

	001
	002
	003
	004
	005
Genomic, metabolic and literature oriented annotation of microbial co-occurrence networks enhances associations confidence level and hypothesis generation	006
	007
	008
	009
	010
	011
	012
	013
	014
Haris Zafeiropoulos <sup>1</sup> , Ermis Ioannis Michail Delopoulos <sup>1</sup> , Andi Erega <sup>2</sup> , Annelies Geirnaert <sup>2</sup> , John Morris <sup>3</sup> , Karoline Faust <sup>1*</sup>	015
<sup>1*</sup> Department of Microbiology, Immunology and Transplantation, Rega Institute for Medical Research , KU Leuven, Herestraat, Leuven, 3000, , Belgium .	016
<sup>2</sup> Institute of Food, Nutrition and Health, ETH Zurich, Street, Zurich, 8092, , Switzerland .	017
<sup>3</sup> Department of Pharmaceutical Chemistry, University of California San Francisco, Street, San Francisco, 94143, California, USA .	018
	019
	020
	021
	022
	023
	024
	025
	026
	027
*Corresponding author(s). E-mail(s): <a href="mailto:karoline.faust@kuleuven.be">karoline.faust@kuleuven.be</a> ;	028
Contributing authors: <a href="mailto:haris.zafeiropoulos@kuleuven.be">haris.zafeiropoulos@kuleuven.be</a> ;	029
<a href="mailto:ermisioannis.michaildelopoulos@student.kuleuven.be">ermisioannis.michaildelopoulos@student.kuleuven.be</a> ;	030
<a href="mailto:andi.erega@hest.ethz.ch">andi.erega@hest.ethz.ch</a> ; <a href="mailto:annelies.geirnaert@hest.ethz.ch">annelies.geirnaert@hest.ethz.ch</a> ;	031
<a href="mailto:scooter@cgl.ucsf.edu">scooter@cgl.ucsf.edu</a> ;	032
	033
	034
<b>Abstract</b>	035
Up to 350 words.	036
The abstract must include the following separate sections:	037
<b>Background:</b> the context and purpose of the study	038
<b>Results:</b> the main findings	039
<b>Conclusions:</b> a brief summary and potential implications	040
	041
	042
	043
	044
	045
	046

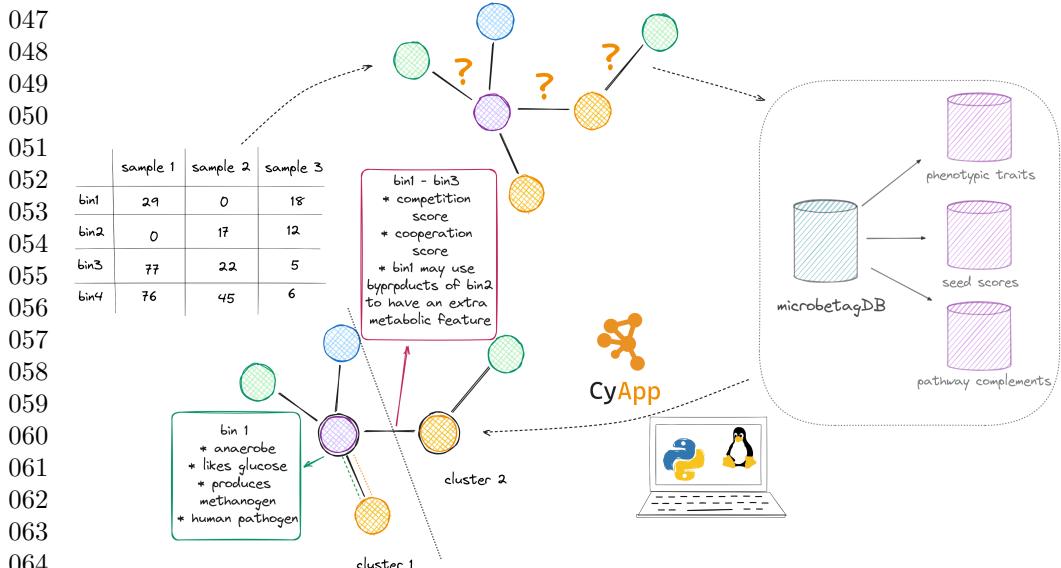


Figure abstract.

065  
066  
067  
068     **Keywords:** microbial associations, enrichment analysis, data integration, pathway  
069     complementarity, seed set

## 072     1 Introduction

073     1 2

- 
- 074     • Role of microbial associations in the study of microbial communities  
075     • HTP sequencing  
076     • cooccurrence networks to analyse HTP data  
077     • challenges of cooccurrence networks / refs: [1]: microbial network inference as a tool  
078     for interaction prediction has highlighted this tool's low accuracy and the biological  
079     implications of network properties are unclear / [2] : modules in microbial co-  
080     occurrence networks may be indicative of ecological processes governing community  
081     structure, such as niche filtering and habitat preference  
082     • current approaches to deal with those challenges  
083     • our contribution

---

084     It is well known that most microbial species live not in isolation but in commu-  
085     nities [3]. Such communities play a crucial role in ecosystem functioning in almost

---

086     <sup>1</sup>We are to submit in the Microbiome journal as a "Software" manuscript, thus we follow [these rules](#)  
087     <sup>2</sup>The introduction should not include subheadings.

any ecosystem type [4, 5]. High-throughput sequencing (HTP) has provided great insight into the diversity and composition of microbial communities [6]. Uncultivated species can now be detected and their features can be inferred through their genomic information [7].	093
NGS facilitates culture-independent sampling of the microorganisms in an area with the potential for both taxonomic and functional annotation;	094
NGS not only provides information about the taxonomic composition of the microbial community but also enables the annotation of their functional capabilities. This combination of taxonomic and functional annotation provides a more detailed and holistic view of the role and activities of microorganisms.	095
Understanding the structure, composition, and dynamics of microbial communities is crucial for unraveling the intricate web of interactions that shape ecosystems. To come up with a comprehensive understanding of these interactions, we need to comprehend the relationships between microorganisms and thus their own interactions.	096
The notion of an interaction varies including cooperation, competition, parasitism, commensalism and amensalism [5].	097
Microbial interactions play a crucial role in shaping ecosystems and influencing various biological processes.	098
These interactions contribute to the overall stability and functioning of ecosystems by influencing nutrient cycling, disease dynamics, and the overall diversity and composition of microbial communities. For example, some microorganisms can produce and release certain compounds that benefit neighboring microbes or inhibit the growth of competing species. Other interactions involve symbiotic relationships, where different microbes rely on each other for survival and perform complementary functions. Understanding and studying microbial interactions is vital for unraveling the complexity of microbial communities and their impact on human health, agriculture, and environmental processes.	099
Physical Interactions: Microorganisms can interact physically by forming biofilms, which are communities of microorganisms attached to surfaces and enclosed in a matrix of extracellular substances. Biofilms provide protection and promote cooperation among the microorganisms within them.	100
Quorum Sensing: Many bacteria use quorum sensing to communicate with each other and coordinate their behavior. They release and detect signaling molecules called autoinducers, which allow them to sense the local population density. This mechanism enables bacteria to coordinate processes like biofilm formation, virulence factor expression, and nutrient acquisition.	101
Antibiosis: Antibiosis refers to the production of antimicrobial substances by microorganisms that inhibit the growth or survival of other microorganisms. This can be a competitive strategy to gain an advantage in a particular environment.	102
metabolic interactions specifically, microorganisms often engage in metabolic exchanges during their interactions. For instance, one microbe may produce metabolites that serve as nutrients or signaling molecules for another microbe. These metabolic interactions can involve the transfer of essential nutrients, the breakdown of complex compounds, or the production of secondary metabolites with antimicrobial properties. Overall, understanding the different types of microbial interactions,	103
	104
	105
	106
	107
	108
	109
	110
	111
	112
	113
	114
	115
	116
	117
	118
	119
	120
	121
	122
	123
	124
	125
	126
	127
	128
	129
	130
	131
	132
	133
	134
	135
	136
	137
	138

139 including metabolic interactions, provides insights into the complexity and dynamics  
140 of microbial communities and their impact on various processes in nature.

141 A widely used approach is the creation of co-occurrence networks based on  
142 community data. To build such networks, there is a great number of approaches: Spear-  
143 man and Pearson correlations, CoNet [8] SparCC [9] SpeicEasi [10], MAGMA [11]  
144 and FlashWeave [12] are just a few of them. However, the outcome is usually  
145 tool-dependent [3, 13, 14].

146 FlashWeave

147 microbeAnnotator [15]

148 Related literature: Karaoz U and Brodie EL (2022) microTrai [16], a computational  
149 pipeline that infers and distills ecologically relevant traits from microbial genome  
150 sequences. It does not apply networks

151

## 152 2 Implementation

153

154 <sup>3</sup>

155

### 156 Genomes included

157 Using the GTDB v202 [metadata files](#), we retrieved the NCBI genome accessions of the  
158 representative genomes of high quality, i.e. completeness  $\geq 95\%$  and contamination  
159  $\leq 5\%$ . That resulted a set of 26,778 covering 22,009 unique NCBI Taxonomy Ids.  
160 Using these accession numbers, we were able to download their corresponding .faa  
161 files when available ([get\\_gtdb\\_faa.py](#)) leading to a set of 16,900 amino acid sequence  
162 files.  
163

### 164 Taxonomy schemes

165

166 microbetag maps the taxonomy of each entry in the abundance table to its correspond-  
167 ing NCBI Taxonomy id and if available its closest GTDB representative genome(s).  
168 Two well established taxonomy schemes are supported. The Genome Taxonomy  
169 DataBase (GTDB) [17] that is being broadly used in bins and/or MAGs taxonomical  
170 classification and the Silva database [18] that has NCBI Taxonomy [19]. The primer  
171 links the representative genomes included to their corresponding NCBI Taxonomy ids  
172 too.

173 There is a great number of taxonomies that are being used in such studies, e.g.  
174 Silva [18], Ribosomal Database Project (RDP) [20], manually curated ones and more,  
175 As a consequence, there is not a standardised format of the taxonomies assigned, from  
176 bioinformatics pipelines used for the analysis of such data. microbetag makes use of  
177 the [fuzzywuzzy](#) library that implements the Levenshtein Distance Metric to get the  
178 closest NCBI taxon name and thus its corresponding NCBI Taxonomy id. ++ ncbi  
179 nodes dump A relatively high similarity score is used (90) to avoid false positives.

180 DADA2 formatted 16S rRNA gene sequences for both bacteria and archaea [21]  
181 were used to trained the TAXID classifier [22] of the DECIPHER package.  
182

183

---

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

<b>Network inference</b>	185
FlashWeave [12]	186
a computational approach based on a flexible Probabilistic Graphical Model framework that integrates metadata and predicts direct microbial interactions from heterogeneous microbial abundance data sets with hundreds of thousands of samples.	187
A flexible Probabilistic Graphical Model framework is used in a computational approach that incorporates metadata and predicts direct microbial interactions. This is done using heterogeneous microbial abundance datasets consisting of hundreds of thousands of samples.	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
<b>Literature oriented node annotation</b>	
Using a set of Tara Oceans samples [23] FAPROTAX [24] estimates the functional potential of the bacterial and archaeal communities, by classifying each taxonomic unit into functional group(s) based on current literature, announcements of cultured representatives and/or manuals of systematic microbiology. In this manually curated approach, a taxon is associated with a function if and only if all the cultured species within the taxon have been shown to exhibit that function. In its current version, FAPROTAX includes more than 80 functions based on 7600 functional annotations and covering more than 4600 taxa. Contrary to gene content based approaches, e.g. PICRUSt2, FAPROTAX estimates metabolic phenotypes based on experimental evidence.	
microbetag invokes the accompanying script of FAPROTAX and converts the taxonomic microbial community profile of the samples included in the user's abundance table or of the taxa present in the provided network, into putative functional profiles. Then, it parses FAPROTAX's subtables to annotate each taxonomic unit present on the user's data with all the functions for which they had a hit. FAPROTAX annotations are not part of the microbetagDB but are computed on the fly.	
<b>Genomic oriented node annotation</b>	
phenDB [25] is a publicly available resource that supports the analysis of bacterial (meta)genomes to identify 47 distinct functional traits. It relies on support vector machines (SVM) trained with manually curated datasets based on gene presence/absence patterns for trait prediction. More specifically, the model for a particular trait is trained using a collection of EggNOG annotated genomes where the knowledge of whether that trait is present or absent among its members is available. The <code>compute-genotype</code> program of phenotrex supports the creation of such tabular <i>genotype</i> files. A <i>genotype</i> file can be used along with a <i>phenotype</i> one, i.e., a file containing true phenotypic trait values for each input genome on which to train the model, and the <code>train</code> program of phenotrex can then be performed. Last, the models can now be used to predict their corresponding traits; based on the completeness/contamination of the genomes, the accuracy varies.	
In the frameowrk of microbetagDB, phenotrex classifiers were re-trained using the genomes provided by phenDB for each trait to sync with the latest version of eggNOG.	

231 Genomes were downloaded from NCBI using the [Batch Entrez](#) program. Then, *geno-*  
 232 *type* files were produced for all the high quality GTDB representative genomes. Each  
 233 model was then used against all the GTDB *genotype* files to annotate each with the  
 234 presence or the absence of the trait.

235

## 236 Pathway complementarity

237

238 For the subset of the 16,900 high quality GTDB representative genomes that a .faa  
 239 was available, *kofamscan* [26] was performed to annotate them with KEGG ORTHOL-  
 240 OGY terms (KOs) [27]. Their KOs were then mapped to their corresponding KEGG  
 241 modules. A KEGG module is defined as a functional unit within the KEGG frame-  
 242 work, that represents a set of enzymes and reactions involved in a specific biological  
 243 process or pathway [28]. A module's definition is a logical expression and consists of  
 244 KOs and the following symbols: a. the space, representing a connection in the pathway  
 245 b. plus sign, representing a molecular complex, c. comma, representing alternatives  
 246 and d. minus sign, designates an optional item in the complex. Both (a) and (b) cases  
 247 should be considered as "AND" logical operators, while (c) would be the "OR".

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

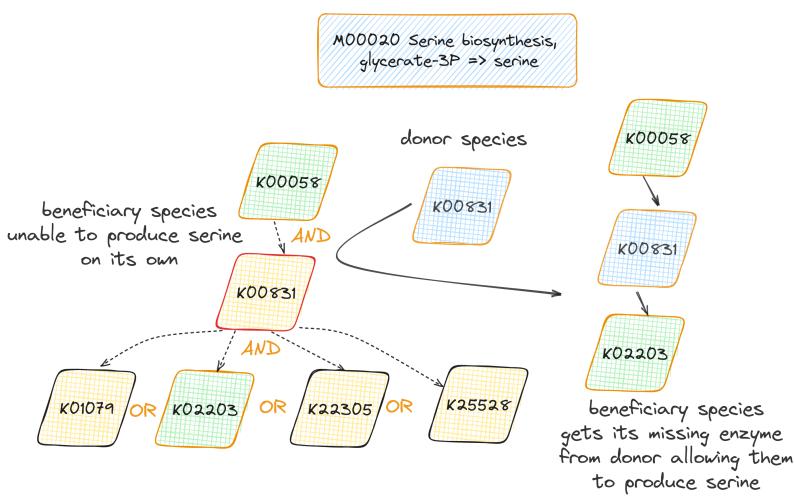
264

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

274

275

276



We define a genome as having a "complete" module if and only if all of the KOs present in any of the module's alternatives are also found among the annotated KOs of the genome. All modules definitions were retrieved using the KEGG API and parsed ([parse\\_module\\_definitions.py](#)). A dictionary was built with all the alternatives, i.e. alternative sets of KOs, for a module to be complete ([module\\_definition\\_map.json](#)). 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 ([pathway\\_complementarity.py](#)). Thanks to the graphical user interface (GUI) of the [KEGG pathway map viewer](#) [29, 30], each complementarity can be visualised as part of the closest KEGG metabolic map; where the KOs coming from the donor are shown with a blue-green colour, while those from the beneficiary's genome itself with rose.

As several GTDB representative genomes might map to the same NCBI Taxonomy Id, all the possible genomes' combinations are annotated in the edge of a pair of species level taxonomically annotated OTUs/ASVs/bins. On top of that, as co-occurrence networks are undirected, both nodes of a suggested association are considered as potential donors and beneficiary species.

## Seed scores using genome scale metabolic reconstructions

A metabolic network's "seed set" is the set of compounds that, based on the network topology, need to be acquired exogenously [31]. 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.

Based on the seed concept, several graph theory-based metrics have been described to predict species interactions directly from their networks' topologies. 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 [32], 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 [33] and NetCompt [34] tools correspondingly. We will be referring to those two indices as "seed scores". Most recently, the PhyloMint Python package [32] was released supporting the calculation of the seed scores of genome scale metabolic network reconstructions (GENREs) in SBML format.

323 In the framework of microbetag, seed scores were computed using PhyloMint and  
324 draft GENREs for all pair-wised combinations of GTDB representative genomes that  
325 have been RAST annotated in the framework of the PATRIC database [35]. GENREs  
326 were reconstructed using the Model SEED pipeline [36] through its Python interface  
327 [ModelSEEDpy](#).

328

### 329 **Clustering network**

330 manta is a heuristic network clustering algorithm that clusters nodes within weighted  
331 networks effectively, leveraging the presence of negative edges and discerning between  
332 weak and microbetag invokes manta [37] to infer clusters from the microbial network.  
333 A taxonomically-informed layout is  
334 strong cluster assignments. ++ taxonomy layout  
335

### 336 **Groups of annotations**

337 Biologically meaningful groups were built using the micrO ontology [38].

338

### 340 **Building the CytoscapeApp**

341 The microbetag CytoscapeApp was build based on the [source code](#) of the scVizNet [39].  
342 Java @Ermis to add  
343 Enrichment analysis is supported. Hypergeometric distribution FDR +++

### 345 **Dependencies, Web server and API**

346 The microbetag web service is container - based and consists of three Docker [40]  
347 (v24.0.2) images: a. the [MySQL](#) database b. an nginx [41] web server and c. the app  
348 itself. The latter uses [Gunicorn](#) (20.1.0) to build an application server which commu-  
349 nicates with the web server using the Web Server Gateway Interface (WSGI) protocol  
350 and handles incoming HTTP requests. microbetag is implemented as a [Flask](#) applica-  
351 tion (v2.3.2); Flask is a micro web framework for developing Python web applications  
352 and RESTful APIs. A thorough description of microbetag's API is available at the  
353 [ReadTheDocs web site](#). The source code of the microbetag web service is available  
354 on [GitHub](#).

355 python 3.11 slim docker image julia 1.7.1 for flashweave mysql.connector 8.0.27  
356 python library pandas 2.1.1. numpy 1.26.0 multiprocessing

357 text processing using awk  
358 KEGG API  
359

360

## 361 **2.1 Running large datasets**

## 362 **3 Results**

363 [4](#)

364

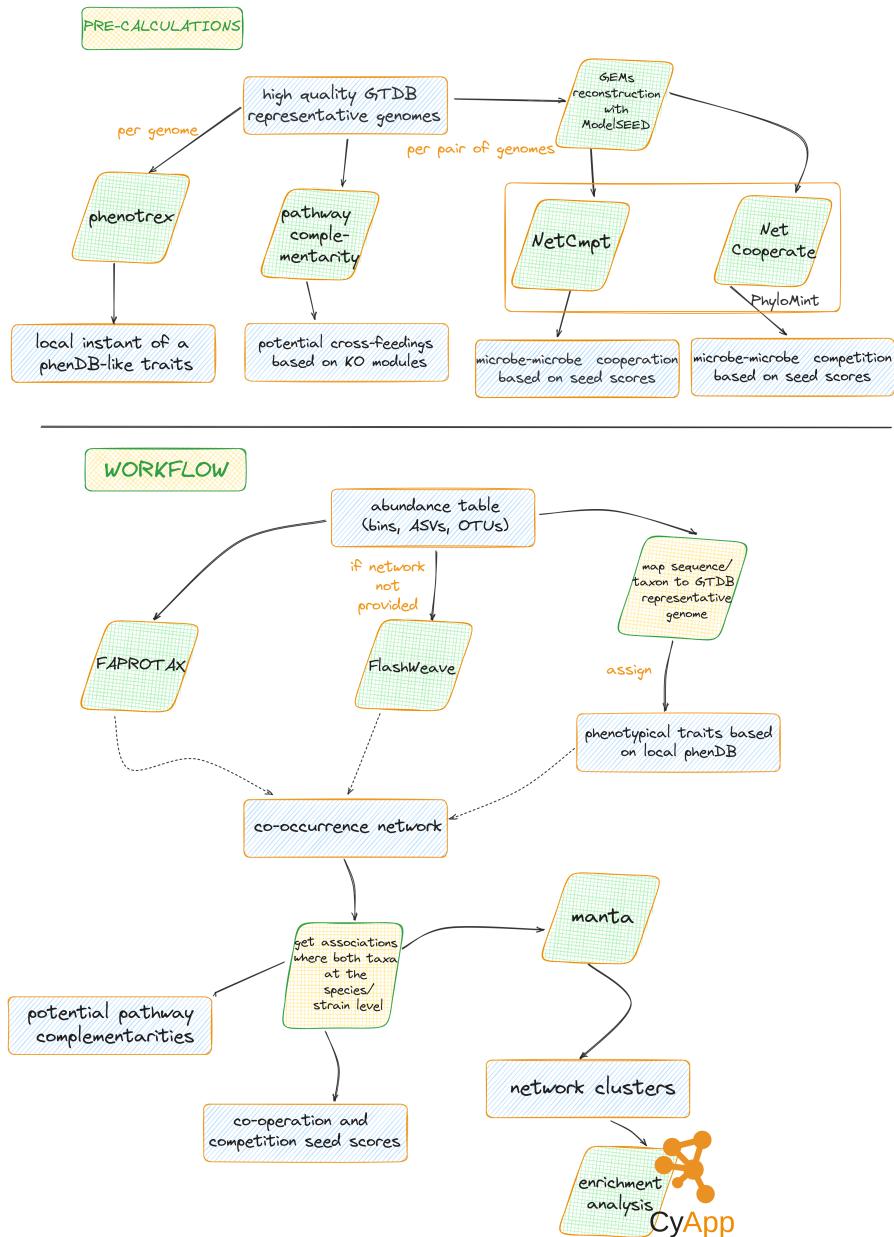
365

---

366 <sup>4</sup>Significant advance over previously published software (usually demonstrated by direct comparison with  
367 available related software) This should include the findings of the study including, if appropriate, results of  
368 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.

## microbetag and microbetagDB

369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414



**Fig. 2:** Diagram of the microbetag pre - calculations and the on the fly workflow. GTDB v207 representative genomes were filtered and for those of high-quality 33 phenotypic traits were predicted using phenotrex [25]. To this end, models were re-trained to sync with recent version of eggNOG [42].

415 microbetag in numbers: 34,608 GTDB representative genomes 32 phen-model-oriented  
 416 metabolic functions 92 FAPROTAX functions 341,568 unique complements involved  
 417 in > 184 million beneficiary - donor pairs' complementarities 30,755 GENREs leading  
 418 to 1 billion competition and complementarity scores

419 annotated network returned in .cyjs format

420 For a computationally efficient way to annotate large networks, a Docker image is  
 421 provided so the user runs a taxonomy assignment using the IDTAXA algorithm [22]  
 422 of the DECIPHER R package [43]. A co-occurrence network is also built using  
 423 FlashWeave [12], as microbetag also does.

424

## 425 **microbetag CytoscapeApp**

426 Overall comment, the CytoscapeApp returns averages and s.d. for example in seed  
 427 scores. If you want the exact values, go through the API.  
 428

429

### 430 A. GTDB-tk: 480 bins

431 Step	Time(sec)	Notes
432 Taxonomy mapping	Cell 1,2	on the fly
433 Network inference	Cell 2,2	on the fly
434 microbetag annotations	Cell 3,2	on the fly
435 manta clustering	Cell 4,2	on the fly

436

437

### 438 C. Silva:

439 Step	Time(sec)	Notes
440 Taxonomy mapping	Cell 1,2	Cell 1,3
441 Network inference	Cell 2,2	Cell 2,3
442 microbetag annotations	Cell 3,2	Cell 3,3
443 manta clustering	Cell 4,2	Cell 4,3

444 **Table 1:** Computing times per step using an abundance table of 400 taxa with  
 445 taxonomy: A. taxonomy scheme B. C. D. <sup>5</sup> specs of the laptop used.

446

447

## 448 **Validation of microbetag potential**

449 vitamin dataset [44]

450

## 452 **Interpetating a real-world network with microbetag**

453 Annelies' dataset.

454

## 455 **4 Discussion**

456 <sup>6</sup>

458

---

459 <sup>6</sup>The user interface should be described and a discussion of the intended uses of the software, and the  
 460 benefits that are envisioned, should be included, together with data on how its performance and functionality

<b>5 Conclusions</b>	461
<sup>7</sup>	462
Data integration	463
Supplementary information. <sup>8</sup>	464
<b>Declarations</b>	465
• <b>Availability of data and materials</b>	466
– Raw sequences for the use case:	467
– Raw data for the validations case:	468
• <b>Funding</b>	469
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.	470
• <b>Conflict of interest/Competing interests</b>	471
The authors declare that they have no other competing interests.	472
• <b>Authors' contributions</b> <sup>9</sup>	473
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., A.E.	474
• <b>Acknowledgements</b>	475
We would like to thank Dr Christina Pavloudi and ++ for the insight on how to organise the trait groups.	476
• <b>Ethics approval</b>	477
Not applicable	478
• <b>Consent to participate</b>	479
Not applicable.	480
• <b>Code availability:</b>	481
– microbetagDB related scripts: <a href="https://github.com/hariszaf/microbetag">https://github.com/hariszaf/microbetag</a>	482
– microbetagApp and webserver: <a href="https://github.com/msysbio/microbetagApp">https://github.com/msysbio/microbetagApp</a> .	483
– CytoscapeApp: <a href="https://github.com/ermisemd/MGG/">https://github.com/ermisemd/MGG/</a>	484
– Validation and use case: {think of having that under the 3D'omics organization}	485
– Documentation web-site: <a href="https://hariszaf.github.io/microbetag/">https://hariszaf.github.io/microbetag/</a>	486
<hr/>	487
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.	488
<sup>7</sup> This should state clearly the main conclusions and provide an explanation of the importance and relevance of the case, data, opinion, database or software reported.	489
<sup>8</sup> If your article has accompanying supplementary file(s) please state so here. E.g. supplementary figures and tables captions.	490
<sup>9</sup> Based on the <a href="#">CRediT system</a> . Current list is indicative.	491
	492
	493
	494
	495
	496
	497
	498
	499
	500
	501
	502
	503
	504
	505
	506

507 **Appendix A Background on seed scores and  
508 complementarities**  
509

510 **A.1 Background on seed scores**  
511

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

515 
$$C_i = 1/\text{seed}'s \text{ group with } i \text{ size} \quad (\text{A1})$$

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

518 
$$MI_{Complementarity} = \frac{|\text{SeedSet}_A \cap \neg\text{SeedSet}_B|}{|\text{SeedSet}_A \cap (\text{SeedSet}_B \cup \neg\text{SeedSet}_B)|} \quad (\text{A2})$$

519 As also described in [32], it is calculated as the proportion of compounds in a  
520 species' seed set that coincide with those in an other's, while also factoring in the  
521 confidence scores associated with seed compounds.

522 
$$MI_{Competition} = \frac{\sum C(\text{SeedSet}_A \cap \text{SeedSet}_B)}{\sum C(\text{SeedSet}_A)} \quad (\text{A3})$$

523 **A.2 Background on pathway complementarity**  
524

525 For example, the definition of the D-Galacturonate degradation in Bacteria (M00631)  
526 is:

527 K01812 K00041 (K01685,K16849+K16850) K00874 (K01625,K17463)  
528 that once breaking down, it leads to 4 alternative sets of KOs (pathways):  
529  
530 K01812 K00041 K01685 K00874 K01625  
531 K01812 K00041 K16849+K16850 K00874 K01625  
532 K01812 K00041 K01685 K00874 K17463  
533 K01812 K00041 K16849+K16850 K00874 K17463

534 **A.3 Complementarities**  
535

536 KEGG compound ModelSEED compounds ModelSEED compounds mapped to  
537 KEGG compounds and kept only those related to KEGG modules.  
538

539 **References**  
540

- 541 [1] Berry, D., Widder, S.: Deciphering microbial interactions and detecting keystone  
542 species with co-occurrence networks. *Frontiers in microbiology* **5**, 219 (2014)  
543

- [2] 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) 553  
554  
555  
556
- [3] Röttjers, L., Faust, K.: From hairballs to hypotheses—biological insights from microbial networks. *FEMS microbiology reviews* **42**(6), 761–780 (2018) 557  
558
- [4] Raes, J., Bork, P.: Molecular eco-systems biology: towards an understanding of community function. *Nature Reviews Microbiology* **6**(9), 693–699 (2008) 559  
560  
561
- [5] Faust, K., Raes, J.: Microbial interactions: from networks to models. *Nature Reviews Microbiology* **10**(8), 538–550 (2012) 562  
563  
564
- [6] 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> 565  
566  
567  
568  
569  
570  
571
- [7] 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) 572  
573  
574  
575
- [8] 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) 576  
577  
578  
579
- [9] Friedman, J., Alm, E.J.: Inferring correlation networks from genomic survey data. *PLoS computational biology* **8**(9), 1002687 (2012) 580  
581
- [10] 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) 582  
583  
584  
585
- [11] Cougoul, A., Bailly, X., Wit, E.C.: Magma: inference of sparse microbial association networks. *BioRxiv*, 538579 (2019) 586  
587  
588
- [12] 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) 589  
590  
591  
592
- [13] 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) 593  
594  
595  
596
- [14] Weiss, S., Van Treuren, W., Lozupone, C., Faust, K., Friedman, J., Deng, Y., 597  
598

- 599 Xia, L.C., Xu, Z.Z., Ursell, L., Alm, E.J., *et al.*: Correlation detection strategies  
600 in microbial data sets vary widely in sensitivity and precision. *The ISME journal*  
601 **10**(7), 1669–1681 (2016)
- 602
- 603 [15] Ruiz-Perez, C.A., Conrad, R.E., Konstantinidis, K.T.: Microbeannotator: a user-  
604 friendly, comprehensive functional annotation pipeline for microbial genomes.  
605 *BMC bioinformatics* **22**, 1–16 (2021)
- 606
- 607 [16] Karaoz, U., Brodie, E.L.: microtrait: a toolset for a trait-based representation of  
608 microbial genomes. *Frontiers in Bioinformatics* **2**, 918853 (2022)
- 609
- 610 [17] Parks, D.H., Chuvochina, M., Rinke, C., Mussig, A.J., Chaumeil, P.-A., Hugen-  
611 holtz, P.: Gtdb: an ongoing census of bacterial and archaeal diversity through  
612 a phylogenetically consistent, rank normalized and complete genome-based  
613 taxonomy. *Nucleic acids research* **50**(D1), 785–794 (2022)
- 614
- 615 [18] Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies,  
616 J., Glöckner, F.O.: The silva ribosomal rna gene database project: improved data  
617 processing and web-based tools. *Nucleic acids research* **41**(D1), 590–596 (2012)
- 618
- 619 [19] Schoch, C.L., Ciufo, S., Domrachev, M., Hotton, C.L., Kannan, S., Khovanskaya,  
620 R., Leipe, D., Mcveigh, R., O'Neill, K., Robbertse, B., *et al.*: Ncbi taxonomy:  
621 a comprehensive update on curation, resources and tools. *Database* **2020**, 062  
622 (2020)
- 623
- 624 [20] Cole, J.R., Wang, Q., Fish, J.A., Chai, B., McGarrell, D.M., Sun, Y., Brown,  
625 C.T., Porras-Alfaro, A., Kuske, C.R., Tiedje, J.M.: Ribosomal database project:  
626 data and tools for high throughput rrna analysis. *Nucleic acids research* **42**(D1),  
627 633–642 (2014)
- 628
- 629 [21] Alishum, A.: DADA2 Formatted 16S rRNA Gene Sequences for Both Bacteria  
630 & Archaea. <https://doi.org/10.5281/zenodo.6655692> . <https://doi.org/10.5281/zenodo.6655692>
- 631
- 632 [22] Murali, A., Bhargava, A., Wright, E.S.: Idtaxa: a novel approach for accurate  
633 taxonomic classification of microbiome sequences. *Microbiome* **6**(1), 1–14 (2018)
- 634
- 635 [23] Sunagawa, S., Coelho, L.P., Chaffron, S., Kultima, J.R., Labadie, K., Salazar,  
636 G., Djahanschiri, B., Zeller, G., Mende, D.R., Alberti, A., *et al.*: Structure and  
637 function of the global ocean microbiome. *Science* **348**(6237), 1261359 (2015)
- 638
- 639 [24] Louca, S., Parfrey, L.W., Doebeli, M.: Decoupling function and taxonomy in the  
640 global ocean microbiome. *Science* **353**(6305), 1272–1277 (2016)
- 641
- 642 [25] Feldbauer, R., Schulz, F., Horn, M., Rattei, T.: Prediction of microbial phenotypes  
643 based on comparative genomics. *BMC bioinformatics* **16**(14), 1–8 (2015)
- 644

- [26] 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) 645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
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
- [27] 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)
- [28] 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)
- [29] Kanehisa, M., Sato, Y.: Kegg mapper for inferring cellular functions from protein sequences. *Protein Science* **29**(1), 28–35 (2020)
- [30] Kanehisa, M., Sato, Y., Kawashima, M.: Kegg mapping tools for uncovering hidden features in biological data. *Protein Science* **31**(1), 47–53 (2022)
- [31] 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)
- [32] 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)
- [33] 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)
- [34] 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)
- [35] 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)
- [36] 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)
- [37] Röttgers, L., Faust, K.: Manta: A clustering algorithm for weighted ecological networks. *Msystems* **5**(1), 10–1128 (2020)

- 691 [38] Blank, C.E., Cui, H., Moore, L.R., Walls, R.L.: Micro: an ontology of pheno-  
692 typic and metabolic characters, assays, and culture media found in prokaryotic  
693 taxonomic descriptions. *Journal of biomedical semantics* **7**(1), 1–10 (2016)  
694
- 695 [39] Choudhary, K., Meng, E.C., Diaz-Mejia, J.J., Bader, G.D., Pico, A.R., Morris,  
696 J.H.: scnetviz: from single cells to networks using cytoscape. *F1000Research* **10**  
697 (2021)
- 698 [40] Merkel, D., *et al.*: Docker: lightweight linux containers for consistent development  
699 and deployment. *Linux j* **239**(2), 2 (2014)
- 701 [41] Reese, W.: Nginx: The high-performance web server and reverse proxy. *Linux J.*  
702 **2008**(173) (2008)
- 703
- 704 [42] Huerta-Cepas, J., Szklarczyk, D., Heller, D., Hernández-Plaza, A., Forslund, S.K.,  
705 Cook, H., Mende, D.R., Letunic, I., Rattei, T., Jensen, L.J., *et al.*: eggNOG 5.0:  
706 a hierarchical, functionally and phylogenetically annotated orthology resource  
707 based on 5090 organisms and 2502 viruses. *Nucleic acids research* **47**(D1), 309–314  
708 (2019)
- 709
- 710 [43] Wright, E.S.: Using deciphér v2. 0 to analyze big biological sequence data in r. *R*  
711 *Journal* **8**(1) (2016)
- 712
- 713 [44] Hessler, T., Huddy, R.J., Sachdeva, R., Lei, S., Harrison, S.T., Diamond, S.,  
714 Banfield, J.F.: Vitamin interdependencies predicted by metagenomics-informed  
715 network analyses and validated in microbial community microcosms. *Nature*  
716 *Communications* **14**(1), 4768 (2023)
- 717
- 718
- 719
- 720
- 721
- 722
- 723
- 724
- 725
- 726
- 727
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736