

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

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

6 ^{1*} Department of Microbiology, Immunology and Transplantation, Rega
7 Institute for Medical Research, KU Leuven, Herestraat, Leuven, 3000, ,
8 Belgium .

9 ² Institute of Food, Nutrition and Health, ETH Zurich, Street, Zurich,
10 8092, , Switzerland .

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

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

14 Contributing authors: haris.zafeiropoulos@kuleuven.be;

15 ermisioannis.michaildelopoulos@student.kuleuven.be;

16 andi.erega@hest.ethz.ch; annelies.geirnaert@hest.ethz.ch;

17 scooter@cgl.ucsf.edu;

18 **Abstract**

19 *

20 Up to 350 words.

21 The abstract must include the following separate sections:

22 **Background:** the context and purpose of the study

23 **Results:** the main findings

24 **Conclusions:** a brief summary and potential implications

*Looks like Chris Quince is [our editor](#).

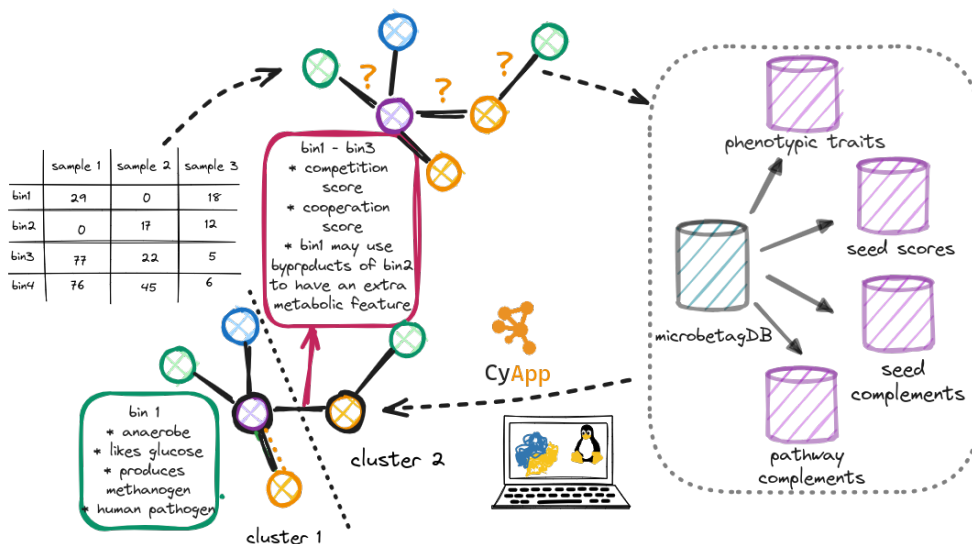


Figure abstract.

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].

Based on the net fitness effects that result for the taxa involved, the notion of an interaction varies including cooperation, competition, parasitism, commensalism and ammensalism [3]. Metabolic interactions can be established through a range of contact-independent- and contact-dependent mechanisms leading to both positive and negative interactions. These interactions can involve either one-way (unidirectional) or two-way (bidirectional) exchanges of metabolites. Depending on the biosynthetic

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

²The Background section should explain the relevant context and the specific issue that the software described is intended to address. No subheadings.

cost borne by the interacting partners, two types of metabolite exchanges occur: by-product cross-feeding, where metabolites result from a selfish act of the producer, and cooperative cross-feeding, where one partner actively invests resources to produce metabolites benefiting the interaction partner [7].

High-throughput sequencing (HTS) has provided great insight into the diversity and composition of microbial communities [8]. Uncultivated species can now be detected, and their features can be inferred through their genomic information [9]. Moreover, the composition of thousands of microbiome samples is now accessible allowing for the inference of patterns among sets of samples. A widely used approach to extract such patterns, is the creation of co-occurrence networks based on metagenomic read data (amplicon and/or shotgun) [10]. A great number of approaches is available for co-occurrence network inference based on a range of statistical concepts such as: correlation (e.g., CoNet [11], SparCC [12]), linear regression (e.g., SpiecEasi [13]) and causal inference (FlashWeave [14]). Nevertheless, microbial co-occurrence networks continue to encounter various challenges [15]. Their inference inherits the challenges of metagenomic data analysis (e.g., compositionality, parameters inference) [16]. As a result, network construction remains a tool-dependent analysis [17, 18]. Moreover, more often than not, the returned network looks like a "hairball" of densely interconnected taxa. Thus, additional analysis is necessary to generate testable hypotheses [15]. Addressing the question of *What can we learn from the hairballs* posed by Röttjers et al. [4] could provide essential insight on the mechanisms of the interactions.

The assessment of interaction predictions derived from microbial co-occurrence networks has underscored their limitations in accuracy for this task [19]. Theoretical principles derived from network studies might provide indications of emergent biological characteristics [4, 20]. For example, modules (highly interconnected nodes) within microbial co-occurrence networks could serve as indicators of ecological processes that govern community structure, including niche filtering and habitat preference [21]. Data integration and clustering have been suggested to address this challenge [15]. Clusters identified in microbial association networks have demonstrated their ability to mirror key drivers of community composition [22] and several algorithms and implementations are available [23]. However, data integration approaches in microbial co-occurrence networks are so-far limited. Here, we present **microbetag**, a microbial co-occurrence network annotator that exploits several channels of information to enhance/diminish the confidence of the associations suggested by the network and generate hypotheses for further investigation both at the taxon pair and the community level.

microbetag serves as a comprehensive platform that provides information on taxa along with their potential metabolic interactions from multiple channels (see Implementation 3). The key concept here is the reverse ecology approach *reverse ecology* [24]. Reverse ecology leverages genomics to explore community ecology with no *a priori* assumptions about the taxa involved. Making the most of the advancements in systems biology and genomic metabolic modeling, as well as system-level analysis of intricate biological networks, the reverse ecology framework enables the prediction of ecological traits for less-understood microorganisms, their interactions with others, and the overall ecology of microbial communities [25].

A metabolic network’s ”seed set” is the set of compounds that, based on the network topology, need to be acquired exogenously [26] (see Figure 2). Such nodes might be independent, i.e. they cannot be activated by any other node in the network, or they can be interdependent forming groups of seed nodes. Seeds are a useful proxy for the habitat of the organism and an essential tool in the framework of reverse ecology [26, 27]. Based on the seed concept, several graph theory-based metrics (indices) have been described to predict species interactions directly from their networks’ topologies [28–31]. Over the last years, the seed approach has been implemented at the Genome-scale metabolic network reconstructions (GENREs) level. GENREs encapsulate mathematical representations capturing the biochemical reactions that could take place within an organism [32–34].

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

However, Bacteria may complement each other not only for getting what is absolutely necessary for them to survive (seeds). For example, microbial species are recognized to exchange metabolites in order to provide support for other advantageous services, such as detoxifying harmful metabolites or offering protection against predators [35, 36]. They can additionally contribute to the production of metabolites essential for the entire community, even if the species itself does not require them [37]. To explore the potential of a species metabolism given they benefit from a partner of theirs, genome annotations combined with collections of functional units to highlight can provide a valid proxy. We present here a naive approach exporting all possible complements between a pair of species based on their KEGG ORTHOLOGY (KOs) annotations and the KEGG MODULES database [38].

microbetag annotates a user’s co-occurrence network by integrating phenotypic traits on the taxa present on the network (nodes) and potential metabolic interactions to their suggested associations (edges). A Graphical User Interface (GUI) is supported as a CytoscapeApp providing a user-friendly environment to investigate annotations in a straightforward way. All annotations present in microbetagDB are also available through an Application Programming Interface (API). **microbetag**’s source code is distributed under a GNU GPL v3 license and available on GitHub. Documentation and further support on how to use **microbetag** is available at [documentation web-site](#). To the best of our knowledge there is not a software with which **microbetag** could be compared with directly. To validate our annotations we used a recently published network with partially known interactions between some pairs of species found associated [39] (see Results section, paragraph 3). To demonstrate **microbetag**’s potential, we present the main features of its interface, and we discuss a real-world use-case (see Discussion section, paragraph 3).

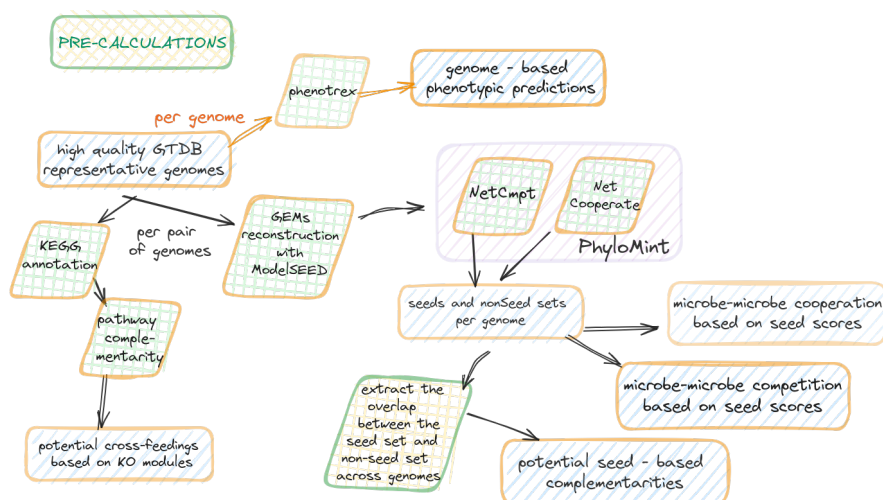


Fig. 1: Diagram of the **microbetag** pre - calculations (top panel) and the on the fly workflow (bottom panel). GTDB v207 representative genomes were filtered and for those of high-quality 33 phenotypic traits were predicted using **phenotrex**. To this end, models were re-trained to sync with recent version of eggNOG.

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 [40] 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 [41]

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.

146 that is being broadly used for bins and/or MAGs taxonomic classification and the Silva
147 database [42] that is widely used in amplicon studies. Both taxonomy schemes link
148 their taxonomies to NCBI Taxonomy Ids [43]. In case none of those two taxonomies
149 was used, and the abundance table contains less than 1,000 taxa, **microbetag** maps
150 the user provided taxonomies to NCBI Taxonomy. To this end, **microbetag** makes
151 use of the **fuzzywuzzy** library that implements the Levenshtein Distance Metric to get
152 the closest NCBI taxon name and thus its corresponding NCBI Taxonomy Id; a rela-
153 tively high similarity score is used (90) to avoid false positives. Also, using the nodes
154 dump file of NCBI Taxonomy, **microbetag** may retrieve the child taxa of a taxon in
155 user's data, along with their corresponding NCBI Taxonomy Ids, if requested by the
156 user. If the user provides their abundance table with taxonomies already mapped to
157 the GTDB taxonomy, **microbetag** will report the best possible annotations in a time
158 efficient manner.

159 Network inference

160 When a co-occurrence network is not provided by the user, **microbetag** exploits
161 FlashWeave [14] to build one on the fly. Yet, **microbetag** supports the annotation of
162 networks built from any algorithm/software, in any format Cytoscape can load.

163 **microbetag** pre-processing

164 In order to aid the user to map their sequences to the GTDB taxonomy, DADA2-
165 formatted 16S rRNA gene sequences for both bacteria and archaea [44] were used to
166 train the IDTAXA classifier of the DECIPHER package [45] and are available through
167 the **microbetag preprocess Docker image**. Likewise, when the abundance table consists
168 of more than 1,000 taxa, providing a network as an input is mandatory. Again, to help
169 the user, **microbetag** preprocess Docker image supports the inference of a network
170 using FlashWeave.

171 For a computationally efficient way to annotate large networks, a Docker image is
172 provided, so the user runs a taxonomy assignment using the IDTAXA algorithm [45]
173 of the DECIPHER R package [46]. A co-occurrence network is also built using
174 FlashWeave [14], as **microbetag** also does.

175 Literature based nodes annotation

176 Using a set of Tara Ocean samples [47] FAPROTAX [48] estimates the functional
177 potential of the bacterial and archaeal communities, by classifying each taxonomic unit
178 into functional group(s) based on current literature, descriptions of cultured represen-
179 tatives and/or manuals of systematic microbiology. In this manually curated approach,
180 a taxon is associated with a function if and only if all the cultured species within the
181 taxon have been shown to exhibit that function. In its current version, FAPROTAX
182 includes more than 80 functions based on 7600 functional annotations and covering
183 more than 4600 taxa. Contrary to gene content based approaches, e.g. PICRUSt2 [49],
184 FAPROTAX estimates metabolic phenotypes based on experimental evidence.

185 `microbetag` invokes the accompanying script of FAPROTAX and converts the
186 taxonomic microbial community profile of the samples included in the user’s abun-
187 dance table or of the taxa present in the provided network, into putative functional
188 profiles. Then, it parses FAPROTAX’s sub-tables to annotate each taxonomic unit
189 present in the user’s data with all the functions for which they had a hit. FAPROTAX
190 annotations are not part of the microbetagDB but are computed on the fly.

191 Genomic based nodes annotation

192 phenDB [50] is a publicly available resource that supports the analysis of bacterial
193 (meta)genomes to identify 47 distinct functional traits, e.g. whether a species is pro-
194 ducing butanol or has a halophilic lifestyle. It relies on support vector machines (SVM)
195 trained with manually curated datasets based on gene presence/absence patterns for
196 trait prediction. More specifically, the model for a particular trait is trained using a
197 collection of EggNOG annotated genomes where the knowledge of whether that trait
198 is present or absent among its members is available. These models (classifiers) are
199 used to predict presence/absence of their corresponding traits in non-studied species.

200 In the framework of microbetagDB, classifiers were re-trained using the genomes
201 provided by phenDB for each trait to sync with the latest version of eggNOG [51]
202 and the `phenotrex` [50] software tool. Genomes were downloaded from NCBI using
203 the `Batch Entrez` program. Then, *genotype* files were produced for all the high quality
204 GTDB representative genomes. Each model was then used against all the GTDB
205 *genotype* files to annotate each with the presence or the absence of the trait. A list of all
206 the phenotypic traits available for the genomes present in microbetagDB is available
207 on `microbetag`’s [documentation site](#). The updated models are also available

208 Pathway complementarity

209 To infer potential pathway complementarities we consider the modules described in
210 KEGG MODULES database [38]. A KEGG module is defined as a functional unit
211 within the KEGG framework that represents a set of enzymes and reactions involved
212 in a specific biological process or pathway [52]. Such a unit consists of several *steps*,
213 each of which may have more than one molecular ways to occur (Figure 2). A module’s
214 definition is a logical expression and consists of KOs that may be coupled with one
215 another as: a. connected steps of the pathway b. parts of a molecular complex, c.
216 alternatives of the same step, and d. optional entities of a complex. Both (a) and
217 (b) cases should be considered as the **AND** logical operator, while (c) would be the **OR**
218 (Figure 2). Given a module’s definition, we will consider as an *alternative* any subset
219 of the KO terms mentioned in the definition, that has exactly one way to perform
220 each step, provided that all the steps of the module are covered. We define a genome
221 as having a *complete* module, if and only if all the KOs of at least one alternative are
222 present on the genome. In Appendix A we show an example of a module along with
223 its alternatives.

224 Within this framework, `kofamscan` [53] was used to annotate with KEGG
225 ORTHOLOGY terms (KOs) the 16,900 high quality GTDB representative genomes

for which a .faa was available [54]. The KOs of each genome were then mapped to their corresponding KEGG modules; a KO may map to more than one module (1 : n).

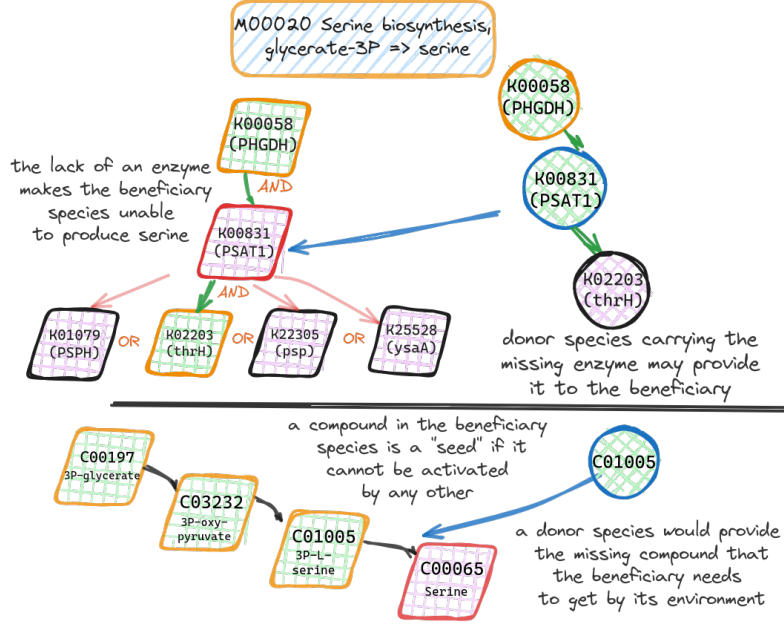


Fig. 2: Pathway complementarity approach. The high quality GTDB genomes were annotated with KEGG ORTHOLOGY (KO) terms. The various ways of getting a KEGG module complete were enumerated and all the possible ways a donor species could "fill" a beneficiary's non-complete module were calculated. In this 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 investigated for potential pathway complementarities, i.e. whether a genome lacking a number of KOs ($genome_A$) to have a complete module ($module_x$) could benefit from another's species genome(s) ($genome_B$). In that case, $genome_B$ does not necessarily have a complete alternative of $module_x$; as long as it has the missing KOs that $genome_A$ needs to complete an alternative of it, $genome_B$ potentially complements $genome_A$ with respect to $module_x$. In total, 341,568 unique complementarities were exported.

Thanks to the graphical user interface (GUI) of the [KEGG pathway map viewer](#) [55, 56], each complementarity can be visualised as part of the closest KEGG

metabolic map; where the KOs contributed by the donor are shown in blue-green whereas those coming from the beneficiary genome are coloured in red.

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

Seed scores and complements using genome scale metabolic reconstructions

The Metabolic Complementarity Index ($MI_{Complementarity}$) measures the degree to which two microbial species can mutually assist each other by complementing each other's biosynthetic capabilities. As described in [57], it is defined as the proportion of seed compounds of a species that can be synthesized by the metabolic network of another, but are not included in the seed set of the latter. $MI_{Complementarity}$ offers an upper bound assessment of the potential for syntrophic interactions between two species. Further, the Metabolic Competition Index ($MI_{Competition}$) represents the similarity in two species' nutritional profiles. This index establishes an upper limit on the level of competition that one species may face from another. Those indices have been thoroughly described and implemented in the NetCooperate [28] and Net-Compt [29] tools correspondingly. We will be referring to those two indices as "seed scores". Recently, the PhyloMint tool [57] was released supporting the calculation of the seed scores of GENREs in SBML format.

In the `microbetag` framework, seed scores were computed using GENREs derived from the high quality GTDB representative genomes and the PhyloMint tool. GENREs were reconstructed using the Model SEED pipeline [58] through its Python interface `ModelSEEDpy`. The latter requires RAST annotated genomes [59]; if available through the PATRIC database [40], annotations were retrieved. For the rest of the genomes, RAST annotation was performed through RASTtk [60].

Moreover, the computed seed and the non-seed (i.e., set of metabolic compounds a genome can build on its own) sets of each genome were used to compute their overlap among all the pairwise combinations of those genomes. More specifically, seed and non-seed compounds of each genome were mapped to their corresponding KO terms and those related to any KEGG MODULE were considered further. Focusing on the KEGG MODULE - related KO terms as terms of interest, the overlap of *seed set_{species_A}* with the *non seed set_{species_B}* was retrieved. Such *seed complementarities* were calculated for all pairwise GENREs and are now available through microbetagDB. Edges of the co-occurrence network where both taxa have been mapped to at least one GTDB genome can be further annotated mentioning all the KEGG maps for which there is at least one seed compound of the potentially beneficiary species.

281 Clustering network

282 **manta** is a heuristic network clustering algorithm that clusters nodes within weighted
283 networks effectively, leveraging the presence of negative edges and discerning between
284 weak and strong cluster assignments. **microbetag** invokes **manta** [23] to infer clusters
285 from the microbial network. In case **manta** is performed, the annotated network inherits
286 the layout that **manta** returns.

287 The microbetag workflow

288 As shown in Figure 3, the **microbetag** workflow expects an abundance table repre-
289 senting either amplicon or shotgun data. If a co-occurrence network is already available
290 the user may provide it too as input. The **microbetag** workflow will first map the taxa
291 present on the abundance table to their corresponding GTDB representative genomes
292 if that is possible, i.e., in case the taxonomy provided does reach the species or the
293 strain level (see paragraph 3). If a network is not provided, **microbetag** will then build
294 one using FlashWeave [14]. Then the abundance table will be used for a literature
295 - based annotation using FAPROTAX [48]. This is the only annotation step that is
296 **microbetagDB** independent in the framework of the web-service workflow. The nodes
297 of the network will be further annotated with phenotypic traits based on the model
298 predictions [50]. Edges linking taxa that have been assigned to the species or strain
299 level will be then annotated with pathway and seed complementarities and with seed
300 scores. Last, a network clustering will be performed assigning each node to a cluster.
301 The annotated network is then returned in a **.cx** format. The user may skip any of
302 these annotation steps if not needed for their analysis.

303 Groups of annotations

304 Biologically meaningful groups were described to group phenotypic traits returned
305 from FAPROTAX and phenDB-like annotation steps. The main groups supported are
306 related to: a. the lifestyle of a species, for example being halophilic or thermophyllic
307 etc., b. the biogeochemical processes a species metabolic potential has been found
308 related to, for example Nitrite-oxidizing bacteria (NOB) bacteria and c. important
309 metabolites a species is suggested to produce, e.g. butanol. Aim of these groups are to
310 facilitate filtering of the taxa present. Enrichment analysis for members of such groups
311 (e.g., based on the findings of a clustering algorithm like **manta**) can be performed
312 through the CytoscapeApp.

313 Software architecture

314 **microbetag** is a Docker-based application. We deployed the **microbetag** application
315 using Docker containers [61] (v24.0.2) managed by Docker Compose (see Supplemen-
316 tary Figure A). Docker Compose is a tool for defining and running multi-container
317 Docker applications using a YAML file to configure the services required for the appli-
318 cation. Containers of three Docker images are being used simultaneously: a. a **MySQL**
319 database including the **microbetagDB** b. a **nginx** [62] web server and c. the applica-
320 tion itself, including the API and the **microbetag** workflow. The latter uses **Gunicorn**

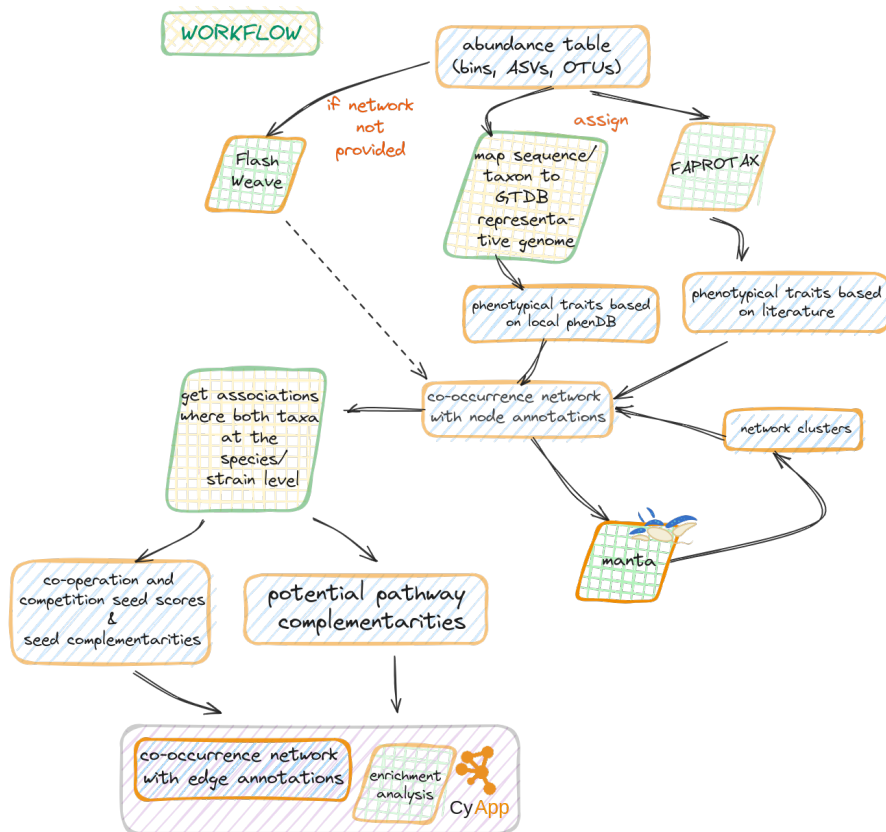


Fig. 3: Diagram of *microbetag*’s on-the-fly workflow. *microbetag* expects either an abundance table only as input and infers a co-occurrence network using FlashWeave or an abundance table along with an already inferred co-occurrence network and after mapping taxa present to GTDB reference genomes, for those possible, phenotypic attributes are assigned on the nodes. Literature-based annotation on the nodes are also using FAPROTAX. On the edges level then, *microbetag* annotates them by assigning the pre-calculated potential complements based on the pathway and the seed complementarities approaches. *microbetag* supports optional network clustering with manta. The annotated network can then be parsed on Cytoscape using the MGG app.

321 (20.1.0) to build an application server which communicates with the web server using
 322 the Web Server Gateway Interface (WSGI) protocol and handles incoming HTTP
 323 requests. *microbetag* is implemented as a [Flask](#) application (v2.3.2); Flask is a micro
 324 web framework for developing Python web applications and RESTful APIs. The API
 325 has a route for performing the *microbetag* workflow, either through any Python
 326 console or the Cytoscape MGG app, but also several other routes that enable quick
 327 and easy access to the *microbetag*DB content, i.e. the genomes present, their pheno-
 328 typic traits predicted annotations, pathway and seed complementarities among specific

329 genomes or NCBI Taxonomy Ids and their corresponding seed scores if available. A
330 thorough description of the `microbetag` API is available at the [ReadTheDocs web](#)
331 [site](#). The source code of the `microbetag` web service is available on [GitHub](#).

332 The MGG CytoscapeApp

333 We developed a Cytoscape app to enable a straightforward, user-friendly way to per-
334 form the `microbetag` workflow and visualise `microbetag` - annotated networks. The
335 `microbetag` CytoscapeApp (called MGG) was built based on the [source code](#) of the
336 `scVizNet` [63]. A visual style was developed to facilitate in distinguishing annotated
337 nodes and edges. Nodes are colored based on the level of the taxonomic assignment
338 with those being annotated highlighted with green. Similarly, edges are light green
339 when they carry a positive weight and red when negative. Black edges denote pathway
340 and/or seed complementarities. The last were not added in the weight edge as the first
341 describes an undirected relationship while the last a directed one. MGG allows the user
342 to import their data, retrieve an annotated network and investigate the annotations
343 through a series of CyPanels both for node and edge annotations. Figure 4 shows an
344 example of the CyPanels. In the nodes panel (4.A), the node name, the taxonomy as
345 well as the NCBI Taxonomy Id and the GTDB genome to which the sequence was
346 mapped can be viewed. Depending on the user's settings and the available annotations
347 for a node, genomic based predictions may be present and/or literature - based ones.
348 Further, the annotation groups mentioned in paragraph 3 are on top of this panel
349 allowing for the selection of the nodes carrying either one among several attributes
350 (OR logical relationship) or all of them (AND). Accordingly, in the edges panel (4.B),
351 the beneficiary taxon is specified along with their corresponding GTDB representa-
352 tive sequence identifier. Pathway and seed complementarities are shown each in a
353 table. Potential metabolic interactions are shown in sub-table entitled with the pair of
354 genomes under consideration, as several GTDB genomes may have been assigned to a
355 node. In case of pathway complementarities, these tables consist of six columns: a. the
356 KEGG MODULE id of the module to be completed, b. its description, c. a more gen-
357 eral metabolism category the module is related to, d. the complement itself as a list of
358 KEGG terms, e. the alternative that is now complete and allows the beneficiary to per-
359 form the module and f. a URL that points to a coloured KEGG map highlighting the
360 complement. If clicked, user's default browser pops-up showing a coloured KEGG map
361 a part of an example is shown in Figure 4.C. Last, MGG supports enrichment analy-
362 sis of the network's nodes based on the phenotypic traits assigned between clusters.
363 Clusters may have been returned from `manta` [23] while performing the `microbetag`
364 workflow or users may assign them on their own or using any other network clustering
365 algorithm. For thorough instructions on how to use MGG and `microbetag` the reader
366 may visit the [ReadTheDocs web site](#).

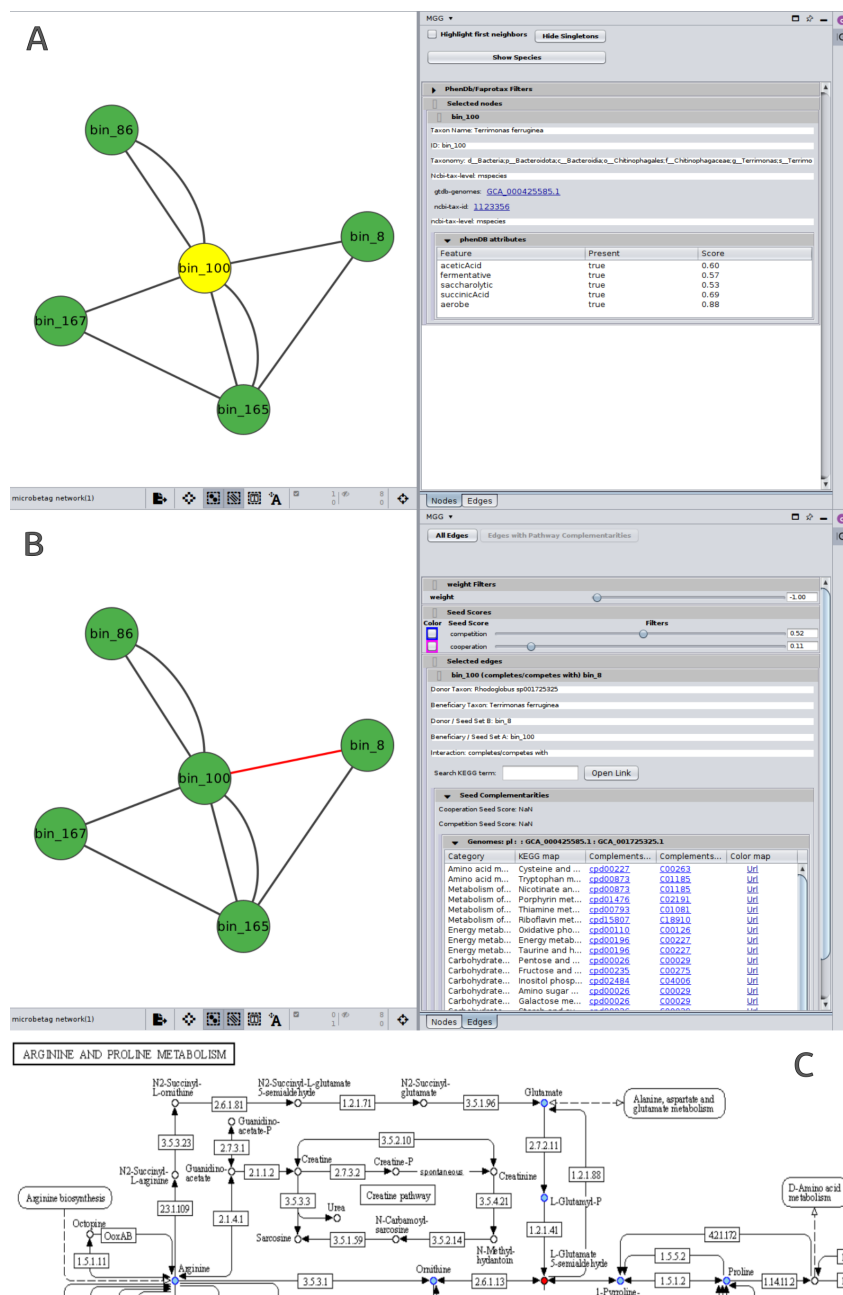


Fig. 4: CyPanels of the MGG CytoscapeApp. **A.** *Nodes* panel display the annotations of each taxon (node) mapped to one or more GTDB genomes. In this example, genomic predicted phenotypic attributes are shown along with their prediction score. **B.** *Edges* panel display the list of potential metabolic complementarities between two nodes, specifying which is the potential donor and the potential beneficiary taxon; thus giving a directed perspective on the graph. There are two cases of complementarities in the *microbetag* framework. *Seed complementarities* shown here are first exported based on ModelSEED complements (column three) and mapped in KEGG COMPOUNDS (column four). In the URL provided, a colored KEGG map is provided. The same applies for the case of the **Pathway complementarities** only there is no ModelSEED ids as they are computed directly from the KEGG annotated genomes and not from the Genome-Scale Metabolic Reconstructions; that is the case for the *Seed complementarities*. **C.** Part of a colored KEGG map returned based on the seed complementarities. Compounds that the beneficiary taxon brings on its own are colored in cyan while the potential complement with red.

367 Results and discussion ⁴

368 Annotating microbial co-occurrence networks with microbetag

369 The `microbetag` software ecosystem consists of five main modules: a. `microbetagDB`
370 including `microbetag` precalculations, b. the `microbetag` workflow to annotate
371 the co-occurrence network, c. a web server hosting both the `microbetagDB` and
372 the `microbetag` application, d. a CytoscapeApp called MGG that enables a user-
373 friendly invoke of the workflow and investigation of the annotated network, and e. a
374 stand-alone pre-processing workflow provided as Docker image for data sets with more
375 than 1,000 sequence identifiers (OTUs/ASVs/bins etc.).

376 Currently, `microbetagDB` includes genome and thus annotations for more than
377 34,000 genomes (Table 1) with the vast majority of those representing bacterial taxa
378 and 364 archaeal. Presence/absence of more than 30 phenotypic traits have been pre-
379 dicted for those genomes. About 1.4 billion potential metabolic interactions leading to
380 pathway or seed complementarities have been precomputed as well. Seed complements
381 are one order of magnitude more than those corresponding to pathway complemen-
382 tarities as for all GENREs present in `microbetagDB` all pairwise complements were
383 calculated ($33,755^2$) and stored even if empty. In case of pathway complementarities,
384 a genome pair is present in the database and thus counted only if a potential com-
385 plementarity was found. Yet, in the first case, the number of genomes with absolutely
386 no potential seed complement ranges from zero to a few dozens. A following paper on
387 the `microbetagDB` content is in preparation. All annotations can be accessed directly
388 from `microbetagDB` through the API. Using GENREs for the seed complements and
389 not the genomes per se supports a more realistic simulation of what the correspond-
390 ing taxa do need to get from the environment to grow (seeds) but also, assuming they
391 grow what they may secrete (non-seeds). However, this comes with its own challenges
392 (see paragraph 3).

393 Running the `microbetag` workflow is straightforward and can be done using a
394 taxonomically assigned abundance table as input. When the taxonomy scheme being
395 used is not one among the GTDB, Silva or the GTDB-oriented taxonomy for 16S
396 rRNA amplicon data (see pre-process paragraph 3), the most time-consuming step
397 of the workflow is the one mapping user's taxonomy to a NCBI Taxonomy Id and
398 from that, to GTDB representative genomes. Network inference can be a computa-
399 tionally intensive step too, particularly as the number of sequences in the abundance
400 table increases. To enable annotation of large data sets, a stand-alone pre-process
401 workflow is provided with `microbetag`. The user can either assign their ampli-
402 con data to the GTDB-oriented taxonomy and/or reconstruct a network locally.
403 Once a network is available and the taxonomy used is among the standard ones

⁴**Results-related:** Significant 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. **Discussion-related:** The user interface should be described and a discussion of the intended uses of the software, and the benefits that are envisioned, should be included, together with data on how its performance and functionality 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.

for **microbetag**, the computational time required for annotation ranges from several seconds to a few minutes based on user’s settings. An annotated network in **cx2** format is returned which can be viewed on Cytoscape. A tutorial, frequently asked questions and hints to address the idiosyncrasy of various data sets are available on the <https://hariszaf.github.io/microbetag/ReadTheDocs> web site while a Gitter community allows users to exchange experience and ask for more specific help. In the following two sections we present a validation and a use case, highlighting the potential of our approach.

Table 1: Summary of the data in microbetagDB

Description	Entries
GTDB representative genomes	34,608
Phen-model-oriented metabolic functions	32
FAPROTAX functions	92
Unique pathway complements	341,568
Pairwise pathway complementarities	184,184,548
GENREs leading	33,755
Seed complements	1,139,400,025
Seed scores	1,105,250,048

Validation of microbetag potential

To validate **microbetag** we used the correlation network of Hessler et al. [39] describing mine tailing-derived laboratory microbial consortia. In this study, *Variovorax*, a thiamine producer, and its co-occurrence with a series of thiamine auxotrophs are discussed. The study was selected as a validation case as the authors tested network’s predictions by performing co-culture experiments measuring the thiamine production. Both bins sequences corresponding to network’s nodes and the original network were retrieved. Using GTDB-tk [64] bins were annotated to GTDB taxonomies. Taxonomies retrieved for each bin were added in the original network which was then annotated with **microbetag**. Figure ?? highlights bin_55 that corresponds to *Variovorax* and its first neighbors. The annotated network is available on **microbetag**’s [GitHub repository](#). GTDB-tk returned GCA_001899795.1 as the one closer to bin_55 assigning it as *Variovorax* sp001899795. **microbetag** then suggested that this specific genome corresponds to an aerobe [65], autotroph if needed *autotrophVariovorax* that utilizes D-glucose, while producing ethanol and lactic acid [66]. Last, Type VI secretion system was suggested to be available on its genome [67]. As shown in Table 2, **microbetag** suggested several thiamine-related potential seed complements between *Variovorax* and their first neighbors on the network (Table 2.A). Further, **microbetag** also suggested potential thiamine-related complements among the neighboring taxa (Table 2.B).

supplementary file ??

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

Thiamine alternative pathway [68, 69]

A. *Variovorax* thiamine-related benefits to its neighbors

Neighboring taxon	node id	KEGG compounds	url
<i>Kapabacteria thiocyanatum</i>	bin_59	C15809	url
<i>Terrimonas ferruginea</i>	bin_100	C15809;C01081	url
<i>Tahibacter</i> sp001725155	bin_167	C15809	url
<i>Microbacterium</i> sp900156455	bin_28	C15809; C20246	url
<i>Sphingobium</i> sp001899715	bin_155	Iminoglycine C15809;	url
<i>Nitrosospora</i> sp001899235	bin_176	None	None
62-47 sp001899255*	bin_233	None	None
<i>Bosea</i> sp001898115	bin_273	C04327;C01279	url
54-19 sp001898225**	bin_41	C15809	url
<i>Rhodoglobus</i> sp001725325	bin_8	C15809	url

B. Potential thiamine-related complements among *Variovorax* neighbors

Beneficiary	Donor	potential complement
<i>T. ferruginea</i>	<i>Tahibacter</i> sp001725155	C01081
<i>T. ferruginea</i>	<i>Rhodoglobus</i> sp001725325	C01081
<i>Nitrosospora</i> sp001899235	<i>Bosea</i> sp001898115	C04327;C01279
Chloroflexi	<i>Bosea</i> sp001898115	C15809
Chloroflexi	Xanthobacteraceae	C15809
Chloroflexi	<i>Nitrosospora</i> sp001899235	C15809

Table 2: Thiamine biosynthesis related seed complements between *Variovorax* and its first closest neighbors on the network of Hessler *et al.* [39] (A), and between pairs of the neighbors (B). Bin sequence files were mapped against GTDB using GTDB-tk. Chloroflexi refers to the GTDB taxonomy of: 54-19 sp001898225

Regarding pantothenate, the *Variovorax* genome mapped from GTDB-tk brings two complete KEGG modules for that, ([M00119](#): Pantothenate biosynthesis, valine/L-aspartate \Rightarrow pantothenate and [M00913](#): Pantothenate biosynthesis, 2-oxoisovalerate/spermine \Rightarrow pantothenate)

thus either some *Variovorax* species can actually produce that or this is a limitation of our method ??

In fact, *Variovorax* seems to be able to benefit some partners of their.

From the hub species:

Interpreting a real-world network with microbetag

Annelies’ dataset.

One last visual component from the use case would be nice to have.

Potential and limitations

The previous paragraph shows the potential the **microbetag** workflow may have in the interpretation of co-occurrence network and how it can be used to generate new hypotheses derived from these. However, **microbetag** benefits the microbiome community in several other ways. The **microbetagDB** provides a vast number of annotations; 31 predicted traits for more than 30,000 genomes, their GENREs along with their corresponding seed sets, potential metabolic complementarities and cooperation/competition scores. Such a resource may benefit a range of studies; from a more theoretical perspective regarding the distribution of the complements among taxonomic groups or how often a complement potentially appears, to more applicable such as eco-evolutionary studies and the investigation of established interactions.

Yet, there is a number of challenges in our approach. First, **microbetag** inherits all the biases and drawbacks of both the data and the software it is based on. For example, regarding the genomes currently supported, and as shown in [70] (see Figure 6b), the original version of CheckM [71] that is still used on GTDB returns lower completeness scores to genomes that correspond to phyla known for having shorter genomes in general, e.g. *Patescibacteria* representative genomes on GTDB have an

average completeness $\sim 65\%$. Thus, only few representatives from these taxonomic groups are present on microbetagDB leading to an important under-representation of Archaea. Functional annotation comes with its own limitations. Some domains boast richer annotations and more comprehensive descriptions compared to others. These areas exhibit a wealth of detail and employ more precise terminology, particularly for widely recognized processes. In our case, pathway complementarity can be as accurate as the KEGG MODULE database goes and the precision of the software annotating genomes with KO terms. It is well known that automated Genome-Scale Metabolic Reconstruction comes with a great number of challenges and different software for this task come with their certain limitations [72]. Using ModelSEED with a complete medium may limit potential metabolic interactions but made those retrieved of higher confidence.

It is also well known that higher-order interactions, i.e. interactions involving more than two species [30] Pairwise relationships do not capture more complex forms of ecological interactions, in which one species depends on (or is influenced by) multiple other species [3]. Further,

Future work

In the near future, we plan to develop two main features: a. the integration of transcriptomics data provided by the user, this would enhance or lower the probability for a potential metabolic interaction to occur based on whether the KO terms involved are present or not, and b. the integration of spatial data; it is well-known that the spatial dimension plays a great role to the extent that an interaction occurs [73], to this end we intend to integrate user's data on how their data are distributed in space. Thus, potential metabolic interactions between taxa that are closer one-another would be more probable to occur.

Last, we already work on a "*for advanced users*" version, a server-independent version of `microbetag` is about to be released, so the user can provide bins/MAGs of theirs and annotations will be held not by mapping taxonomies to reference genomes but using their sequencing data directly. This would require important computing resources and time and cannot be supported in an app-framework like the one presented here. In this case, one will be again able to investigate the annotated network returned through Cytoscape and the MGG app. ⁵

Conclusions ⁶

Co-occurrence networks are widely used in microbiome studies to infer associations [4]. Both their inference and their interpretation though come with a range of challenges [15]. Metabolic exchanges among microbial taxa is considered ubiquitous [74] at least in a great range of environments. In our study, we exploit reverse-ecology approaches and publicly available genomic data and software to predict phenotypic traits and metabolic interactions and annotate with those co-occurrence networks

⁵Could be part of this release; time will tell

⁶This should state clearly the main conclusions and provide an explanation of the importance and relevance of the case, data, opinion, database or software reported.

502 derived from amplicon or shotgun data. Our annotation was in-line with the study
 503 of Hessler et al. [citehessler2023vitamin](#) predicting thiamine-related metabolic interac-
 504 tions among *Variovorax* and its closest neighbors, suggesting several ways to achieve
 505 them. Using [Enrichment analysis using them combined with network cluster-](#)
 506 ing algorithms can further benefit their interpretation. Both the microbetagDB and
 507 **microbetag** workflow may benefit microbiome studies, both as a resource and as a
 508 hypothesis generation tool.

use
case

509 **Supplementary information.** List of supplementary figures and tables.

510 **Supplementary Figure 1:** **microbetag** software ecosystem architecture.

511 **Supplementary Table 1:** *Variovorax* genomes present on microbetagDB and
 512 their corresponding complete/incomplete presence of the pantothenate - related
 513 KEGG modules

514 **Supplementary Table 2:** Computing times per step of the **microbetag** workflow
 515 using four different data sets.

516 Declarations

517 • Availability of data and materials

- 518 – Raw sequences for the use case:
- 519 – Raw data for the validations case:

520 • Funding

521 This work was initiated thanks to an EMBO Scientific Exchange Grant to HZ. It
 522 was then supported by the 3D'omics Horizon project (101000309). We would also
 523 like to thank the National Resource for Network Biology (NRNB) and the Google
 524 Summer of Code 2023 for the support of E.I.M.D.

525 • Conflict of interest/Competing interests

526 The authors declare that they have no other competing interests.

527 • Authors' contributions ⁷

528 Conceptualization: K.F. Methodology: K.F. and H.Z. Software: H.Z., E.I.M.D. and
 529 J.M. Validation: H.Z. and K.F. Formal analysis: H.Z. and K.F. Investigation: H.Z.
 530 Resources: K.F., A.E. and A.G. Data Curation: H.Z. Writing - Original Draft: H.Z.
 531 and K.F. Writing - Review & Editing: all Visualization: H.Z. Supervision: K.F.,
 532 H.Z. and S.M. Project administration: K.F. Funding acquisition: K.F., H.Z.

533 • Acknowledgements

534 We would like to thank Dr. Christina Pavloudi for the insight on how to organise the
 535 trait groups. We would also like to thank Dr. Hessler and Prof. Jillian F. Banfield
 536 for sharing both the bins and the network of their study [39].

537 • Ethics approval

538 Not applicable

539 • Consent to participate

540 Not applicable.

541 • Code availability:

⁷Based on the [CRediT](#) system.

- 542 – microbetagDB related scripts: <https://github.com/hariszaf/microbetag>
- 543 – **microbetag** application: <https://github.com/msysbio/microbetagApp>.
- 544 – MGG CytoscapeApp: <https://github.com/ermismd/MGG/>
- 545 – Validation and use case: ⁸
- 546 – Documentation web-site: <https://hariszaf.github.io/microbetag/>

547 Appendix A

548 Background on pathway and seed complementarity

549 For a genome to have a KEGG module *complete* means it affords at least one com-
 550 plete *alternative*. Based on the module’s definition, alternatives are considered as the
 551 unique combinations of KOs that will enable the module. For example, the definition
 552 of the D-Galacturonate degradation in Bacteria ([M00631](#)) is:

553 K01812 K00041 (K01685,K16849+K16850) K00874 (K01625,K17463)

554 Once breaking down, it leads to 4 alternative sets of KOs (pathways):

555 K01812 K00041 K01685 K00874 K01625
 556 K01812 K00041 K16849+K16850 K00874 K01625
 557 K01812 K00041 K01685 K00874 K17463
 558 K01812 K00041 K16849+K16850 K00874 K17463

560 In alternatives two and four, the K16849+K16850 is a *complex*, meaning both KO
 561 terms are required for the step to be available.

562 In case of seed complementarity, in the framework of **microbetag** we focus on
 563 the effect that a metabolic exchange between two taxa might have if the seed of the
 564 beneficiary taxon is related to a KEGG MODULE. Therefore, the KOs that were found
 565 related to modules were mapped to ModelSEED ids. The initial seed and non-seed sets
 566 that were exported as sets of ModelSEED ids were then mapped to KOs too. When
 567 the non-seed set of a genome (donor) was able to provide a seed that is related to a
 568 KEGG module to another genome (beneficiary) was considered as potential metabolic
 569 interaction.

Table A1: *Variovorax* genomes present on microbetagDB and their corresponding complete/incomplete presence of the pantothenate-related KEGG modules.

Genome	md:M00119	md:M00913
GCA_004210915.1	incomplete	complete
GCA_902506565.1	incomplete	incomplete
GCF_000184745.1	complete	complete
GCF_000282635.1	complete	complete
GCF_000463015.1	complete	complete
GCF_000834655.1	complete	complete
GCF_001424835.1	complete	complete
GCF_001425205.1	complete	complete
GCF_001426505.1	complete	complete
GCF_001577265.1	incomplete	incomplete
GCF_002157355.1	complete	complete
GCF_002754375.1	complete	complete
GCF_003019815.1	incomplete	complete
GCF_003852515.1	complete	complete
GCF_003951285.1	complete	complete
GCF_003952165.1	complete	complete
GCF_003952185.1	complete	complete
GCF_003984625.1	complete	complete
GCF_003984645.1	complete	complete
GCF_006438845.1	complete	complete
GCF_007828835.1	complete	complete
GCF_009498455.1	complete	complete
GCF_009755665.1	complete	complete
GCF_010499245.1	complete	complete
GCF_013376045.1	complete	complete
GCF_014170375.1	complete	complete
GCF_014302995.1	complete	complete
GCF_014303735.1	incomplete	incomplete
GCF_901827175.1	complete	complete
GCF_901827205.1	complete	complete

570

571

Validation

software development

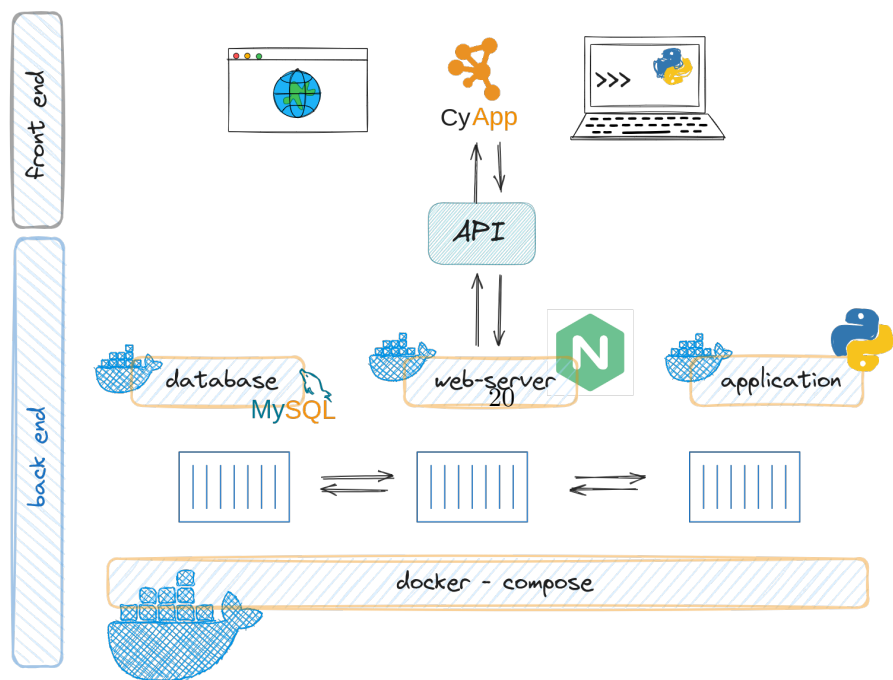


Fig. A1: microbetag software ecosystem.

⁸Consider moving that under the 3D'omics organization

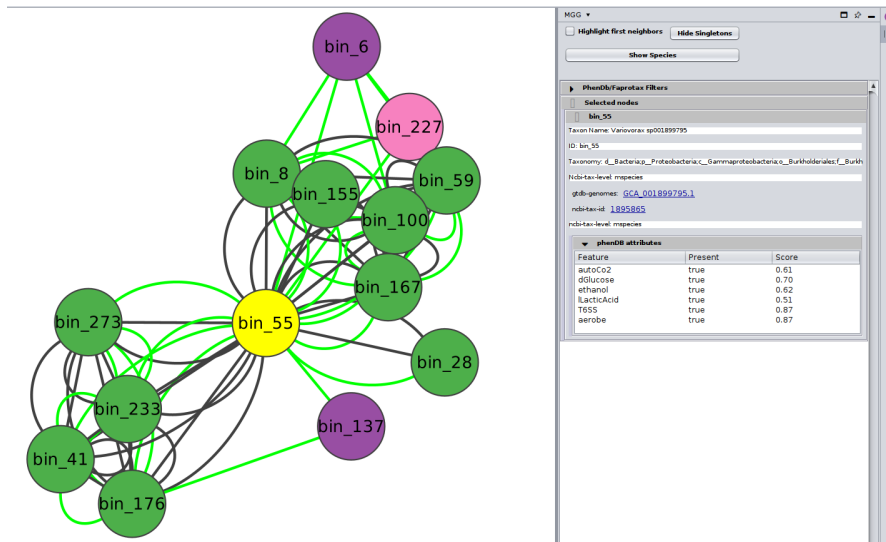


Fig. A2: *Variovorax* node and its neighbors **microbetag** annotated. All but three of them were not mapped to a GTDB representative genome. All edges

References

- [1] Yuan, M.M., Guo, X., Wu, L., Zhang, Y., Xiao, N., Ning, D., Shi, Z., Zhou, X., Wu, L., Yang, Y., *et al.*: Climate warming enhances microbial network complexity and stability. *Nature Climate Change* **11**(4), 343–348 (2021)
- [2] Raes, J., Bork, P.: Molecular eco-systems biology: towards an understanding of community function. *Nature Reviews Microbiology* **6**(9), 693–699 (2008)
- [3] Faust, K., Raes, J.: Microbial interactions: from networks to models. *Nature Reviews Microbiology* **10**(8), 538–550 (2012)
- [4] Röttjers, L., Faust, K.: From hairballs to hypotheses—biological insights from microbial networks. *FEMS microbiology reviews* **42**(6), 761–780 (2018)
- [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., *et al.*: Millions of reads, thousands of taxa: microbial community structure and associations analyzed via marker genes. *FEMS microbiology reviews* **40**(5), 686–700 (2016)
- [6] Robinson, C.J., Bohannan, B.J., Young, V.B.: From structure to function: the ecology of host-associated microbial communities. *Microbiology and Molecular Biology Reviews* **74**(3), 453–476 (2010)
- [7] D’Souza, G., Shitut, S., Preussger, D., Yousif, G., Waschina, S., Kost, C.: Ecology and evolution of metabolic cross-feeding interactions in bacteria. *Natural Product Reports* **35**(5), 455–488 (2018)
- [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., Pafilis, 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]. *F1000Research* **13**(50) (2024) <https://doi.org/10.12688/f1000research.144515.1>
- [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., *et al.*: A new view of the tree of life. *Nature microbiology* **1**(5), 1–6 (2016)
- [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. *Computational and structural biotechnology journal* **19**, 2687–2698 (2021)
- [11] Faust, K., Sathirapongsasuti, J.F., Izard, J., Segata, N., Gevers, D., Raes, J., Huttenhower, C.: Microbial co-occurrence relationships in the human microbiome. *PLoS computational biology* **8**(7), 1002606 (2012)

- [12] Friedman, J., Alm, E.J.: Inferring correlation networks from genomic survey data. *PLoS computational biology* **8**(9), 1002687 (2012)
- [13] Kurtz, Z.D., Müller, C.L., Miraldi, E.R., Littman, D.R., Blaser, M.J., Bonneau, R.A.: Sparse and compositionally robust inference of microbial ecological networks. *PLoS computational biology* **11**(5), 1004226 (2015)
- [14] Tackmann, J., Rodrigues, J.F.M., Mering, C.: Rapid inference of direct interactions in large-scale ecological networks from heterogeneous microbial sequencing data. *Cell systems* **9**(3), 286–296 (2019)
- [15] Faust, K.: Open challenges for microbial network construction and analysis. *The ISME Journal* **15**(11), 3111–3118 (2021)
- [16] Cao, H.-T., Gibson, T.E., Bashan, A., Liu, Y.-Y.: Inferring human microbial dynamics from temporal metagenomics data: Pitfalls and lessons. *BioEssays* **39**(2), 1600188 (2017)
- [17] Kishore, D., Birzu, G., Hu, Z., DeLisi, C., Korolev, K.S., Segrè, D.: Inferring microbial co-occurrence networks from amplicon data: a systematic evaluation. *Msystems*, 00961–22 (2023)
- [18] Weiss, S., Van Treuren, W., Lozupone, C., Faust, K., Friedman, J., Deng, Y., Xia, L.C., Xu, Z.Z., Ursell, L., Alm, E.J., *et al.*: Correlation detection strategies in microbial data sets vary widely in sensitivity and precision. *The ISME journal* **10**(7), 1669–1681 (2016)
- [19] Berry, D., Widder, S.: Deciphering microbial interactions and detecting keystone species with co-occurrence networks. *Frontiers in microbiology* **5**, 219 (2014)
- [20] Guo, B., Zhang, L., Sun, H., Gao, M., Yu, N., Zhang, Q., Mou, A., Liu, Y.: Microbial co-occurrence network topological properties link with reactor parameters and reveal importance of low-abundance genera. *npj Biofilms and Microbiomes* **8**(1), 3 (2022)
- [21] Ma, B., Wang, Y., Ye, S., Liu, S., Stirling, E., Gilbert, J.A., Faust, K., Knight, R., Jansson, J.K., Cardona, C., *et al.*: Earth microbial co-occurrence network reveals interconnection pattern across microbiomes. *Microbiome* **8**, 1–12 (2020)
- [22] Guidi, L., Chaffron, S., Bittner, L., Eveillard, D., Larhlimi, A., Roux, S., Darzi, Y., Audic, S., Berline, L., Brum, J.R., *et al.*: Plankton networks driving carbon export in the oligotrophic ocean. *Nature* **532**(7600), 465–470 (2016)
- [23] Röttjers, L., Faust, K.: Manta: A clustering algorithm for weighted ecological networks. *Msystems* **5**(1), 10–1128 (2020)
- [24] Levy, R., Borenstein, E.: Reverse ecology: from systems to environments and

- back. In: *Evolutionary Systems Biology*, pp. 329–345. Springer, ??? (2012)
- [25] 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)
- [26] 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)
- [27] Parter, M., Kashtan, N., Alon, U.: Environmental variability and modularity of bacterial metabolic networks. *BMC evolutionary biology* **7**, 1–8 (2007)
- [28] 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)
- [29] Kreimer, A., Doron-Faigenboim, A., Borenstein, E., Freilich, S.: Netcmt: a network-based tool for calculating the metabolic competition between bacterial species. *Bioinformatics* **28**(16), 2195–2197 (2012)
- [30] 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)
- [31] Belcour, A., Frioux, C., Aite, M., Bretaudeau, A., Hildebrand, F., Siegel, A.: Metage2metabo, microbiota-scale metabolic complementarity for the identification of key species. *Elife* **9**, 61968 (2020)
- [32] Thiele, I., Palsson, B.Ø.: A protocol for generating a high-quality genome-scale metabolic reconstruction. *Nature protocols* **5**(1), 93–121 (2010)
- [33] Durot, M., Bourguignon, P.-Y., Schachter, V.: Genome-scale models of bacterial metabolism: reconstruction and applications. *FEMS microbiology reviews* **33**(1), 164–190 (2008)
- [34] 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
- [35] Little, A.E., Robinson, C.J., Peterson, S.B., Raffa, K.F., Handelsman, J.: Rules of engagement: interspecies interactions that regulate microbial communities. *Annu. Rev. Microbiol.* **62**, 375–401 (2008)

- [36] Zientz, E., Dandekar, T., Gross, R.: Metabolic interdependence of obligate intracellular bacteria and their insect hosts. *Microbiology and Molecular Biology Reviews* **68**(4), 745–770 (2004)
- [37] Kallus, Y., Miller, J.H., Libby, E.: Paradoxes in leaky microbial trade. *Nature communications* **8**(1), 1361 (2017)
- [38] 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) <https://doi.org/10.1021/ci3005379> <https://doi.org/10.1021/ci3005379>. PMID: 23384306
- [39] Hessler, T., Huddy, R.J., Sachdeva, R., Lei, S., Harrison, S.T., Diamond, S., Banfield, J.F.: Vitamin interdependencies predicted by metagenomics-informed network analyses and validated in microbial community microcosms. *Nature Communications* **14**(1), 4768 (2023)
- [40] Wattam, A.R., Davis, J.J., Assaf, R., Boisvert, S., Brettin, T., Bun, C., Conrad, N., Dietrich, E.M., Disz, T., Gabbard, J.L., *et al.*: Improvements to patric, the all-bacterial bioinformatics database and analysis resource center. *Nucleic acids research* **45**(D1), 535–542 (2017)
- [41] Parks, D.H., Chuvochina, M., Rinke, C., Mussig, A.J., Chaumeil, P.-A., Hugenholtz, P.: Gtdb: an ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and complete genome-based taxonomy. *Nucleic acids research* **50**(D1), 785–794 (2022)
- [42] Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies, J., Glöckner, F.O.: The silva ribosomal rna gene database project: improved data processing and web-based tools. *Nucleic acids research* **41**(D1), 590–596 (2012)
- [43] Schoch, C.L., Ciufo, S., Domrachev, M., Hotton, C.L., Kannan, S., Khovanskaya, R., Leipe, D., Mcveigh, R., O’Neill, K., Robbertse, B., *et al.*: Ncbi taxonomy: a comprehensive update on curation, resources and tools. *Database* **2020**, 062 (2020)
- [44] Alishum, A.: DADA2 Formatted 16S rRNA Gene Sequences for Both Bacteria & Archaea. <https://doi.org/10.5281/zenodo.6655692> . <https://doi.org/10.5281/zenodo.6655692>
- [45] Murali, A., Bhargava, A., Wright, E.S.: Idtaxa: a novel approach for accurate taxonomic classification of microbiome sequences. *Microbiome* **6**(1), 1–14 (2018)
- [46] Wright, E.S.: Using decipher v2. 0 to analyze big biological sequence data in r. *R Journal* **8**(1) (2016)

- [47] Sunagawa, S., Coelho, L.P., Chaffron, S., Kultima, J.R., Labadie, K., Salazar, G., Djahanschiri, B., Zeller, G., Mende, D.R., Alberti, A., *et al.*: Structure and function of the global ocean microbiome. *Science* **348**(6237), 1261359 (2015)
- [48] Louca, S., Parfrey, L.W., Doebeli, M.: Decoupling function and taxonomy in the global ocean microbiome. *Science* **353**(6305), 1272–1277 (2016)
- [49] 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)
- [50] Feldbauer, R., Schulz, F., Horn, M., Rattei, T.: Prediction of microbial phenotypes based on comparative genomics. *BMC bioinformatics* **16**(14), 1–8 (2015)
- [51] 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)
- [52] 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)
- [53] 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)
- [54] 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)
- [55] Kanehisa, M., Sato, Y.: Kegg mapper for inferring cellular functions from protein sequences. *Protein Science* **29**(1), 28–35 (2020)
- [56] Kanehisa, M., Sato, Y., Kawashima, M.: Kegg mapping tools for uncovering hidden features in biological data. *Protein Science* **31**(1), 47–53 (2022)
- [57] Lam, T.J., Stambouliau, M., Han, W., Ye, Y.: Model-based and phylogenetically adjusted quantification of metabolic interaction between microbial species. *PLoS computational biology* **16**(10), 1007951 (2020)
- [58] 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)
- [59] Overbeek, R., Olson, R., Pusch, G.D., Olsen, G.J., Davis, J.J., Disz, T., Edwards,

- 750 R.A., Gerdes, S., Parrello, B., Shukla, M., *et al.*: The seed and the rapid anno-
751 tation of microbial genomes using subsystems technology (rast). *Nucleic acids*
752 *research* **42**(D1), 206–214 (2014)
- 753 [60] Brettin, T., Davis, J.J., Disz, T., Edwards, R.A., Gerdes, S., Olsen, G.J., Olson,
754 R., Overbeek, R., Parrello, B., Pusch, G.D., *et al.*: Rasttk: a modular and extensi-
755 ble implementation of the rast algorithm for building custom annotation pipelines
756 and annotating batches of genomes. *Scientific reports* **5**(1), 1–6 (2015)
- 757 [61] Merkel, D., *et al.*: Docker: lightweight linux containers for consistent development
758 and deployment. *Linux j* **239**(2), 2 (2014)
- 759 [62] Reese, W.: Nginx: The high-performance web server and reverse proxy. *Linux J.*
760 **2008**(173) (2008)
- 761 [63] Choudhary, K., Meng, E.C., Diaz-Mejia, J.J., Bader, G.D., Pico, A.R., Morris,
762 J.H.: scnetviz: from single cells to networks using cytoscape. *F1000Research* **10**
763 (2021)
- 764 [64] Chaumeil, P.-A., Mussig, A.J., Hugenholtz, P., Parks, D.H.: GTDB-Tk: a toolkit
765 to classify genomes with the Genome Taxonomy Database. Oxford University
766 Press (2020)
- 767 [65] Carbajal-Rodríguez, I., Stöveken, N., Satola, B., Wübbeler, J.H.,
768 Steinbüchel, A.: Aerobic degradation of mercaptosuccinate by the gram-
769 negative bacterium *Halovivax paradoxus* strain b4. *Journal of*
770 *Bacteriology* **193**(2), 527–539 (2011) <https://doi.org/10.1128/jb.00793-10>
771 <https://journals.asm.org/doi/pdf/10.1128/jb.00793-10>
- 772 [66] Sun, J., Matsumoto, K., Nduko, J.M., Ooi, T., Taguchi, S.: Enzymatic charac-
773 terization of a depolymerase from the isolated bacterium *Halovivax* sp. c34 that
774 degrades poly(enriched lactate-co-3-hydroxybutyrate). *Polymer Degradation and*
775 *Stability* **110**, 44–49 (2014) [https://doi.org/10.1016/j.polymdegradstab.2014.08.](https://doi.org/10.1016/j.polymdegradstab.2014.08.013)
776 [013](https://doi.org/10.1016/j.polymdegradstab.2014.08.013)
- 777 [67] Astafyeva, Y., Gurschke, M., Qi, M., Bergmann, L., Indenbirken,
778 D., Grahl, I., Katzowitsch, E., Reumann, S., Hanelt, D., Alawi, M.,
779 Streit, W.R., Krohn, I.: Microalgae and bacteria interaction—evidence
780 for division of labour in the alga microbiota. *Microbiology Spec-*
781 *trum* **10**(4), 00633–22 (2022) <https://doi.org/10.1128/spectrum.00633-22>
782 <https://journals.asm.org/doi/pdf/10.1128/spectrum.00633-22>
- 783 [68] Llaveró-Pasquina, M., Geisler, K., Holzer, A., Mehrshahi, P., Mendoza-Ochoa,
784 G.I., Newsad, S.A., Davey, M.P., Smith, A.G.: Thiamine metabolism genes
785 in diatoms are not regulated by thiamine despite the presence of predicted
786 riboswitches. *New Phytologist* **235**(5), 1853–1867 (2022)

- 787 [69] Romine, M.F., Rodionov, D.A., Maezato, Y., Osterman, A.L., Nelson, W.C.:
788 Underlying mechanisms for syntrophic metabolism of essential enzyme cofactors
789 in microbial communities. *The ISME journal* **11**(6), 1434–1446 (2017)
- 790 [70] Chklovski, A., Parks, D.H., Woodcroft, B.J., Tyson, G.W.: Checkm2: a rapid,
791 scalable and accurate tool for assessing microbial genome quality using machine
792 learning. *Nature Methods* **20**(8), 1203–1212 (2023)
- 793 [71] Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P., Tyson, G.W.:
794 Checkm: assessing the quality of microbial genomes recovered from isolates, single
795 cells, and metagenomes. *Genome research* **25**(7), 1043–1055 (2015)
- 796 [72] Mendoza, S.N., Olivier, B.G., Molenaar, D., Teusink, B.: A systematic assessment
797 of current genome-scale metabolic reconstruction tools. *Genome biology* **20**(1),
798 1–20 (2019)
- 799 [73] Dal Co, A., Vliet, S., Kiviet, D.J., Schlegel, S., Ackermann, M.: Short-range
800 interactions govern the dynamics and functions of microbial communities. *Nature*
801 *ecology & evolution* **4**(3), 366–375 (2020)
- 802 [74] Kost, C., Patil, K.R., Friedman, J., Garcia, S.L., Ralser, M.: Metabolic exchanges
803 are ubiquitous in natural microbial communities. *Nature Microbiology*, 1–9 (2023)