enzymes and reactions involved in a specific biological process or pathway [50]. Such a unit consists of several *steps*, each of which may have more than one molecular ways to occur (Figure 1). A module's definition is a logical expression and consists of KOs that may be coupled with one another as: a. connected steps of the pathway b. parts of a molecular complex, c. alternatives of the same step, and d. optional entities of a complex. Both (a) and (b) cases should be considered as the AND logical operator, while (c) would be the OR (Figure 1). Given a module's definition, we will consider as an *alternative* any subset of the KO terms mentioned in the definition, that has exactly one way to perform each step, provided that all the steps of the module are covered. We define a genome as having a *complete* module, if and only if all of the KOs of at least one alternative are present on the genome.

Within this framework, *kofamscan* [51] was used to annotate with KEGG ORTHOLOGY terms (KOs) the 16,900 high quality GTDB representative genomes for which a .faa was available [52]. The KOs of each genome were then mapped to their corresponding KEGG modules; a KO may map to more than one modules (1 : n).

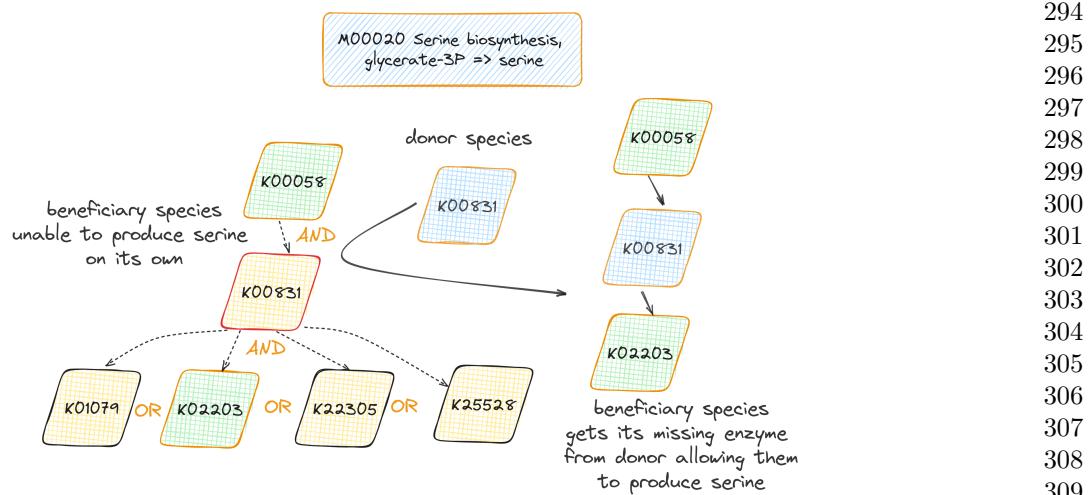


Fig. 1: Pathway complementarity approach. The high quality GTDB genomes were annotated with KEGG ORTHOLOGY (KO) terms. The various ways of getting a KEGG module complete were enumerated and all the possible ways a donor species could "fill" a beneficiary's non-complete module were calculated. In this case, there are 4 unique ways for having the serine biosynthesis module complete; in all of them K00831 is required. However, it is missing from the beneficiary species that supports the 2 out of the 3 steps of the module's definition. A donor species having and potentially sharing the corresponding enzyme of K00831 may enable the beneficiary species to produce serine.

All module definitions were retrieved using the KEGG API and parsed to enumerate their alternatives. Each pair of the KEGG annotated genomes was then

323 investigated for potential pathway complementarities, i.e. whether a genome lacking a
324 number of KOs ($genome_A$) to have a complete module ($module_x$) could benefit from
325 another's species genome(s) ($genome_B$). In that case, $genome_B$ does not necessarily
326 have a complete alternative of $module_x$; as long as it has the missing KOs that
327 $genome_A$ needs to complete an alternative of it, $genome_B$ potentially complements
328 $genome_A$ with respect to $module_x$. In total, 341,568 unique complementarities were
329 exported.

330 Thanks to the graphical user interface (GUI) of the [KEGG pathway map](#)
331 [viewer](#) [53, 54], each complementarity can be visualised as part of the closest KEGG
332 metabolic map; where the KOs contributed by the donor are shown in blue-green
333 whereas those coming from the beneficiary genome are coloured in rose.

334 `microbetag` annotates the edges of a co-occurrence network by identifying pairs
335 where both taxa map to an annotated genome present on microbetagDB. Since
336 co-occurrence networks are undirected, both nodes of a suggested association are
337 considered as potential donors and beneficiary species. When more than one GTDB
338 representative genome map to the same NCBI Taxonomy Id all the possible genome
339 combinations are considered. Finally, two edges are added in such pairs of taxa in the
340 annotated network: one considering $species_A$ as the potential beneficiary and $species_B$
341 as the potential donor species, and one vice-versa.
342

343 Seed scores using genome scale metabolic reconstructions

344 The Metabolic Complementarity Index ($MI_{Complementarity}$) measures the degree to
345 which two microbial species can mutually assist each other by complementing each
346 other's biosynthetic capabilities. As described in [55], it is defined as the proportion
347 of seed compounds of a species that can be synthesized by the metabolic network of
348 another, but are not included in the seed set of the latter. $MI_{Complementarity}$ offers
349 an upper bound assessment of the potential for syntrophic interactions between two
350 species. Further, the Metabolic Competition Index ($MI_{Competition}$) represents the sim-
351 ilarity in two species' nutritional profiles. This index establishes an upper limit on the
352 level of competition that one species may face from another. Those indices have been
353 thoroughly described and implemented in the NetCooperate [27] and NetCompt [28]
354 tools correspondingly. We will be referring to those two indices as "seed scores".
355

356 Recently, the PhyloMinttool [55] was released supporting the calculation of the
357 seed scores of GENREs in SBML format.

358 In the framework of `microbetag`, seed scores were computed using `PhyloMint` and
359 draft GENREs for all pairwise combinations of GTDB representative genomes that
360 have been RAST annotated in the framework of the PATRIC database [39]. GENREs
361 were reconstructed using the Model SEED pipeline [56] through its Python interface
362 [ModelSEEDpy](#).

363 Moreover, the computed seed and the non-seed (i.e., set of metabolic compounds
364 a genome can build on its own) sets of each genome were used to compute their
365 overlap among all the pairwise combinations of those genomes. More specifically, the
366 overlap of $seed\ set_{species_A}$ with the $non\ seed\ set_{species_B}$ was retrieved. `microbetag`
367 then annotates again the edges of the co-occurrence network where both taxa have
368

been mapped to a at least one GTDB genome, mentioning all the KEGG maps for which there is at least one seed compound of the potentially beneficiary species	369 370 371 372 373 374 375 376 377 378 379 380 381 382
Clustering network	371 372 373 374 375 376 377 378 379 380 381 382
manta is a heuristic network clustering algorithm that clusters nodes within weighted networks effectively, leveraging the presence of negative edges and discerning between weak and microbetag invokes manta [22] to infer clusters from the microbial network.	373 374 375 376 377 378 379 380 381 382
A taxonomically-informed layout is	375 376 377 378 379 380 381 382
strong cluster assignments. ++ taxonomy layout	376 377 378 379 380 381 382
Groups of annotations	379 380 381 382
Biologically meaningful groups were built using the micrO ontology [57].	380 381 382
Building the CytoscapeApp	383 384
The <code>microbetag</code> CytoscapeApp was build based on the source code of the scVizNet [58]. Java @Ermis to add	385 386
Enrichment analysis is supported. Hypergeometric distribution FDR +++	387 388
Dependencies, Web server and API	389 390
The <code>microbetag</code> web service is container - based and consists of three Docker [59] (v24.0.2) images: a. the MySQL database b. an nginx [60] web server and c. the app itself. The latter uses Gunicorn (20.1.0) to build an application server which communicates with the web server using the Web Server Gateway Interface (WSGI) protocol and handles incoming HTTP requests. <code>microbetag</code> is implemented as a Flask application (v2.3.2); Flask is a micro web framework for developing Python web applications and RESTful APIs. A thorough description of <code>microbetag</code> 's API is available at the ReadTheDocs web site . The source code of the <code>microbetag</code> web service is available on GitHub .	391 392 393 394 395 396 397 398 399
python 3.11 slim docker image julia 1.7.1 for flashweave mysql.connector 8.0.27	400
python library pandas 2.1.1. numpy 1.26.0 multiprocessing	401
text processing using awk	402
KEGG API	403
	404
	405
	406
	407
	408
	409
	410
	411
	412
	413
	414

415 **Results**⁴

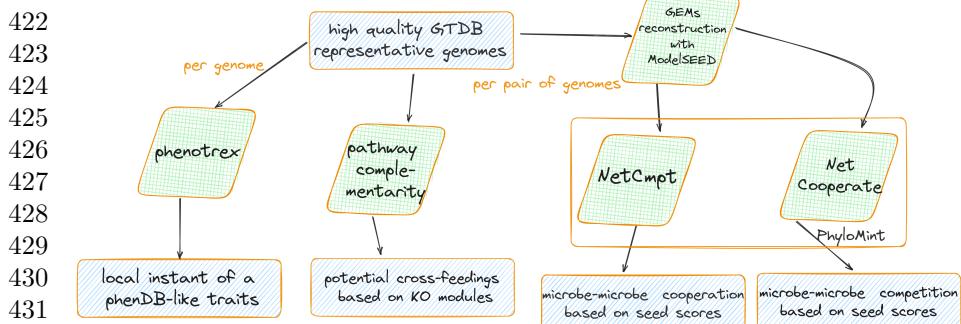
416 **microbetag and microbetagDB**

418

419

420 **PRE-CALCULATIONS**

421



422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459 **Fig. 2:** Diagram of the **microbetag** pre - calculations and the on the fly workflow.

460 GTDB v207 representative genomes were filtered and for those of high-quality 33

461 phenotypic traits were predicted using **phenotrex**. To this end, models were re-trained

462 to sync with recent version of eggNOG.

10

⁴Significant advance over previously published software (usually demonstrated by direct comparison with available related software) This should include the findings of the study including, if appropriate, results of statistical analysis which must be included either in the text or as tables and figures. This section may be combined with the Discussion section for Software articles.

Table 1: Summary of Data⁵

Description	Entries		
GTDB representative genomes	34,608	461	
Phen-model-oriented metabolic functions	32	462	
FAPROTAX functions	92	463	
Unique complement*	341,568	464	
GENREs leading to ~ 1 billion competition and complementarity scores	30,755	465	
		466	
		467	
		468	
		469	
annotated network returned in .cyjs format		470	
For a computationally efficient way to annotate large networks, a Docker image is provided so the user runs a taxonomy assignment using the IDTAXA algorithm [44] of the DECIPHER R package [61]. A co-occurrence network is also built using FlashWeave [14], as microbetag also does.		471	
		472	
		473	
		474	
		475	
		476	
		477	
microbetag CytoscapeApp		478	
Overall comment, the CytoscapeApp returns averages and s.d. for example in seed scores. If you want the exact values, go through the API.		479	
		480	
A. GTDB-tk: 480 bins		481	
Step	Time(sec)	Notes	
Taxonomy mapping	Cell 1,2	on the fly	
Network inference	Cell 2,2	on the fly	
microbetag annotations	Cell 3,2	on the fly	
manta clustering	Cell 4,2	on the fly	
B. GTDB 16S: 3000 ASVs		482	
Step	Time(sec)	Notes	
Taxonomy assignment		Docker image on HP ⁶	483
Taxonomy mapping	Cell 1,2	Cell 1,3	484
Network inference	Cell 2,2	Cell 2,3	485
microbetag annotations	Cell 3,2	Cell 3,3	486
manta clustering	Cell 4,2	Cell 4,3	487
			488
C. Silva:			489
Step	Time(sec)	Notes	
Taxonomy mapping	Cell 1,2	Cell 1,3	
Network inference	Cell 2,2	Cell 2,3	
microbetag annotations	Cell 3,2	Cell 3,3	
manta clustering	Cell 4,2	Cell 4,3	
D. fuzzywuzzy:		490	
Step	Time(sec)	Notes	
Taxonomy mapping	Cell 1,2	Cell 1,3	491
Network inference	Cell 2,2	Cell 2,3	492
microbetag annotations	Cell 3,2	Cell 3,3	493
manta clustering	Cell 4,2	Cell 4,3	494
			495
Table 2: Computing times per step using an abundance table of 400 taxa with taxonomy: A. taxonomy scheme B. C. D. ⁶ specs of the laptop used.		496	
		497	
		498	
The app was based on the StringApp and supported by the NRNB group.		499	
		500	
Validation of microbetag potential		501	
vitamin dataset [38]		502	
Metagenomic or metabarcoding data are often used to predict microbial interactions in complex communities, but these predictions are rarely explored experimentally. Here, we use an organism abundance correlation network to investigate factors		503	
		504	
		505	
		506	

Study
those
2 to
under-
stand
our
find-
ings

507 that control community organization in mine tailings-derived laboratory microbial
508 consortia grown under dozens of conditions.

509 The network is overlaid with metagenomic information about functional capacities
510 to generate testable hypotheses.

511 Thiamine alternative pathway [62, 63]

512

513 Discussion ⁷

514

515 Interpreting a real-world network with **microbetag**

516 Annelies' dataset.

517

518 **microbetag** as a resource

519 Limitations

520

521 As shown in [64] (see Figure 6b), the original version of CheckM [65] that is still used on
522 GTDB returns lower completeness scores to genomes that correspond to phyla known
523 for having shorter genomes in general, e.g. Patescibacteria representative genomes on
524 GTDB have an average completeness 65%. **microbetag** inherits this in the filtering
525 process for getting only high quality genomes and thus, only few representatives from
526 these taxonomic groups are present on microbetagDB.

527 It is well known that higher-order interactions, i.e. interactions involving more
528 than two species [29]. Pairwise relationships do not capture more complex forms of
529 ecological interactions, in which one species depends on (or is influenced by) multiple
530 other species. [3]

531

532 Future work

533

534 Further indices using the seed concept have been also presented such as the metabolic
535 interaction potential (*MIP*) and the metabolic resource overlap (*MRO*). *MIP* is
536 defined as the difference between the minimal number of components required for the
537 growth of all members in a noninteracting community and an interacting community,
538 i.e. the maximum number of essential nutritional components that a community can
539 provide for itself through interspecies metabolic exchanges [29]. Similarly, *MRO* is
540 defined as the maximum possible overlap between the minimal nutritional require-
541 ments of all member species [29]. Regression and association rule mining [66] can be
542 applied to address this challenge.

- 543
- 544 • pathway and seed complementarities for higher-order interactions
 - 545 • spatial dimension
 - 546 • transcriptomics data integration: compare potential complementarities with what
 - 547 is going on
 - 548 •

549

550 ⁷The user interface should be described and a discussion of the intended uses of the software, and the
551 benefits that are envisioned, should be included, together with data on how its performance and functionality
552 compare with, and improve, on functionally similar existing software. A case study of the use of the software
may be presented. The planned future development of new features, if any, should be mentioned.

Conclusions	553
8	554
Data integration	555
Supplementary information.	556
9	557
	558
Declarations	559
• Availability of data and materials	560
– Raw sequences for the use case:	561
– Raw data for the validations case:	562
• Funding	563
This work was initiated thanks to an EMBO Scientific Exchange Grant to HZ. It was then supported by the 3D'omics Horizon project (101000309). We would also like to thank the National Resource for Network Biology (NRNB) and the Google Summer of Code 2023 for the support of E.I.M.D.	564
• Conflict of interest/Competing interests	565
The authors declare that they have no other competing interests.	566
• Authors' contributions ¹⁰	567
Conceptualization: K.F. Methodology: K.F. and H.Z. Software: H.Z., E.I.M.D. and J.M Validation: H.Z. and K.F. Formal analysis: H.Z. and K.F. Investigation: H.Z. Resources: K.F., A.E. and A.G. Data Curation: H.Z. Writing - Original Draft: H.Z. and K.F. Writing - Review & Editing: all Visualization: H.Z. Supervision: K.F., H.Z. and S.M. Project administration: K.F. Funding acquisition: K.F., H.Z.	568
• Acknowledgements	569
We would like to thank Dr Christina Pavloudi and ++ for the insight on how to organise the trait groups.	570
• Ethics approval	571
Not applicable	572
• Consent to participate	573
Not applicable.	574
• Code availability:	575
– microbetagDB related scripts: https://github.com/hariszaf/microbetag	576
– microbetagApp and webserver: https://github.com/msysbio/microbetagApp .	577
– CytoscapeApp: https://github.com/ermismd/MGG/	578
– Validation and use case: {think of having that under the 3D'omics organization}	579
– Documentation web-site: https://hariszaf.github.io/microbetag/	580

⁸This should state clearly the main conclusions and provide an explanation of the importance and relevance of the case, data, opinion, database or software reported.

⁹If your article has accompanying supplementary file(s) please state so here. E.g. supplementary figures and tables captions.

¹⁰Based on the [CRediT system](#). Current list is indicative.

599 **Appendix A Mappings**

600
601 $n : 1 n : n$ etc

602
603 **Appendix B Background on seed scores and**
604 **complementarities**
605

606 **B.1 Background on seed scores**

607
608 In that case, once a seed is assured, it activates all the rest of that group. Therefore,
609 a confidence level (C) ranging from 0 to 1, has been previously described to quantify
610 the relevance of each seed:

611
612
$$C_i = 1 / \text{seed}'s \text{ group with } i \text{ size} \quad (\text{B1})$$

613 $C = 0$ corresponds to a non-seed node, while $C = 1$ represents an independent
614 node.

615
616
$$MI_{\text{Complementarity}} = \frac{|\text{SeedSet}_A \cap \neg \text{SeedSet}_B|}{|\text{SeedSet}_A \cap (\text{SeedSet}_B \cup \neg \text{SeedSet}_B)|} \quad (\text{B2})$$

617 As also described in [55], it is calculated as the proportion of compounds in a
618 species' seed set that coincide with those in an other's, while also factoring in the
619 confidence scores associated with seed compounds.

620
621
$$MI_{\text{Competition}} = \frac{\sum C(\text{SeedSet}_A \cap \text{SeedSet}_B)}{\sum C(\text{SeedSet}_A)} \quad (\text{B3})$$

622 **B.2 Background on pathway complementarity**

623 For example, the definition of the D-Galacturonate degradation in Bacteria ([M00631](#))
624 is:

625 K01812 K00041 (K01685,K16849+K16850) K00874 (K01625,K17463)
626 that once breaking down, it leads to 4 alternative sets of KOs (pathways):
627
628 K01812 K00041 K01685 K00874 K01625
629 K01812 K00041 K16849+K16850 K00874 K01625
630 K01812 K00041 K01685 K00874 K17463
631 K01812 K00041 K16849+K16850 K00874 K17463
632
633
634
635
636
637
638
639

640 **B.3 Complementarities**

641 KEGG compound ModelSEED compounds ModelSEED compounds mapped to
642 KEGG compounds and kept only those related to KEGG modules.
643

References	645
[1] Yuan, M.M., Guo, X., Wu, L., Zhang, Y., Xiao, N., Ning, D., Shi, Z., Zhou, X., Wu, L., Yang, Y., <i>et al.</i> : Climate warming enhances microbial network complexity and stability. <i>Nature Climate Change</i> 11 (4), 343–348 (2021)	646
[2] Raes, J., Bork, P.: Molecular eco-systems biology: towards an understanding of community function. <i>Nature Reviews Microbiology</i> 6 (9), 693–699 (2008)	647
[3] Faust, K., Raes, J.: Microbial interactions: from networks to models. <i>Nature Reviews Microbiology</i> 10 (8), 538–550 (2012)	648
[4] Röttjers, L., Faust, K.: From hairballs to hypotheses—biological insights from microbial networks. <i>FEMS microbiology reviews</i> 42 (6), 761–780 (2018)	649
[5] Bálint, M., Bahram, M., Eren, A.M., Faust, K., Fuhrman, J.A., Lindahl, B., O’Hara, R.B., Öpik, M., Sogin, M.L., Unterseher, M., <i>et al.</i> : Millions of reads, thousands of taxa: microbial community structure and associations analyzed via marker genes. <i>FEMS microbiology reviews</i> 40 (5), 686–700 (2016)	650
[6] Robinson, C.J., Bohannan, B.J., Young, V.B.: From structure to function: the ecology of host-associated microbial communities. <i>Microbiology and Molecular Biology Reviews</i> 74 (3), 453–476 (2010)	651
[7] Louca, S., Jacques, S.M., Pires, A.P., Leal, J.S., Srivastava, D.S., Parfrey, L.W., Farjalla, V.F., Doebeli, M.: High taxonomic variability despite stable functional structure across microbial communities. <i>Nature ecology & evolution</i> 1 (1), 0015 (2016)	652
[8] Finn, R., Balech, B., Burgin, J., Chua, P., Corre, E., Cox, C., Donati, C., Santos, V., Fosso, B., Hancock, J., Heil, K., Ishaque, N., Kale, V., Kunath, B., Médigue, C., Paflis, E., Pesole, G., Richardson, L., Santamaria, M., Van Den Bossche, T., Vizcaíno, J., Zafeiropoulos, H., Willassen, N., Pelletier, E., Batut, B.: Establishing the elixir microbiome community [version 1; peer review: awaiting peer review]. <i>F1000Research</i> 13 (50) (2024) https://doi.org/10.12688/f1000research.144515.1	653
[9] Hug, L.A., Baker, B.J., Anantharaman, K., Brown, C.T., Probst, A.J., Castelle, C.J., Butterfield, C.N., Hernsdorf, A.W., Amano, Y., Ise, K., <i>et al.</i> : A new view of the tree of life. <i>Nature microbiology</i> 1 (5), 1–6 (2016)	654
[10] Matchado, M.S., Lauber, M., Reitmeier, S., Kacprowski, T., Baumbach, J., Haller, D., List, M.: Network analysis methods for studying microbial communities: A mini review. <i>Computational and structural biotechnology journal</i> 19 , 2687–2698 (2021)	655
[11] Faust, K., Sathirapongsasuti, J.F., Izard, J., Segata, N., Gevers, D., Raes, J., Huttenhower, C.: Microbial co-occurrence relationships in the human microbiome.	656
	657
	658
	659
	660
	661
	662
	663
	664
	665
	666
	667
	668
	669
	670
	671
	672
	673
	674
	675
	676
	677
	678
	679
	680
	681
	682
	683
	684
	685
	686
	687
	688
	689
	690

- 691 PLoS computational biology **8**(7), 1002606 (2012)
- 692
- 693 [12] Friedman, J., Alm, E.J.: Inferring correlation networks from genomic survey data.
- 694 PLoS computational biology **8**(9), 1002687 (2012)
- 695
- 696 [13] Kurtz, Z.D., Müller, C.L., Miraldi, E.R., Littman, D.R., Blaser, M.J., Bon-
- 697 neau, R.A.: Sparse and compositionally robust inference of microbial ecological
- 698 networks. PLoS computational biology **11**(5), 1004226 (2015)
- 699
- 700 [14] Tackmann, J., Rodrigues, J.F.M., Mering, C.: Rapid inference of direct interac-
- 701 tions in large-scale ecological networks from heterogeneous microbial sequencing
- 702 data. Cell systems **9**(3), 286–296 (2019)
- 703
- 704 [15] Faust, K.: Open challenges for microbial network construction and analysis. The
- 705 ISME Journal **15**(11), 3111–3118 (2021)
- 706
- 707 [16] Kishore, D., Birzu, G., Hu, Z., DeLisi, C., Korolev, K.S., Segrè, D.: Inferring
- 708 microbial co-occurrence networks from amplicon data: a systematic evaluation.
- 709 Msystems, 00961–22 (2023)
- 710
- 711 [17] Weiss, S., Van Treuren, W., Lozupone, C., Faust, K., Friedman, J., Deng, Y.,
- 712 Xia, L.C., Xu, Z.Z., Ursell, L., Alm, E.J., *et al.*: Correlation detection strategies
- 713 in microbial data sets vary widely in sensitivity and precision. The ISME journal
- 714 **10**(7), 1669–1681 (2016)
- 715
- 716 [18] Berry, D., Widder, S.: Deciphering microbial interactions and detecting keystone
- 717 species with co-occurrence networks. Frontiers in microbiology **5**, 219 (2014)
- 718
- 719 [19] Guo, B., Zhang, L., Sun, H., Gao, M., Yu, N., Zhang, Q., Mou, A., Liu, Y.: Micro-
- 720 bial co-occurrence network topological properties link with reactor parameters
- 721 and reveal importance of low-abundance genera. npj Biofilms and Microbiomes
- 722 **8**(1), 3 (2022)
- 723
- 724 [20] Ma, B., Wang, Y., Ye, S., Liu, S., Stirling, E., Gilbert, J.A., Faust, K., Knight, R.,
- 725 Jansson, J.K., Cardona, C., *et al.*: Earth microbial co-occurrence network reveals
- 726 interconnection pattern across microbiomes. Microbiome **8**, 1–12 (2020)
- 727
- 728 [21] Guidi, L., Chaffron, S., Bittner, L., Eveillard, D., Larhlimi, A., Roux, S., Darzi,
- 729 Y., Audic, S., Berline, L., Brum, J.R., *et al.*: Plankton networks driving carbon
- 730 export in the oligotrophic ocean. Nature **532**(7600), 465–470 (2016)
- 731
- 732 [22] Röttjers, L., Faust, K.: Manta: A clustering algorithm for weighted ecological
- 733 networks. Msystems **5**(1), 10–1128 (2020)
- 734
- 735 [23] Levy, R., Borenstein, E.: Reverse ecology: from systems to environments and
- 736 back. In: Evolutionary Systems Biology, pp. 329–345. Springer, ??? (2012)
- 737
- 738 [24] Levy, R., Borenstein, E.: Metagenomic systems biology and metabolic modeling of

the human microbiome: From species composition to community assembly rules. Gut Microbes 5 (2), 265–270 (2014)	737 738 739
[25] Borenstein, E., Kupiec, M., Feldman, M.W., Ruppin, E.: Large-scale reconstruction and phylogenetic analysis of metabolic environments. Proceedings of the National Academy of Sciences 105 (38), 14482–14487 (2008)	740 741 742
[26] Parter, M., Kashtan, N., Alon, U.: Environmental variability and modularity of bacterial metabolic networks. BMC evolutionary biology 7 , 1–8 (2007)	743 744 745
[27] Levy, R., Carr, R., Kreimer, A., Freilich, S., Borenstein, E.: Netcooperate: a network-based tool for inferring host-microbe and microbe-microbe cooperation. BMC bioinformatics 16 (1), 1–6 (2015)	746 747 748 749
[28] Kreimer, A., Doron-Faigenboim, A., Borenstein, E., Freilich, S.: Netcmpt: a network-based tool for calculating the metabolic competition between bacterial species. Bioinformatics 28 (16), 2195–2197 (2012)	750 751 752 753
[29] Zelezniak, A., Andrejev, S., Ponomarova, O., Mende, D.R., Bork, P., Patil, K.R.: Metabolic dependencies drive species co-occurrence in diverse microbial communities. Proceedings of the National Academy of Sciences 112 (20), 6449–6454 (2015)	754 755 756 757 758
[30] Belcour, A., Frioux, C., Aite, M., Bretaudéau, A., Hildebrand, F., Siegel, A.: Metage2metabo, microbiota-scale metabolic complementarity for the identification of key species. Elife 9 , 61968 (2020)	759 760 761 762 763
[31] Thiele, I., Palsson, B.Ø.: A protocol for generating a high-quality genome-scale metabolic reconstruction. Nature protocols 5 (1), 93–121 (2010)	764
[32] Durot, M., Bourguignon, P.-Y., Schachter, V.: Genome-scale models of bacterial metabolism: reconstruction and applications. FEMS microbiology reviews 33 (1), 164–190 (2008)	765 766 767 768
[33] Cerk, K., Ugalde-Salas, P., Nedjad, C.G., Lecomte, M., Muller, C., Sherman, D.J., Hildebrand, F., Labarthe, S., Frioux, C.: Community-scale models of microbiomes: Articulating metabolic modelling and metagenome sequencing. Microbial Biotechnology n/a(n/a) , 14396 https://doi.org/10.1111/1751-7915.14396 https://ami-journals.onlinelibrary.wiley.com/doi/pdf/10.1111/1751-7915.14396 . e14396 MICROBIO-2023-392.R1	769 770 771 772 773 774 775
[34] Mori, M., Ponce-de-León, M., Peretó, J., Montero, F.: Metabolic complementation in bacterial communities: necessary conditions and optimality. Frontiers in Microbiology 7 , 1553 (2016)	776 777 778 779
[35] Zientz, E., Dandekar, T., Gross, R.: Metabolic interdependence of obligate intracellular bacteria and their insect hosts. Microbiology and Molecular Biology	780 781 782

- 783 Reviews **68**(4), 745–770 (2004)
- 784
- 785 [36] Kallus, Y., Miller, J.H., Libby, E.: Paradoxes in leaky microbial trade. *Nature*
786 communications **8**(1), 1361 (2017)
- 787
- 788 [37] Muto, A., Kotera, M., Tokimatsu, T., Nakagawa, Z., Goto, S., Kanehisa, M.:
789 Modular architecture of metabolic pathways revealed by conserved sequences of
790 reactions. *Journal of Chemical Information and Modeling* **53**(3), 613–622 (2013)
791 <https://doi.org/10.1021/ci3005379> <https://doi.org/10.1021/ci3005379>. PMID:
792 23384306
- 793
- 794 [38] Hessler, T., Huddy, R.J., Sachdeva, R., Lei, S., Harrison, S.T., Diamond, S.,
795 Banfield, J.F.: Vitamin interdependencies predicted by metagenomics-informed
796 network analyses and validated in microbial community microcosms. *Nature*
797 Communications **14**(1), 4768 (2023)
- 798
- 799 [39] Wattam, A.R., Davis, J.J., Assaf, R., Boisvert, S., Brettin, T., Bun, C., Conrad,
800 N., Dietrich, E.M., Disz, T., Gabbard, J.L., *et al.*: Improvements to patric, the
801 all-bacterial bioinformatics database and analysis resource center. *Nucleic acids*
802 research **45**(D1), 535–542 (2017)
- 803
- 804 [40] Parks, D.H., Chuvochina, M., Rinke, C., Mussig, A.J., Chaumeil, P.-A., Hugen-
805 holtz, P.: Gtdb: an ongoing census of bacterial and archaeal diversity through
806 a phylogenetically consistent, rank normalized and complete genome-based
807 taxonomy. *Nucleic acids research* **50**(D1), 785–794 (2022)
- 808
- 809 [41] Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies,
810 J., Glöckner, F.O.: The silva ribosomal rna gene database project: improved data
811 processing and web-based tools. *Nucleic acids research* **41**(D1), 590–596 (2012)
- 812
- 813 [42] Schoch, C.L., Ciufo, S., Domrachev, M., Hotton, C.L., Kannan, S., Khovanskaya,
814 R., Leipe, D., Mcveigh, R., O'Neill, K., Robbertse, B., *et al.*: Ncbi taxonomy:
815 a comprehensive update on curation, resources and tools. *Database* **2020**, 062
816
- 817 [43] Alishum, A.: DADA2 Formatted 16S rRNA Gene Sequences for Both Bacteria
818 & Archaea. <https://doi.org/10.5281/zenodo.6655692> . <https://doi.org/10.5281/zenodo.6655692>
- 819
- 820
- 821 [44] Murali, A., Bhargava, A., Wright, E.S.: Idtaxa: a novel approach for accurate
822 taxonomic classification of microbiome sequences. *Microbiome* **6**(1), 1–14 (2018)
- 823
- 824 [45] Sunagawa, S., Coelho, L.P., Chaffron, S., Kultima, J.R., Labadie, K., Salazar,
825 G., Djahanschiri, B., Zeller, G., Mende, D.R., Alberti, A., *et al.*: Structure and
826 function of the global ocean microbiome. *Science* **348**(6237), 1261359 (2015)
- 827
- 828 [46] Louca, S., Parfrey, L.W., Doebeli, M.: Decoupling function and taxonomy in the

- global ocean microbiome. *Science* **353**(6305), 1272–1277 (2016) 829
830
- [47] Douglas, G.M., Maffei, V.J., Zaneveld, J.R., Yurgel, S.N., Brown, J.R., Taylor, C.M., Huttenhower, C., Langille, M.G.: Picrust2 for prediction of metagenome functions. *Nature biotechnology* **38**(6), 685–688 (2020) 831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
- [48] Feldbauer, R., Schulz, F., Horn, M., Rattei, T.: Prediction of microbial phenotypes based on comparative genomics. *BMC bioinformatics* **16**(14), 1–8 (2015)
- [49] Huerta-Cepas, J., Szklarczyk, D., Heller, D., Hernández-Plaza, A., Forslund, S.K., Cook, H., Mende, D.R., Letunic, I., Rattei, T., Jensen, L.J., *et al.*: eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic acids research* **47**(D1), 309–314 (2019)
- [50] Muto, A., Kotera, M., Tokimatsu, T., Nakagawa, Z., Goto, S., Kanehisa, M.: Modular architecture of metabolic pathways revealed by conserved sequences of reactions. *Journal of chemical information and modeling* **53**(3), 613–622 (2013)
- [51] Aramaki, T., Blanc-Mathieu, R., Endo, H., Ohkubo, K., Kanehisa, M., Goto, S., Ogata, H.: Kofamkoala: Kegg ortholog assignment based on profile hmm and adaptive score threshold. *Bioinformatics* **36**(7), 2251–2252 (2020)
- [52] Kanehisa, M., Goto, S., Sato, Y., Furumichi, M., Tanabe, M.: Kegg for integration and interpretation of large-scale molecular data sets. *Nucleic acids research* **40**(D1), 109–114 (2012)
- [53] Kanehisa, M., Sato, Y.: Kegg mapper for inferring cellular functions from protein sequences. *Protein Science* **29**(1), 28–35 (2020)
- [54] Kanehisa, M., Sato, Y., Kawashima, M.: Kegg mapping tools for uncovering hidden features in biological data. *Protein Science* **31**(1), 47–53 (2022)
- [55] Lam, T.J., Stamboulian, M., Han, W., Ye, Y.: Model-based and phylogenetically adjusted quantification of metabolic interaction between microbial species. *PLoS computational biology* **16**(10), 1007951 (2020)
- [56] Henry, C.S., DeJongh, M., Best, A.A., Frybarger, P.M., Lindsay, B., Stevens, R.L.: High-throughput generation, optimization and analysis of genome-scale metabolic models. *Nature biotechnology* **28**(9), 977–982 (2010)
- [57] Blank, C.E., Cui, H., Moore, L.R., Walls, R.L.: Micro: an ontology of phenotypic and metabolic characters, assays, and culture media found in prokaryotic taxonomic descriptions. *Journal of biomedical semantics* **7**(1), 1–10 (2016)
- [58] Choudhary, K., Meng, E.C., Diaz-Mejia, J.J., Bader, G.D., Pico, A.R., Morris, J.H.: scnetviz: from single cells to networks using cytoscape. *F1000Research* **10**

- 875 (2021)
- 876
- 877 [59] Merkel, D., *et al.*: Docker: lightweight linux containers for consistent development
878 and deployment. *Linux j* **239**(2), 2 (2014)
- 879
- 880 [60] Reese, W.: Nginx: The high-performance web server and reverse proxy. *Linux J.*
881 **2008**(173) (2008)
- 882
- 883 [61] Wright, E.S.: Using decipher v2. 0 to analyze big biological sequence data in r. *R*.
884 *Journal* **8**(1) (2016)
- 885
- 886 [62] Llavero-Pasquina, M., Geisler, K., Holzer, A., Mehrshahi, P., Mendoza-Ochoa,
887 G.I., Newsad, S.A., Davey, M.P., Smith, A.G.: Thiamine metabolism genes
888 in diatoms are not regulated by thiamine despite the presence of predicted
889 riboswitches. *New Phytologist* **235**(5), 1853–1867 (2022)
- 890
- 891 [63] Romine, M.F., Rodionov, D.A., Maezato, Y., Osterman, A.L., Nelson, W.C.:
892 Underlying mechanisms for syntrophic metabolism of essential enzyme cofactors
893 in microbial communities. *The ISME journal* **11**(6), 1434–1446 (2017)
- 894
- 895 [64] Chklovski, A., Parks, D.H., Woodcroft, B.J., Tyson, G.W.: Checkm2: a rapid,
896 scalable and accurate tool for assessing microbial genome quality using machine
897 learning. *Nature Methods* **20**(8), 1203–1212 (2023)
- 898
- 899 [65] Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P., Tyson, G.W.:
900 Checkm: assessing the quality of microbial genomes recovered from isolates, single
901 cells, and metagenomes. *Genome research* **25**(7), 1043–1055 (2015)
- 902
- 903 [66] Agrawal, R., Imieliński, T., Swami, A.: Mining association rules between sets of
904 items in large databases. In: Proceedings of the 1993 ACM SIGMOD International
905 Conference on Management of Data. SIGMOD '93, pp. 207–216. Association for
906 Computing Machinery, New York, NY, USA (1993). <https://doi.org/10.1145/170035.170072> . <https://doi-org.kuleuven.e-bronnen.be/10.1145/170035.170072>
- 907
- 908
- 909
- 910
- 911
- 912
- 913
- 914
- 915
- 916
- 917
- 918
- 919
- 920