

## 1 **OrtSuite – from genomes to prediction of microbial interactions within targeted ecosystem** 2 **processes**

3 João Pedro Saraiva<sup>1</sup>, Alexandre Bartholomäus<sup>2</sup>, René Kallies<sup>1</sup>, Marta Gomes<sup>3</sup>, Marcos Bicalho<sup>1</sup>,  
4 Carsten Vogt<sup>1</sup>, Antonis Chatzinotas<sup>1,4,5</sup>, Peter Stadler<sup>6,7,8,9,10</sup>, Oscar Dias<sup>3</sup>, Ulisses Nunes da  
5 Rocha<sup>1\*</sup>

6  
7 <sup>1</sup>Department of Environmental Microbiology, Helmholtz Centre for Environmental Research-  
8 UFZ, Leipzig, Germany

9 <sup>2</sup> GFZ German Research Centre for Geosciences, Section Geomicrobiology, Potsdam, Germany

10 <sup>3</sup> Centre of Biological Engineering, University of Minho, Portugal

11 <sup>4</sup>Institute of Biology, Leipzig University, Leipzig, Germany

12 <sup>5</sup>German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Leipzig,  
13 Germany

14 <sup>6</sup>Bioinformatics Group, Department of Computer Science, Interdisciplinary Center for  
15 Bioinformatics, and Competence Center for Scalable Data Services and Solutions  
16 Dresden/Leipzig, University of Leipzig, Leipzig, Germany

17 <sup>7</sup>Max Planck Institute for Mathematics in the Sciences, Leipzig, Germany

18 <sup>8</sup>Institute for Theoretical Chemistry, University of Vienna, Wien, Austria

19 <sup>9</sup>Facultad de Ciencias, Universidad Nacional de Colombia, Bogotá, Colombia

20 <sup>10</sup>Santa Fe Institute, Santa Fe, U.S.A.

21  
22 \*Correspondence: [ulisses.rocha@ufz.de](mailto:ulisses.rocha@ufz.de)  
23

24 Running title: Mining interactions with OrtSuite  
25

26 **Abstract:** The high complexity found in microbial communities makes the identification of  
27 microbial interactions challenging. To address this challenge, we present OrtSuite, a flexible  
28 workflow to predict putative microbial interactions based on genomic content of microbial  
29 communities and targeted to specific ecosystem processes. The pipeline is composed of three user-  
30 friendly bash commands. OrtSuite combines ortholog clustering with genome annotation strategies  
31 limited to user-defined sets of functions allowing for hypothesis-driven data analysis such as  
32 assessing microbial interactions in specific ecosystems. OrtSuite matched, on average, 96 % of  
33 experimentally verified KEGG orthologs involved in benzoate degradation in a known group of  
34 benzoate degraders. Identification of putative synergistic species interactions was evaluated using  
35 the sequenced genomes of an independent study which had previously proposed potential species  
36 interactions in benzoate degradation. OrtSuite is an easy to use workflow that allows for rapid  
37 functional annotation based on a user curated database and can easily be extended to ecosystem  
38 processes where connections between genes and reactions are known. OrtSuite is an open-source  
39 software available at <https://github.com/mdsufz/OrtSuite>.

40 **Keywords:** functional annotation/microbial interactions/microbial modelling/orthologs/partial  
41 genome-scale models.

## 42 Introduction

43 In environments where microorganisms play a key role, the microbial community functional  
44 potential encompasses the building blocks for all possible interspecies interactions (Maestre *et al*,  
45 2012; Mulder *et al*, 2001). For example, in environments rich in methane, microbial communities  
46 are dominated by species with genes encoding proteins involved in methanogenesis (Lyu *et al*,  
47 2018). Soil microbes, especially those in the rhizosphere are genetically adapted to support plants  
48 in the resistance against pathogens and tolerance to stress (Mendes *et al*, 2018). In this context,  
49 natural ecosystems are populated by an enormous number of microbes (Locey & Lennon, 2016).  
50 For example, soil environments can contain more than  $10^{10}$  organisms per gram of soil which are  
51 distributed in a heterogeneous way making a global search for interspecies interactions unfeasible  
52 (Raynaud & Nunan, 2014). The exponential increase in high-throughput sequencing data and the  
53 development of computational sciences and bioinformatics pipelines has advanced our  
54 understanding of microbial community composition and distribution in complex ecosystems (Roh  
55 *et al*, 2010). This knowledge increased our ability to reconstruct and functionally characterize  
56 genomes in complex communities, for example by the recovery of metagenome-assembled  
57 genomes (MAGs) (Parks *et al*, 2017; Pasolli *et al*, 2019; Tully *et al*, 2018). While several tools  
58 have been developed to improve the reconstruction of MAGs, the same cannot be said for  
59 predicting interspecies interactions (Morin *et al*, 2018). Studies by Parks (Parks *et al*, 2017) and  
60 Tully (Tully *et al*, 2018), while advancing the reconstruction of MAGs, did not perform any  
61 functional characterization or prediction of interspecies interactions. Pasolli and collaborators  
62 (Pasolli *et al*, 2019) performed functional annotation of representative species in their study by  
63 employing several tools such as EggNOG (Huerta-Cepas *et al*, 2017), KEGG (Kanehisa *et al*,  
64 2004) and DIAMOND (Buchfink *et al*, 2015). However, the sheer number of representative  
65 genomes (4930) and the lack of focus on specific ecosystem processes makes predicting  
66 interspecies interactions a challenge. Furthermore, the challenge of predicting interspecies  
67 interactions increases due to the multitude of potential interactions not only between species in  
68 microbial communities but also between microbes and their hosts (e.g., plants, animals and  
69 microeukaryotes) (Slade *et al*, 2017). An integrated pipeline for annotation and visualization of  
70 metagenomes (MetaErg) developed by Dong and Strous (Dong & Strous, 2019) attempt to address  
71 some of the challenges in metagenome annotation such as the inference of biological functions  
72 and integration of expression data. MetaErg performs comprehensive annotation and visualization  
73 of MAGs by integrating data from multiple sources such as Pfam (Mistry *et al*, 2021), KEGG  
74 (Kanehisa *et al*, 2004) and FOAM (Prestat *et al*, 2014). However, MetaErg's full genome  
75 annotation requires elevated processing times and computational resources due to its untargeted  
76 approach. Furthermore, there is a lack of a user-friendly tool to explore the results tables and graphs  
77 to extract pathway specific information tied to each MAG and thus infer potential species  
78 interactions based on their functional profiles.

79 Genome-based modelling approaches have routinely been used to study single organisms as well  
80 as microbial communities (Gottstein *et al*, 2016). For example, constraint-based models are highly  
81 employed in the study and prediction of metabolic networks (Heirendt *et al*, 2019). These models  
82 are generated upon the premise that any given function is feasible as long as the protein-encoding  
83 gene is present. Although species may lack the genetic potential to perform all functions necessary

84 to survive in a given ecosystem, in nature microbes do not exist in isolation and may benefit from  
85 their interaction with other species. By assessing the genomic content of individual species, we are  
86 able to identify groups of microbes whose combined content may account for complete ecosystem  
87 functioning. However, generating full genome metabolic networks for each species in a microbial  
88 community is time consuming and requires information not easily obtained for each community  
89 member such as biomass composition and nutritional requirements.

90 In order to decrease complexity and facilitate analysis, the search of interactions can be limited to  
91 groups of organisms (e.g. microbe-microbe or host-microbe) or specific ecosystem processes (e.g.  
92 nitrification or deadwood decomposition). A network-based tool for predicting metabolic  
93 capacities of microbial communities and interspecies interactions (NetMet) was recently  
94 developed by Tal et al., (Tal *et al.*, 2020). The tool only requires a list of species-specific enzyme  
95 identifiers and a list of compounds required for a given environment. However, besides the  
96 necessity of previous annotation of genomes, NetMet does not consider the rules that govern each  
97 reaction (e.g. protein complexes). Accurate annotation of gene function from sequencing data is  
98 essential to predict, ecosystem processes potentially performed by microbial communities,  
99 particularly in cases where an ecosystem process is performed by the synergy of two or more  
100 species. Simple methods for the annotation of genomes rely, for instance, on the search for  
101 homologous sequences. Computational tools such as BLAST (Altschul *et al.*, 1990) and  
102 DIAMOND (Buchfink *et al.*, 2015) allow the comparison of nucleotide or protein sequences to  
103 those present in databases. These approaches allow inferring the function of uncharacterized  
104 sequences from their homologous pairs whose function is already known. The degree of  
105 confidence in the assignment of biological function is increased if this has been validated by, for  
106 example, experimental data. Approaches based on orthology are increasingly used for genome-  
107 wide functional annotation (Huerta-Cepas *et al.*, 2017). Orthologs are homologous sequences that  
108 descend from the same ancestor separated after a speciation event retaining the same function  
109 (Koonin, 2005). OrthoMCL (Li *et al.*, 2003), CD-HIT (Li & Godzik, 2006) and OrthoFinder  
110 (Emms & Kelly, 2015, 2019) are just a few tools that identify homologous relationships between  
111 sequences using orthology. OrthoFinder has been shown to be more accurate than several other  
112 orthogroup inference methods since it considers gene length in the detection of ortholog groups by  
113 introducing a score transformation step (Emms & Kelly, 2015). However, OrthoFinder, due to its  
114 all-versus-all sequence alignment approach, requires intensive computational resources resulting  
115 in long running times when using large data sets for clustering. Because of the enormous number  
116 of potential combination, limiting the scope of research to specific ecosystem processes may  
117 reduce the computational and resource costs associated with the integration of ortholog clustering  
118 tools and functional annotation strategies.

119 In this study, we developed OrtSuite; a workflow that can: (i) perform accurate ortholog based  
120 functional annotation, (ii) reveal putative microbial synergistic interactions, and (iii) digest and  
121 present results for pathway and community driven biological questions. These different features  
122 can be achieved with the use of three bash commands in a reasonable computational time. This  
123 research question / hypothesis targeted approach integrates a user-defined database – Ortholog-  
124 Reaction Association database (ORAdb) – with up-to-date ortholog clustering tools. OrtSuite  
125 allows the search for putative microbial interactions by calculating the combined genomic

126 potential of individual species in specific ecosystem processes. OrtSuite also provides a visual  
127 representation of the species genetic potential mapped to each of the reactions defined by the user.  
128 We evaluate this workflow using a clearly defined set of reactions involved in the well-described  
129 benzoate-to-Acetyl-CoA (BTA) conversion. Further, we used this workflow to functionally  
130 characterize a set of known benzoate degraders. OrtSuite's ability to identify putative interspecies  
131 interactions was evaluated on species whose potential interactions have been previously predicted  
132 under controlled conditions (Fetzer *et al*, 2015).

133

## 134 **Results**

### 135 **Ortsuite is a flexible and user-friendly pipeline**

136 One of the motivations to develop Ortsuite was to facilitate the targeted analysis of the genomic  
137 potential of microbial communities including the prediction of putative synergistic interspecies  
138 interactions. To achieve this, OrtSuite was developed to integrate ortholog clustering tools (Emms  
139 & Kelly, 2019) with sequence alignment programs (Buchfink *et al*, 2015). To increase user-  
140 friendliness, three scripts were created that, collectively, perform all five tasks associated with  
141 OrtSuite: (1) download of sequences to populate ORAdb, (2) generation of Gene-Protein-Reaction  
142 (GPR) rules, (3) clustering of orthologs, (4) targeted functional annotation and (5) prediction of  
143 putative synergistic interspecies interactions (Figure 1). Additional control is also given to the user  
144 such as establishing thresholds in the minimum e-values (during sequence alignment of sequences  
145 in ortholog clusters to ORAdb). Other constraints include restricting the number of putative  
146 microbial interactions based on the presence of transporters and subsets of reactions to be  
147 performed by individual species (Supplementary data – Table S1). Since data in public repositories  
148 is frequently being added or updated and to include personal knowledge the user can manually  
149 curate the files in the ORAdb and GPR rules, with the latter being strongly advised.

150 A git repository for OrtSuite (<https://github.com/mdsufz/OrtSuite>) was also generated. This  
151 repository provides users with an easy-to-follow detailed guide covering installation to the running  
152 of the three scripts and generated outputs.

153

### 154 **Computing time of OrtSuite stages**

155 The runtime of each OrtSuite step was evaluated on a set of genomes whose genomic potential in  
156 the conversion of benzoate to acetyl-CoA was known (Table 1). The same set of genomes is used  
157 in the OrtSuite's GitHub tutorial page  
158 ([https://github.com/mdsufz/OrtSuite/blob/master/OrtSuite\\_tutorial.md](https://github.com/mdsufz/OrtSuite/blob/master/OrtSuite_tutorial.md)) with a total of 75.5  
159 Megabytes of data. OrtSuite was used to analyze this data on a laptop with 4 cores and 16  
160 Gigabytes of RAM. All OrtSuite steps were run on default settings, and the total runtime of each  
161 step was recorded (Table 2). The total workflow was completed in 3 h 50 min and the longest  
162 single step runtime consisted of 2 h and 47 min which involved the construction of the ORAdb.  
163 The user does have the option to modify the number of cores used during functional annotation  
164 which should further decrease run times.

165

## 166 **Higher recall rates during clustering of orthologs with DIAMOND**

167 We performed an evaluation of the effects of point mutations during clustering of orthologs using  
168 OrthoFinder (Emms & Kelly, 2019). OrthoFinder allows users to choose between DIAMOND  
169 (Buchfink *et al*, 2015) and BLAST (Altschul *et al*, 1990) as sequence aligners. To test which  
170 sequence aligner yielded the best results we performed ortholog clustering of a dataset consisting  
171 of the original target genomes as well as a set of artificially mutated genomes (Supplementary data  
172 - Test\_genome\_set) using both aligners. The results showed 0.01 difference between OrthoFinder  
173 and DIAMOND precision (Table 2). However, DIAMOND showed a 9.5% higher recall than that  
174 observed for OrthoFinder what suggests DIAMOND may have higher sensitivity in the clustering  
175 of sequences with the same function. All artificially mutated sequences (even those with mutation  
176 rates of 25%) were clustered together with their non-mutated ortholog. In parallel, we also  
177 performed sequence alignment using NCBI's BLASTp (Madden, 2003) between the protein  
178 sequences of the DNA-mutated and un-mutated genes. E-values of sequence alignments in all  
179 species ranged from 0 to  $5e^{-180}$  and percentage of identity from 61.32 to 98.84% (Supplementary  
180 data – Table S2). For validation of the OrtSuite workflow, clustering of protein orthologs was  
181 repeated using only the original unmutated 18 genomes and the default aligner (DIAMOND). A  
182 complete overview of the results generated during the clustering of orthologs (e.g. number of genes  
183 in ortholog clusters, number of unassigned genes and number of ortholog clusters) was also  
184 obtained (Supplementary data - Table S3).

185

## 186 **High rate of KEGG annotations predicted by OrtSuite**

187 The third step of OrtSuite consists of performing cluster annotation in a two-stage process. In the  
188 first, only 50% of sequences are used in the alignment to the sequences ORAdb. Those with a  
189 minimum e-value proceed to the second stage where all sequences contained in this cluster will be  
190 aligned. At the end, annotation of clusters will take into consideration additional parameters such  
191 as bit scores. To evaluate the thresholds used in the annotation of ortholog clusters we used one  
192 relaxed (0.001) and four restrictive ( $1e^{-4}$ ,  $1e^{-6}$ ,  $1e^{-9}$  and  $1e^{-16}$ ) e-value cutoffs. An overview of the  
193 results (e.g. number of clusters containing orthologs from ORAdb, number of ortholog clusters  
194 with annotated sequences) is shown in (Supplementary data – Table S4). The performance of  
195 OrtSuite in the functional annotation of the genomes in the *Test\_genome\_set* is shown in  
196 (Supplementary data – Table S5). On average, 96% of the annotations assigned by KEGG were  
197 also identified by OrtSuite. The complete list of results of functional annotation using the different  
198 e-value cutoffs are available in the Supplementary data - Table S6, Table S7, Table S8 and Table  
199 S9. Similarly, the mapping of species with the genetic potential for each reaction (considering the  
200 GPR rules) using the different e-value cutoffs can be found in the Supplementary data – Table  
201 S10, Table S11, Table S12 and Table S13. In terms of annotation, no striking difference was  
202 observed between the four different e-value cutoffs used during the restrictive search stage.  
203 However, the largest decrease in the number of ortholog clusters that transits from the relaxed  
204 search to the restrictive occurs when using an e-value cutoff of  $1e^{-16}$  (Supplementary data – Table

205 S4). The difference in computing time between lower and higher e-value thresholds was negligible  
206 (< 2 min). Other annotation tools, such as NCBI's BLAST tool (Altschul *et al*, 1990),  
207 BlastKOALA (Kanehisa *et al*, 2016) and Prokka (Seemann, 2014), can annotate full genomes, the  
208 latter at a relatively fast pace. On average, full genome annotation of our genomes in the  
209 *Test\_genome\_set* dataset using Prokka required 12 mins on a customary laptop with 16 Gigabytes  
210 of RAM and four CPUs to complete. BlastKOALA required approximately 3 hours to annotate a  
211 single genome.

212

### 213 **Identifying genetic potential to perform a pathway**

214 To test OrtSuite's ability to identify species with the genetic potential to perform a pathway  
215 individually we defined sets of reactions that are used in three alternative pathways for the  
216 conversion of benzoate to acetyl-CoA (Supplementary data – Table S14). Next, we compared the  
217 results to the species' known genomic content in each alternative pathway (Supplementary data –  
218 Table S15). OrtSuite matched KEGG's predictions in species' ability to perform each alternative  
219 benzoate degradation pathway in all but two species - *Azoarcus* sp. DN11 and *Thauera* sp. MZ1T.  
220 Furthermore, OrtSuite identified five species capable of performing conversion pathways not  
221 contemplated in KEGG. *Azoarcus* sp. KH32C, *Aromatoleum aromaticum* EbN1,  
222 *Magnetospirillum* sp. XM-1 and *Sulfuritalea hydrogenivorans* sk43H have the genetic potential to  
223 perform both pathways involving the anaerobic conversion of benzoate to acetyl-CoA while  
224 *Azoarcus* sp. CIB has to genetic potential to perform all alternative pathways (except when using  
225 an e-value cutoff of  $1e^{-16}$ ). No genes in *Thauera* sp. MZ1T involved in the conversion of crotonyl-  
226 CoA to 3-Hydroxybutanoyl-CoA (R03026) were identified by OrtSuite which impedes the  
227 anaerobic conversion of benzoate to acetyl-CoA. The default e-value for the restrictive search was  
228 set to  $1e^{-9}$  since OrtSuite's performance did not change significantly between all tested e-value  
229 cutoffs but showed a greater drop in the number of consistent orthogroups (i.e. clusters of orthologs  
230 whose sequences are all annotated with the same function) from  $1e^{-9}$  to  $1e^{-16}$ .

231

### 232 **Using OrtSuite to predict interspecies interactions**

233 In this study, we tested the ability of OrtSuite in identifying interspecies interactions involved in  
234 the conversion of benzoate to acetyl-CoA where experimental data were available. Prediction of  
235 synergistic interspecies interactions was assessed on a set of sequenced isolates (Supplementary  
236 data - Fetzer\_genome\_set.zip): Monocultures of these isolates and randomly assembled  
237 communities of one to 12 species including these benzoate-degraders and additional species  
238 incapable of directly using benzoate as a carbon source were analyzed previously (Fetzer *et al*,  
239 2015) under three different environmental conditions (low substrate concentration: 1g/L benzoate,  
240 high substrate concentration: 6g/L benzoate and high substrate concentration + additional osmotic  
241 stress: 6g/L benzoate supplemented with 15g /L of NaCl). In that study, Fetzer et al investigated  
242 if the presence or absence of a particular species positively or negatively affected biomass  
243 production. Since under specific conditions the presence of a degrader alone was not sufficient for  
244 community biomass production, they further analyzed if potential species interactions could be of

245 relevance. Briefly, they defined for all environmental conditions minimal communities, which  
246 showed community growth without the need to include other species and identified whether the  
247 presence of a single species alone or potential interaction between the specific species (and thus  
248 potential partners) present in these minimal communities stimulated biomass production (Fetzer  
249 *et al*, 2015). Using OrtSuite, we aimed to identify which potential species interactions predicted  
250 by Fetzer and collaborators could be a result of their combined genetic potential.

251 Our dataset contained 69,193 protein sequences distributed across the 12 species resulting in a  
252 total of 59 Megabytes of data. More than 84% of all genes were placed in 9,533 ortholog clusters.  
253 In addition, 541 clusters were composed of sequences obtained from all 12 species (Supplementary  
254 data - Table S16). OrtSuite's annotation stage resulted in 326 ortholog clusters with annotated  
255 sequences from ORAdb (Supplementary data - Table S17). The mapping of KOs to each species  
256 in the *Fetzer\_genome\_set* is available as supplementary data (Table S18). The genomic potential  
257 of each species for aerobic and anaerobic benzoate metabolizing pathways is shown in Figure 2.  
258 The complete mapping of reactions to each species is available in the supplementary data (Table  
259 S19). Based on the 326 ortholog clusters and the Gene-Protein-Reaction (GPR) rules  
260 (Supplementary data - Table S20), five species (*Cupriavidus necator* JMP134, *Pseudomonas*  
261 *putida* ATCC17514, *Rhodococcus sp.* Isolate UFZ (Umweltforschung Zentrum), *Rhodococcus*  
262 *ruber* BU3 and *Sphingobium yanoikuyae* DSM6900) contained all protein-encoding genes  
263 required to perform aerobic conversion of benzoate to acetyl-CoA. In the Fetzer study,  
264 *Rhodococcus sp.* Isolate UFZ and *S. yanoikuyae* did not show growth in a medium containing  
265 benzoate. The incomplete functional potential of *C. testosteroni* ATCC 17713 and *P. putida*  
266 ATCC17514 to perform aerobic conversion of benzoate to acetyl-CoA is at odds with their  
267 reported growth as monocultures in the presence of benzoate as shown in the Fetzer study. The  
268 number of species with the genetic potential for each reaction involved in the aerobic benzoate  
269 degradation pathway (P3) is shown in (Supplementary data - Table S21). Usually, all species with  
270 the complete genomic potential to perform a complete pathway are excluded when calculating  
271 interspecies interactions since they do not require the presence of others. However, species  
272 identified by OrtSuite with the complete functional potential to perform each defined pathway  
273 were also included to compare to the results in the Fetzer study presented above. A total of 2382  
274 combinations of species interactions were obtained whose combined genetic potential covered all  
275 reactions. The complete list of potentially interacting species is available in the supplementary data  
276 (Table S22).

277 In the anaerobic degradation pathways (P1 and P2) no species presented the genomic content to  
278 encode proteins involved in the conversion of benzoyl-CoA to Cyclohexa-1,5-diene-1-carboxyl-  
279 CoA (R02451) (Supplementary data - Table S23). This reaction requires the presence of a protein  
280 complex either composed of four subunits (K04112, K04113, K04114, K04115) or composed of  
281 two subunits (K19515, K19516). No species was annotated with all subunits in either protein  
282 complex. Therefore, no species interactions were identified that would allow the complete  
283 anaerobic conversion of benzoate to acetyl-CoA. In the low substrate environment, OrtSuite  
284 identified 826 of 830 (99.5%) species combinations showing growth. In the high substrate  
285 environment, OrtSuite predicted 644 of 646 (99.7%). In the high substrate+salt stress environment,

286 OrtSuite predicted all 271 (100%) combinations of species exhibiting growth (Supplementary data  
287 - Table S24).

## 288 **Discussion**

289 We designed OrtSuite to allow hypothesis-driven exploration of microbial interactions in a user-  
290 friendly manner. This was achieved by integrating up-to-date clustering tools with faster sequence  
291 alignment methods and limiting the scope to user-defined ecosystem processes or metabolic  
292 functions. Using only three bash commands required to run the complete workflow, OrtSuite is a  
293 user-friendly tool capable of running in a customary computer (four cores and 16GB of RAM)  
294 with even faster runtimes when using high performance computing.

295 The clustering of orthologs by OrthoFinder using DIAMOND (Buchfink *et al*, 2015) showed  
296 higher sensitivity and lower runtime compared to BLAST (Altschul *et al*, 1990) which has also  
297 been shown by Hernández-Salmerón and Moreno-Hagelsieb (Hernández-Salmerón & Moreno-  
298 Hagelsieb, 2020). Furthermore, low e-values and medium to high identity percentages in the  
299 sequence alignments between mutated and original genes indicates that the mutated genes still  
300 share enough sequence similarity to the original protein sequence. These results suggest that  
301 mutation rates of up to 25% of single DNA base pairs will not have an observable effect on the  
302 clustering of orthologs. OrthoFinder's algorithm removes the gene length bias from the sequence  
303 alignment process, which may also explain why mutated genes were clustered with the original.  
304 Although it has been suggested that most genetic variations are neutral, changes in single base  
305 pairs can have a drastic effect on protein function (e.g. depending on the location of the mutation)  
306 (Ng & Henikoff, 2006). To this purpose, experimental functional studies can be used to validate  
307 previously unannotated orthologs. Furthermore, this study case does not consider the distribution  
308 of mutations across species and gene families which can also have different effects on the  
309 clustering of orthologs (Khanal *et al*, 2015). Therefore, future studies increasing the rates of DNA  
310 base pair substitutions and other types of mutations as well as experiments targeting protein  
311 function in ortholog clusters are needed.

312 Next, we aimed to improve and facilitate functional annotation and prediction of synergistic  
313 microbial interactions. Exploring the great amount of data generated from full genome annotation  
314 of individual species from complex microbial communities is a daunting task. This is evident in a  
315 study by Singleton and collaborators (Singleton *et al*, 2021) where the connection between  
316 structure and function required the analysis of metagenomics data, 16S and molecular techniques  
317 such as fluorescent in situ hybridization and Raman spectroscopy. Looking solely to functional  
318 annotation, two challenges, among others, arise: First, performing all vs all sequence alignments  
319 in complex communities is resource consuming (time and computational power) and second,  
320 manual inspection of each annotated genome for target genes or pathways is required. Identifying  
321 interspecies interactions based on the microbe's complete genomic potential is also challenging.  
322 For example, network approaches are increasingly employed in ecology but selection of the most  
323 appropriate approach is not always straightforward and easy to implement (Delmas *et al*, 2019).  
324 OrtSuite overcomes these challenges by first performing cluster annotation in a two-stage process  
325 and limited to user-defined set of functions of interest which decreases the number of sequence  
326 alignments to be performed. The user-defined database coupled with the scripts for automated



327 identification of interspecies interactions contained in OrtSuite decreases the time required to  
328 generate the data and facilitates its interpretation by the user. Additionally, OrtSuite generates a  
329 graphical representation of the network further facilitating analysis of the whole microbial  
330 community ([https://github.com/mdsufz/OrtSuite/blob/master/network\\_example.png](https://github.com/mdsufz/OrtSuite/blob/master/network_example.png)).

331 OrtSuite not only confirmed all but two of KEGG's predictions in species' ability to perform each  
332 alternative benzoate degradation pathway used in this study but also identified five species capable  
333 of performing conversion pathways not contemplated in KEGG. On average, an additional 18.3  
334 KO identifiers were mapped to genes not previously annotated in our species. The use of e-value  
335 and bit score as the filtering criteria rather than sequence identity, employed by KEGG, may  
336 explain the increase in functionally annotated genes. For example, the alignment of a sequence of  
337 *A. defluvii* (adv: AWL30228.1) to the sequences in ORAdb annotated as K04105 (conversion of  
338 benzoate to benzoyl-CoA) showed high bit-scores (200.7) and low e-values ( $2e^{-54}$ ) but the identity  
339 percentage did not exceed 28.6%. The use of e-values and bit scores to infer function has been  
340 nicely reviewed by Pearson (Pearson, 2013) who suggests that e-values and bit scores are more  
341 sensitive and reliable than identity percentages in finding homology since they take into account  
342 evolutionary distance of aligned sequences, the sequence lengths and the scoring matrix.

343 To test the prediction of putative synergistic microbial interactions we used data from an  
344 independent study performed by Fetzer and collaborators (Fetzer *et al*, 2015); hereafter Fetzer  
345 study. In the Fetzer study five species showed biomass growth (estimated by optical density at  
346 590nm wave length) in medium containing benzoate. We evaluated whether these species  
347 possessed the complete genomic content to encode all proteins required for each benzoate to  
348 acetyl-CoA conversion pathway. The remaining seven species were not able to grow as  
349 monocultures in media with benzoate as sole carbon source. Therefore, we evaluated whether the  
350 lack of growth was confirmed by lack of essential protein-encoding genes involved in conversion  
351 of benzoate to acetyl-CoA. Fetzer study also showed that, under specific nutrient and stress  
352 conditions, total biomass production was influenced by the presence of non-degrading species.  
353 Thus, we evaluated whether putative species interactions identified by OrtSuite fit the results  
354 obtained by in the Fetzer study. OrtSuite confirmed the functional potential for aerobic conversion  
355 of benzoate to acetyl-CoA in three of the five species whose growth in monocultures was observed  
356 during their study. In Fetzer's study, two species, *S. yanoikuyae* (accession number  
357 GCA\_903797735.1) and *Rhodococcus sp.* (accession number GCA\_903819475.1), were not able  
358 to grow as monoculture in the presence of benzoate. However, OrtSuite predicted that both  
359 possessed the functional potential to aerobically convert benzoate to acetyl-CoA. In their study, in  
360 a medium containing 1g/L of benzoate, growth was considered when optical densities (OD) were  
361 above 0.094. The OD measured for *S. yanoikuyae* was 0.0916. The annotation of genes with the  
362 ability to perform the complete aerobic conversion of benzoate to acetyl-CoA combined with the  
363 small difference in OD to the minimum threshold suggests that *S. yanoikuyae* indeed can grow on  
364 low benzoate containing medium but at perhaps at lower growth rates. In the case of *Rhodococcus*  
365 *sp.* Isolate UFZ, the OD was never measured above 0.022 which, again, might indicate slow  
366 growing species. Another possible explanation is that although these two species possess the genes  
367 necessary for aerobic benzoate degradation they are not active. In Fetzer's study, the observed  
368 growth of *Comamonas testosteroni* ATCC11996 and *Pseudomonas fluorescens* DSM6290 in the

369 low benzoate environment was not confirmed by OrtSuite. To note, benzoate conversion  
370 intermediates were not determined in the Fetzer experiment. Hence, it is possible that these two  
371 species utilize reactions or pathways that were not included in the benzoate degradation pathways  
372 used in our study. Despite the presence of benzoate degraders, another possible explanation as to  
373 the unobserved growth in Fetzer's study for certain experimental conditions is the lack of tolerance  
374 of these species to high benzoate concentrations. For example, *C. necator* growth was shown to  
375 be stimulated at low benzoic acid concentrations but inhibited at high concentrations (Wang *et al*,  
376 2014). In addition, the set of genes used in our study did not consider the presence of stress related  
377 factors. To assess these effects, stress-resistance associated genes and reactions such as those  
378 involved in medium acidification (Kitko *et al*, 2009) could be added as constraints. Similar results  
379 were obtained when using a high substrate+salt stress medium. Under these conditions, presence  
380 of benzoate degraders alone was not sufficient to achieve growth of species combinations.  
381 Benzoate degradation has been shown to decrease in hyperosmotic environments (Bazire *et al*,  
382 2007) therefore, additional constraints such as genes that confer resistance to environmental  
383 stressors or adverse conditions sodium chloride (NaCl) could be included during the identification  
384 of interspecies interactions under different or changing environmental conditions.

385 No single species or combination of species possessed the complete genomic potential to  
386 anaerobically convert benzoate to acetyl-CoA via the two proposed pathways (P1 and P2). Since  
387 all growth experiments were conducted in aerobic conditions, it is possible that the species in  
388 question are only capable of using benzoate as a carbon source in aerobic environments. To fully  
389 explore all the species potential to convert benzoate, additional degradation pathways could be  
390 checked in the future using a multi-omics approach. Furthermore, the only constraints  
391 added were related to the reactions that composed each pathway. Additional constraints can be  
392 included in future studies, such as potential mandatory transport-associated reactions, to increase  
393 confidence in the proposed interspecies interactions. OrtSuite confirmed that most interspecies  
394 interactions (> 99%) identified by Fetzer and collaborators were possible due to their combined  
395 metabolic potential to aerobically degrade benzoate to acetyl-CoA but not under anoxic conditions.

396 In this study, we ran OrtSuite on a dataset comprised of 18 genomes (Table 1). To determine if  
397 this range would be within the number of genomes in regular microbiome studies we calculated  
398 the average number of MAGs from different studies focusing on their recovery. A study performed  
399 by Parks and collaborators (Parks *et al*, 2017) analyzed sequencing data from 149 projects. Most  
400 projects (91%) consisted of less than 20 samples. On average, they recovered 5.3 metagenome-  
401 assembled genomes (MAGs) per metagenome. Work performed by Pasolli and collaborators  
402 (Pasolli *et al*, 2019) on microbial diversity in the human microbiome recovered, on average, 16  
403 MAGs per metagenomic library. From the 46 studies used in their work, 30 consisted of less than  
404 200 samples. Another study by Tully and collaborators focusing on marine environments (Tully  
405 *et al*, 2018) recovered 2631 MAGs from 234 samples (average of 11 MAGs per sample). Our  
406 analysis demonstrates that the average number of MAGs recovered from a metagenome currently  
407 range from five to 16. Therefore, performing targeted functional annotation and interspecies  
408 interactions predictions using OrtSuite in average sized metagenome samples is still feasible using  
409 a customary laptop.

410 In summary, OrtSuite allows hypothesis-driven exploration of potential interactions between  
411 microbial genomes by limiting the search universe to a user-defined set of ecosystem processes.  
412 This is achieved by rapidly assessing the genetic potential of a microbial community for a given  
413 set of reactions considering the relationships between genes and proteins. The two-step annotation  
414 of clusters of orthologs with a personalized ORAdb decreases the overall number of sequence  
415 alignments that need to be computed. User-specified constraints, such as the presence of  
416 transporter genes, further reduces the search space for putative microbial interactions. Users have  
417 substantial control over several steps of OrtSuite: from manual curation of ORAdb, custom  
418 sequence similarity cutoffs to the addition of constraints for inference of putative microbial  
419 interactions. The reduction of the search space of synergistic interactions by OrtSuite will also  
420 allow more comprehensive and computationally demanding tasks to be performed such as  
421 (Community) Flux Balance Analysis which depend heavily on genome-scale metabolic models  
422 (Thommes *et al*, 2019; Ravikrishnan & Raman, 2021). As long as links between genes, proteins  
423 and reactions exist, the flexibility and easy usage of OrtSuite allow its application to the study of  
424 any given ecosystem process.

425

## 426 **Materials and Methods**

### 427 **OrtSuite workflow**

428 The OrtSuite workflow consists of three main steps performed by the use of three bash commands  
429 (Figure 1). Briefly, the first step consists in the generation of a user defined ortholog-reaction  
430 associated database (ORAdb) and collection of the gene-protein-reaction (GPR) rules. This task  
431 takes as input a list of KEGG identifiers which will be used to download all protein sequences  
432 associated with a set of reactions/pathway of interest. Next, all gene-protein-rules (GPRs)  
433 associated with each reaction will be downloaded from KEGG Modules. In the second step  
434 OrtSuite employs OrthoFinder (Emms & Kelly, 2015) to generate ortholog clusters. This step  
435 takes as input a folder with the location of the genomic sequences. The third step consists of the  
436 functional annotation of species, identification of putative synergistic interspecies interactions and  
437 generation of visual representations of the results.

438

### 439 **OrtSuite step 1 (green box, Figure 1) – User defined Ortholog-Reaction Association database 440 (ORAdb) and Gene-Protein-Reaction (GPR) rules file**

441 The ORAdb used for functional annotation consists of sets of protein sequences involved in the  
442 enzymatic reactions that compose a pathway/function of interest defined by the user. This database  
443 is generated during the execution of the *DB\_construction.sh* script in OrtSuite requiring only the  
444 user to provide:

- 445 • a location of the project folder where all results will be stored
- 446 • a text file with a list of KEGG identifiers (one identifier per line)
- 447 • the full path to the OrtSuite installation folder

448 The list of identifiers can be KEGG reactions (RID) (e.g. R11353, R02451), enzyme commission  
449 (EC) numbers (e.g. 1.3.7.8, 4.1.1.103) or KEGG ortholog identifiers (e.g. K07539, K20941). This  
450 file is used by OrtSuite to automatically retrieve the KEGG Ortholog identifiers (KO) (in case the  
451 identifiers provided are not KO identifiers) and to download all their associated protein sequences  
452 (Kanehisa *et al*, 2004). OrtSuite makes use of the python library *grequests* which allows multiple  
453 queries in KEGG subsequently decreasing the time required for retrieving the ortholog associated  
454 sequences. The user-defined ORAdb will be composed of KO-specific sequence files in FASTA  
455 format associated with all reactions/enzymes of interest. Users also have the opportunity to  
456 manually add or edit the sets of reactions and the associated protein sequences in the ORAdb. This  
457 feature is of particular importance since many reactions associated with ecosystem processes are  
458 constantly being discovered and updated and might not be included in the latest version of KEGG.  
459 In addition, during the execution of the *DB\_construction.sh* OrtSuite performs the automated  
460 download of the gene-protein-reaction (GPR) rules from KEGG Modules. This feature is vital  
461 since many reactions can be catalyzed by enzymes with a single (i.e., one protein) or multiple  
462 subunits (i.e., protein complexes). Despite the automated process, it is strongly advised to  
463 manually curate the final table to guarantee accurate results. An example of the final GPR table is  
464 shown in the Supplementary data (Table S20).

465

## 466 **OrtSuite step 2 (purple box, Figure 1) - Generation of protein ortholog clusters**

467 The second step of OrtSuite, takes a set of protein sequences and generates clusters of orthologs.  
468 This set of protein sequences can originate from single isolates or from the complete set of protein  
469 sequences recovered from metagenomes or metagenome-assembled genomes. Indeed, the use of  
470 protein sequences from isolates, metagenome-assembled genomes and co-culture experiments will  
471 benefit greatly from OrtSuite's reduction of the universe of potential microbial interactions based  
472 on the user defined ORAdb. Orthology considers that phylogenetically distinct species can share  
473 functional similarities based on a common ancestor (Gabaldón & Koonin, 2013). Potentially, genes  
474 with equal function will be grouped together. To perform this task the OrtSuite pipeline uses  
475 OrthoFinder (Emms & Kelly, 2015). Two sequence aligners are available in OrthoFinder –  
476 DIAMOND (Buchfink *et al*, 2015) and BLAST (Altschul *et al*, 1990). DIAMOND is used by  
477 default due to its improved trade-off between execution time and sensitivity (Emms & Kelly,  
478 2019). This step is performed by running the command *orthofinder* located in the installation folder  
479 of OrthoFinder. This command takes as input the full path to the folder containing the protein  
480 sequences to be clustered and the full path to the folder where results are to be stored.

481

## 482 **OrtSuite step 3 (yellow box, Figure 1) - Functional annotation of ortholog clusters**

483 The third step of OrtSuite consists in the assignment of functions to protein sequences contained  
484 in the ortholog clusters. Functional annotation of these clusters consists of a two-step process  
485 termed relaxed and restrictive search, respectively. The goal of the relaxed search is to decrease  
486 the number of alignments required to assign functions to sequences in the ortholog clusters. Here,  
487 50% of the total number of sequences from each cluster are randomly selected and aligned to all

488 sequences associated to each reaction present in the ORAdb. Only the e-value is considered during  
489 this stage. Ortholog clusters where e-values meet a user-defined threshold to sequences in the  
490 ORAdb proceed to the restrictive search. The default e-value was set to 0.001, as the main  
491 objective of the relaxed search is to capture as many sequences for annotation as possible while  
492 avoiding an exaggerated number of sequence alignments. In the restrictive search, all sequences  
493 in the transitioned ortholog clusters are aligned to all the sequences in the reaction set(s) present  
494 in the ORAdb to which they had a hit during the relaxed search. Again, the query sequence is only  
495 assigned to the function of a reference sequence if the e-value is below a determined threshold  
496 (default  $1e^{-9}$ ). Next, an additional filter is applied based on annotation bit score values (default 50).  
497 Although we established default values for the relaxed and restrictive search as well as bit score,  
498 the user has the option to define the thresholds for all individual parameters.

499 The identification of putative interactions between species is based on all combinations of bacterial  
500 isolates with the genomic content to perform the user-defined pathway defined in the ORAdb. The  
501 input for this task consists of: (1) a binary table generated at the end of the functional annotation  
502 which indicates the presence or absence of sequences annotated to each reaction in the ORAdb in  
503 each species (e.g. Supplementary Table S10); (2) a set of Gene-Protein-Reaction (GPR) rules for  
504 all reactions considered (e.g. Supplementary data - Table S20); and (3) a user-defined tab-  
505 delimited file where the sets of reactions for complete pathways, subsets of reactions required to  
506 be performed by single species and transporter-associated genes (e.g. Supplementary data – Table  
507 S1) are described. To further reduce the vast amount of putative microbial interactions and to  
508 increase confidence in the results manual filtering can be performed to reflect available knowledge  
509 (e.g. known cross-feeding relationship between species) and/or the likelihood of biologically  
510 feasible species interactions). The user also may have interest in assessing subsets of microbial  
511 interactions using specific criteria. Therefore, additional constraints can be applied to the list of  
512 putative microbial interactions further reducing the search space. These include the degree of  
513 completeness of a pathway, the number of reactions expected to be performed by a single species  
514 or the presence or absence of transporter genes. Additionally, a graphical network visualization is  
515 also produced during this step. Graphical network visualization is implemented in R using the  
516 packages visNetwork (v2.0.9), reshape2 (v1.4.3) and RColorBrewers (v1.1-2) but also requires the  
517 pandoc linux library. Graphical visualization was implemented with R v3.6 but tested also with  
518 v4.0. The visualization creates a HTML file that allows interactive exploration of the network and  
519 provides hyperlinks to KEGG if available.

520 All tasks - functional annotation, prediction of putative microbial interactions and generation of  
521 graphical visualizations - are performed by running the script *annotate\_and\_predict.sh* included  
522 in OrtSuite ([https://github.com/mdsufz/OrtSuite/blob/master/annotate\\_and\\_predict.sh](https://github.com/mdsufz/OrtSuite/blob/master/annotate_and_predict.sh)). OrtSuite's  
523 predictions of individual species and combinations of species with the genetic potential to perform  
524 each defined pathway is stored in text files located in a folder termed “interactions”.

525

526 **Conversion of benzoate to acetyl-CoA as a model pathway**

527 We selected three alternative pathways involved in the conversion of benzoate to acetyl-CoA  
528 (BTA) to test the functional annotation and prediction of putative synergistic microbial interactions  
529 using OrtSuite (Supplementary data - Table S14). Two pathways consisted in the anaerobic  
530 degradation of benzoate to acetyl-CoA via benzoyl-CoA differing only in the reactions required  
531 for transformation of glutaryl-CoA to crotonyl-CoA (hereafter, respectively, P1 and P2). P1 first  
532 converts glutaryl-CoA to glutaconyl-CoA and then to crotonoyl-CoA while P2 directly converts  
533 glutaryl-CoA to crotonoyl-CoA. One pathway consisted in the aerobic degradation of benzoate via  
534 catechol (hereafter P3). The complete number of reactions, enzymes, KO identifiers and KO-  
535 associated sequences in each alternative pathway is shown in the supplementary data  
536 (Supplementary data - Table S25).

537

### 538 **Species selection for testing functional annotation**

539 To assess the performance of OrtSuite, we selected the transformation of benzoate to acetyl-CoA  
540 as a model pathway and a set of previously characterized species known to be involved in this  
541 pathway (Table 1). This set of species was divided in two groups. The first group contained  
542 sequenced genomes of species whose ability to convert benzoate to acetyl-CoA has been  
543 demonstrated by KEGG (Kanehisa *et al*, 2004) and were selected as positive controls. These  
544 species were classified according to their genomic potential: complete, if all protein encoding  
545 genes required for a BTA pathway were present in their genome or partial, if not all protein  
546 encoding genes were present. The second group consisted of species known to lack all required  
547 protein encoding genes and were selected as negative controls. In total, we selected 18 species as  
548 positive controls. Seven of them have the genetic potential to perform the alternative P2 pathway;  
549 eight have the genetic potential to perform alternative path P3 (positive controls); and, none able  
550 to completely perform the alternative path P1. To note that species *Thauera sp.* MZ1T has the  
551 genetic potential to perform P2 and P3 pathways. Four organisms were selected as negative  
552 controls. Using their genomes, we evaluated the performance of OrtSuite based on precision and  
553 recall rates for clustering of orthologs and the correct functional annotation of sequences. Also, a  
554 set of genomes from the species containing the genetic potential to degrade benzoate were  
555 artificially mutated at the nucleotide level at different rates in order to determine how levels of  
556 point mutations in open reading frames (ORFs) affected clustering of ortholog groups.

557

### 558 **Species selection for validation of putative interspecies interactions**

559 In a study performed by Fetzer and collaborators (Fetzer *et al*, 2015) community biomass  
560 production of mono- and mixed-cultures was assessed in medium containing benzoate. The authors  
561 used this data to infer potential species interactions. This set of genomes was processed with  
562 OrtSuite to determine the species' genetic potential to degrade benzoate, either individually or as  
563 a result of their interaction. Our results were compared to those obtained by Fetzer and  
564 collaborators and used to assess whether the study's inferred potential interactions could be derived  
565 from their combined genetic potential.

566

## 567 **Evaluation of ortholog clustering**

568 The clustering of orthologs was evaluated by measuring the pairwise precision and recall.  
569 Clustering precision measures how many pairs of sequences associated with the same molecular  
570 function are grouped together and is calculated by dividing the number of correctly clustered  
571 sequences by the total number of clustered sequences (Equation 1).

572

$$\text{Clustering precision} = \text{correctly clustered sequences} / \text{total number of clustered sequences} \quad (1)$$

573

574 where, correctly clustered sequences refers to the pairs of sequences that share the same function  
575 and are clustered together and total number of clustered sequences refers to all pairs of sequences  
576 that are clustered together irrespective of sharing the same function.

577

578 Clustering recall measures how many pairs of sequences with the same molecular function are not  
579 clustered together. Recall is calculated by dividing the number of correctly clustered sequences by  
580 the total true sequence clusters (Equation 2).

581

$$\text{Clustering recall} = \text{correctly clustered sequences} / \text{total true sequence clusters} \quad (2)$$

582

583 where, correctly clustered sequences refers to the pairs of sequences that share the same function  
584 and are clustered together and total true sequence clusters refers to all pairs of sequences that have  
585 the same function.

586

## 587 **Evaluation of sequence aligner used for clustering of orthologs**

588 Changes of a single DNA base can result in the production of a different amino acid which might  
589 result in a different protein. To determine the impact of mutations on the clustering of orthologs a  
590 single gene from three species was artificially mutated at different rates. These mutations were  
591 introduced in the nucleotide sequences of each gene. Only substitutions were considered since  
592 these are the most commonly studied (Lynch, 2010) and none of the mutations were allowed to  
593 occur on the first and last codon. When, during the mutation, new stop or/and start codons were  
594 introduced, the translation was made for all the possible proteins and the largest was selected.

595 *Burkholderia vietnamiensis* G4 was mutated on the gene K05783, *Azoarcus sp.* CIB on the gene  
596 K07537 and *Aromatoleum aromaticum* EbN1 on the gene K07538. Each gene was mutated at rates  
597 of 0.01, 0.03, 0.05, 0.1, 0.15 and 0.25. Each mutation rate resulted in an in silico strain of the

598 original genome (e.g., *Burkholderia vietnamiensis* G4 strain K05783\_25, where “K05783” is the  
599 KEGG ortholog identifier and “25” is the rate of mutation). A total of 18 strains were generated  
600 (six in silico mutated strains per genome). The complete set of original and artificially mutated  
601 genomes is available in a compressed file (Supplementary data - Test\_genomes\_set.zip).

602

### 603 **Evaluation of functional annotation**

604 Functional annotation was evaluated based on the data collected from KEGG (Altschul *et al*,  
605 1990). Annotation performance is calculated by dividing the number of matching annotated  
606 sequences by the total number of annotations (Equation 3).

607

$$\text{Annotation performance} = \text{matching annotated sequences} / \text{total number of annotations} \quad (3)$$

608

609 where, matching annotated sequences refers to the number of sequences annotated by KEGG  
610 annotations predicted by OrtSuite and total number of annotations refers to the all sequences that  
611 were assigned a function by KEGG.

612

### 613 **Evaluation of microbial interaction predictions**

614 We evaluated the prediction of putative microbial interactions using a genome set from an  
615 independent study (Fetzer *et al*, 2015) containing species with exhibited growth in medium  
616 containing benzoate (defined as Fetzer\_genome\_set). The authors do not identify specific potential  
617 interactions in the transformation of benzoate but infer interspecific interactions in an environment  
618 containing benzoate as the major carbon source. For the complete set of species combinations and  
619 benzoate degradation capabilities and effects identified by Fetzer and collaborators, see (Fetzer *et*  
620 *al*, 2015) (Supplementary data - Table S24).

#### 621 *Bacterial cultures and sequencing*

622 Bacterial cryo-cultures of the different isolates were revived on LB agar plates. Single colonies  
623 were picked and grown overnight in 2 ml LB medium at 37°C. The cells were pelleted by  
624 centrifugation. Cells were lysed and genomic DNA was extracted using a Nucleospin Tissue Kit  
625 (Machery and Nagel). Approximately 150 to 1000 ng of DNA were used for fragmentation (insert  
626 size: 300 – 700 bp) and sequencing libraries were prepared following the NEB Ultra II FS Kit  
627 protocol (New England Biolabs). Libraries were quantified using a JetSeq Library Quantification  
628 Lo-ROX Kit (BioLine) and quality-checked by Bioanalyzer (Agilent). These libraries were  
629 sequenced on an Illumina MiSeq instrument with a final concentration of 8 pM using the v3 600  
630 cycles chemistry and 5% PhiX.

#### 631 *Genome assembly and Open Reading Frame prediction*



632 The sequenced reads were quality checked using Trim Galore v0.4.4\_dev. Next, genomes were  
633 assembled using the Spades Assembler v3.15.2 and their quality assessed using CheckM.  
634 Taxonomic classification was performed using Genome Taxonomy Database (GTDBTk) release  
635 95. Open Reading Frames (ORFs) were predicted using Prodigal v2.6.3. Translation of sequences  
636 to amino acid format was performed using faTrans from kentUtils ([https://github.com/ENCODE-](https://github.com/ENCODE-DCC/kentUtils/tree/master/src/Utils/faTrans)  
637 [DCC/kentUtils/tree/master/src/Utils/faTrans](https://github.com/ENCODE-DCC/kentUtils/tree/master/src/Utils/faTrans)).

638

639 **Author Contributions:** JS, OD, PS and UNR developed the concept of OrtSuite. JS, MG, AB and  
640 UNR developed the OrtSuite workflow. JS, MG, AB and UNR performed the benchmarks. CV  
641 provided information and data for defining benzoate to acetyl-CoA conversion pathways. RK  
642 sequenced bacterial isolates that were provided by AC. AB created the interactive network  
643 visualization module. JS and UNR wrote the manuscript. All authors read and commented on  
644 different versions of the manuscript and approved the final manuscript.

645

646 **Funding:** This work was funded by the Helmholtz Young Investigator grant VH-NG-1248 Micro  
647 ‘Big Data’.

648

#### 649 **Data Availability:**

650 The datasets and computer code produced in this study are available in the following databases:

- 651 • The genomes used to test the workflow are available at National Centre for Biotechnology  
652 Information (<https://www.ncbi.nlm.nih.gov/>) under the accession identifiers [CP029389-](#)  
653 [CP029397](#), [GCF\\_000001735](#), [AP012304](#), [AP012305](#), [CP021731](#), [CP011072](#), [CP007785-](#)  
654 [CP007787](#), [CP000614-CP000621](#), [CP003230](#), [CP005996](#), [CP003108](#), [CR555306-](#)  
655 [CR5553068](#), [GCF\\_000225785](#), [LN997848-LN997849](#), [CP022989-CP022996](#), [CP024315](#),  
656 [AP012547](#), [CP022046-CP022047](#) and [CP001281-CP001282](#).
- 657 • The genome assemblies used to predict interspecies interactions are available at National  
658 Centre for Biotechnology Information (<https://www.ncbi.nlm.nih.gov/>) with the study  
659 accession PRJEB38476: (<https://www.ncbi.nlm.nih.gov/bioproject/648592>).
- 660 • OrtSuite scripts: GitHub (<https://github.com/mdsufz/OrtSuite>).

661

662 **Acknowledgments:** We thank the early users of OrtSuite Sandra Silva, Felipe Côrrea and Jonas  
663 Kasmanas for their help with debugging and for workflow suggestions. We also thank Diogo Lima  
664 and Emanuel Cunha for their assistance in the implementation of the script required to generate  
665 the Gene-Protein-Reaction (GPR) rules; and Nicole Steinbach for her work in the sequencing of  
666 the isolates used as the test set.

667 **Conflicts of Interest:** The authors declare no conflict of interest.

668 **References**

- 669 Altschul SF, Gish W, Miller W, Myers EW & Lipman DJ (1990) Basic local alignment search tool. *Journal*  
670 *of molecular biology* 215: 403–10
- 671 Bazire A, Diab F, Jebbar M & Haras D (2007) Influence of high salinity on biofilm formation and benzoate  
672 assimilation by *Pseudomonas aeruginosa*. *Journal of Industrial Microbiology and Biotechnology*  
673 34: 5–8
- 674 Buchfink B, Xie C & Huson DH (2015) Fast and sensitive protein alignment using DIAMOND. *Nat*  
675 *Methods* 12: 59–60
- 676 Delmas E, Besson M, Brice M-H, Burkle LA, Riva GVD, Fortin M-J, Gravel D, Guimarães PR, Hembry  
677 DH, Newman EA, et al (2019) Analysing ecological networks of species interactions. *Biological*  
678 *Reviews* 94: 16–36
- 679 Devanadera A, Vejarano F, Zhai Y, Suzuki-Minakuchi C, Ohtsubo Y, Tsuda M, Kasai Y, Takahata Y,  
680 Okada K & Nojiri H (2019) Complete Genome Sequence of an Anaerobic Benzene-Degrading  
681 Bacterium, *Azoarcus* sp. Strain DN11. *Microbiol Resour Announc* 8
- 682 Dong X & Strous M (2019) An Integrated Pipeline for Annotation and Visualization of Metagenomic  
683 Contigs. *Front Genet* 10
- 684 Emms DM & Kelly S (2015) OrthoFinder: solving fundamental biases in whole genome comparisons  
685 dramatically improves orthogroup inference accuracy. *Genome biology* 16: 157–157
- 686 Emms DM & Kelly S (2019) OrthoFinder: phylogenetic orthology inference for comparative genomics.  
687 *Genome Biology* 20: 238
- 688 Fetzer I, Johst K, Schäwe R, Banitz T, Harms H & Chatzinotas A (2015) The extent of functional  
689 redundancy changes as species' roles shift in different environments. *Proc Natl Acad Sci USA*  
690 112: 14888–14893
- 691 Gabaldón T & Koonin EV (2013) Functional and evolutionary implications of gene orthology. *Nature*  
692 *Reviews Genetics* 14: 360–366
- 693 Gottstein W, Olivier BG, Bruggeman FJ & Teusink B (2016) Constraint-based stoichiometric modelling  
694 from single organisms to microbial communities. *Journal of The Royal Society Interface* 13:  
695 20160627
- 696 Heirendt L, Arreckx S, Pfau T, Mendoza SN, Richelle A, Heinken A, Haraldsdóttir HS, Wachowiak J,  
697 Keating SM, Vlasov V, et al (2019) Creation and analysis of biochemical constraint-based models  
698 using the COBRA Toolbox v.3.0. *Nature Protocols* 14: 639–702
- 699 Hernández-Salmerón JE & Moreno-Hagelsieb G (2020) Progress in quickly finding orthologs as reciprocal  
700 best hits: comparing blast, last, diamond and MMseqs2. *BMC Genomics* 21: 741
- 701 Hu Y, Feng Y, Zhang X & Zong Z (2017) *Acinetobacter defluvii* sp. nov., recovered from hospital sewage.  
702 *International Journal of Systematic and Evolutionary Microbiology*, 67: 1709–1713

- 703 Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, von Mering C & Bork P (2017) Fast  
704 Genome-Wide Functional Annotation through Orthology Assignment by eggNOG-Mapper. *Mol*  
705 *Biol Evol* 34: 2115–2122
- 706 Jenul C, Sieber S, Daepfen C, Mathew A, Lardi M, Pessi G, Hoepfner D, Neuburger M, Linden A,  
707 Gademann K, et al (2018) Biosynthesis of fragin is controlled by a novel quorum sensing signal.  
708 *Nat Commun* 9: 1–13
- 709 Junghare M, Patil Y & Schink B (2015) Draft genome sequence of a nitrate-reducing, o-phthalate  
710 degrading bacterium, *Azoarcus* sp. strain PA01T. *Standards in Genomic Sciences* 10: 90
- 711 Kanehisa M, Goto S, Kawashima S, Okuno Y & Hattori M (2004) The KEGG resource for deciphering the  
712 genome. *Nucleic acids research* 32: D277–D280
- 713 Kanehisa M, Sato Y & Morishima K (2016) BlastKOALA and GhostKOALA: KEGG Tools for Functional  
714 Characterization of Genome and Metagenome Sequences. *J Mol Biol* 428: 726–731
- 715 Khanal A, Yu McLoughlin S, Kershner JP & Copley SD (2015) Differential Effects of a Mutation on the  
716 Normal and Promiscuous Activities of Orthologs: Implications for Natural and Directed  
717 Evolution. *Mol Biol Evol* 32: 100–108
- 718 Kitko RD, Cleeton RL, Armentrout EI, Lee GE, Noguchi K, Berkmen MB, Jones BD & Slonczewski JL  
719 (2009) Cytoplasmic Acidification and the Benzoate Transcriptome in *Bacillus subtilis*. *PLOS ONE*  
720 4: e8255
- 721 Koonin EV (2005) Orthologs, Paralogs, and Evolutionary Genomics. *Annual Review of Genetics* 39: 309–  
722 338
- 723 Lee Y, Lee Y & Jeon CO (2019) Biodegradation of naphthalene, BTEX, and aliphatic hydrocarbons by  
724 *Paraburkholderia aromaticivorans* BN5 isolated from petroleum-contaminated soil. *Sci Rep* 9: 860
- 725 Li L, Stoeckert CJ & Roos DS (2003) OrthoMCL: identification of ortholog groups for eukaryotic genomes.  
726 *Genome Res* 13: 2178–2189
- 727 Li W & Godzik A (2006) Cd-hit: a fast program for clustering and comparing large sets of protein or  
728 nucleotide sequences. *Bioinformatics* 22: 1658–1659
- 729 Locey KJ & Lennon JT (2016) Scaling laws predict global microbial diversity. *PNAS*: 201521291
- 730 Lynch M (2010) Evolution of the mutation rate. *Trends Genet* 26: 345–352
- 731 Lyu Z, Shao N, Akinyemi T & Whitman WB (2018) Methanogenesis. *Curr Biol* 28: R727–R732
- 732 Madden T (2003) The BLAST Sequence Analysis Tool National Center for Biotechnology Information  
733 (US)
- 734 Maestre FT, Castillo-Monroy AP, Bowker MA & Ochoa-Hueso R (2012) Species richness effects on  
735 ecosystem multifunctionality depend on evenness, composition and spatial pattern. *Journal of*  
736 *Ecology* 100: 317–330

- 737 Mendes LW, Raaijmakers JM, de Hollander M, Mendes R & Tsai SM (2018) Influence of resistance  
738 breeding in common bean on rhizosphere microbiome composition and function. *ISME J* 12: 212–  
739 224
- 740 Messina E, Denaro R, Crisafi F, Smedile F, Cappello S, Genovese M, Genovese L, Giuliano L, Russo D,  
741 Ferrer M, et al (2016) Genome sequence of obligate marine polycyclic aromatic hydrocarbons-  
742 degrading bacterium *Cycloclasticus* sp. 78-ME, isolated from petroleum deposits of the sunken  
743 tanker Amoco Milford Haven, Mediterranean Sea. *Marine Genomics* 25: 11–13
- 744 Meyer-Cifuentes I, Fiedler S, Müller JA, Kappelmeyer U, Mäusezahl I & Heipieper HJ (2017) Draft  
745 Genome Sequence of *Magnetospirillum* sp. Strain 15-1, a Denitrifying Toluene Degrader Isolated  
746 from a Planted Fixed-Bed Reactor. *Genome Announc* 5
- 747 Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL, Tosatto SCE, Paladin L,  
748 Raj S, Richardson LJ, et al (2021) Pfam: The protein families database in 2021. *Nucleic Acids*  
749 *Research* 49: D412–D419
- 750 Morin M, Pierce EC & Dutton RJ (2018) Changes in the genetic requirements for microbial interactions  
751 with increasing community complexity. *eLife* 7: e37072
- 752 Mrozik A & Labuzek S (2002) A comparison of biodegradation of phenol and homologous compounds by  
753 *Pseudomonas vesicularis* and *Staphylococcus sciuri* strains. *Acta Microbiol Pol* 51: 367–378
- 754 Mulder CPH, Uliassi DD & Doak DF (2001) Physical stress and diversity-productivity relationships: The  
755 role of positive interactions. *PNAS* 98: 6704–6708
- 756 Ng PC & Henikoff S (2006) Predicting the effects of amino acid substitutions on protein function. *Annu*  
757 *Rev Genomics Hum Genet* 7: 61–80
- 758 O’Sullivan LA, Weightman AJ, Jones TH, Marchbank AM, Tiedje JM & Mahenthiralingam E (2007)  
759 Identifying the genetic basis of ecologically and biotechnologically useful functions of the  
760 bacterium *Burkholderia vietnamiensis*. *Environmental Microbiology* 9: 1017–1034
- 761 Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, Hugenholtz P & Tyson GW  
762 (2017) Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree  
763 of life. *Nat Microbiol* 2: 1533–1542
- 764 Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, Beghini F, Manghi P, Tett A, Ghensi P, et  
765 al (2019) Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000  
766 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* 176: 649–662.e20
- 767 Pearson WR (2013) An Introduction to Sequence Similarity (“Homology”) Searching. *Curr Protoc*  
768 *Bioinformatics* 0 3
- 769 Peng T, Luo A, Kan J, Liang L, Huang T & Hu Z (2018) Identification of A Ring-Hydroxylating  
770 Dioxygenases Capable of Anthracene and Benz[a]anthracene Oxidization from *Rhodococcus* sp.  
771 P14. *MMB* 28: 183–189
- 772 Prestat E, David MM, Hultman J, Taş N, Lamendella R, Dvornik J, Mackelprang R, Myrold DD,  
773 Jumpponen A, Tringe SG, et al (2014) FOAM (Functional Ontology Assignments for

- 774 Metagenomes): a Hidden Markov Model (HMM) database with environmental focus. *Nucleic*  
775 *Acids Res* 42: e145
- 776 Rabus R, Boll M, Heider J, Meckenstock RU, Buckel W, Einsle O, Ermler U, Golding BT, Gunsalus RP,  
777 Kroneck PMH, et al (2016) Anaerobic Microbial Degradation of Hydrocarbons: From Enzymatic  
778 Reactions to the Environment. *MMB* 26: 5–28
- 779 Ravikrishnan A & Raman K (2021) Unraveling microbial interactions in the gut microbiome. *bioRxiv*:  
780 2021.05.17.444446
- 781 Raynaud X & Nunan N (2014) Spatial Ecology of Bacteria at the Microscale in Soil. *PLOS ONE* 9: e87217
- 782 Robertson WJ, Franzmann PD & Mee BJ (2000) Spore-forming, Desulfosporosinus-like sulphate-reducing  
783 bacteria from a shallow aquifer contaminated with gasoline. *Journal of Applied Microbiology* 88:  
784 248–259
- 785 Roh SW, Abell GCJ, Kim K-H, Nam Y-D & Bae J-W (2010) Comparing microarrays and next-generation  
786 sequencing technologies for microbial ecology research. *Trends Biotechnol* 28: 291–299
- 787 Seemann T (2014) Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30: 2068–2069
- 788 Singleton CM, Petriglieri F, Kristensen JM, Kirkegaard RH, Michaelsen TY, Andersen MH, Kondrotaitė Z,  
789 Karst SM, Dueholm MS, Nielsen PH, et al (2021) Connecting structure to function with the  
790 recovery of over 1000 high-quality metagenome-assembled genomes from activated sludge using  
791 long-read sequencing. *Nat Commun* 12: 2009
- 792 Slade EM, Kirwan L, Bell T, Philipson CD, Lewis OT & Roslin T (2017) The importance of species identity  
793 and interactions for multifunctionality depends on how ecosystem functions are valued. *Ecology*  
794 98: 2626–2639
- 795 Sperfeld M, Diekert G & Studenik S (2019) Anaerobic aromatic compound degradation in *Sulfuritalea*  
796 *hydrogenivorans* sk43H. *FEMS Microbiol Ecol* 95
- 797 Suvorova IA & Gelfand MS (2019) Comparative Genomic Analysis of the Regulation of Aromatic  
798 Metabolism in Betaproteobacteria. *Front Microbiol* 10
- 799 Tal O, Selvaraj G, Medina S, Ofaim S & Freilich S (2020) NetMet: A Network-Based Tool for Predicting  
800 Metabolic Capacities of Microbial Species and their Interactions. *Microorganisms* 8: 840
- 801 Thommes M, Wang T, Zhao Q, Paschalidis IC & Segrè D (2019) Designing Metabolic Division of Labor in  
802 Microbial Communities. *mSystems* 4
- 803 Tully BJ, Graham ED & Heidelberg JF (2018) The reconstruction of 2,631 draft metagenome-assembled  
804 genomes from the global oceans. *Scientific Data* 5: 170203
- 805 Valderrama JA, Durante-Rodríguez G, Blázquez B, García JL, Carmona M & Díaz E (2012) Bacterial  
806 Degradation of Benzoate: CROSS-REGULATION BETWEEN AEROBIC AND ANAEROBIC  
807 PATHWAYS. *J Biol Chem* 287: 10494–10508
- 808 Wang B, Lai Q, Cui Z, Tan T & Shao Z (2008) A pyrene-degrading consortium from deep-sea sediment of  
809 the West Pacific and its key member *Cycloclasticus* sp. P1. *Environmental Microbiology*

810 Wang W, Yang S, Hunsinger GB, Pienkos PT & Johnson DK (2014) Connecting lignin-degradation  
811 pathway with pre-treatment inhibitor sensitivity of *Cupriavidus necator*. *Frontiers in*  
812 *Microbiology* 5: 247

813

## 814 **Figure legends**

815 Figure 1 - OrtSuite workflow. OrtSuite takes a text file containing a list of identifiers for each  
816 reaction in the pathway of interest supplied by the user to retrieve all protein sequences from  
817 KEGG Orthology and are stored in ORAdb. Subsequently the same list of identifiers is used to  
818 obtain the Gene-Protein-Reaction (GPR) rules from KEGG Modules (Step 1). Protein sequence  
819 from samples supplied by the user are clustered using OrthoFinder (Step 2). In step 3, the tasks of  
820 functional annotation, identification of putative synergistic species interactions and graphical  
821 visualization of the network are performed. Functional annotation consists of a two-stage process  
822 (relaxed and restrictive search). Relaxed search performs sequence alignments between 50% of  
823 randomly selected sequences from each generated cluster. Clusters whose representative  
824 sequences share a minimum E-value of 0.001 to sequences in the reaction set(s) in ORAdb  
825 transition to the restrictive search. Here, all sequences from the cluster are aligned to all sequences  
826 in the corresponding reaction set(s) to which they had a hit (default E-value =  $1e^{-9}$ ). Next, the  
827 annotated sequences are further filtered to those with a bit score greater than 50 and are used to  
828 identify putative microbial interactions based on their functional potential. Constraints can also be  
829 added to reduce the search space of microbial interactions (e.g. subsets of reactions required to be  
830 performed by single species, transport-related reactions). Additionally, an interactive network  
831 visualization of the results is produced and accessed via a HTML file.

832 Figure 2 - Mapping of the genomic potential of each species from the Fetzer\_genome\_set dataset  
833 to each reaction in aerobic (yellow) and anaerobic (blue) benzoate-to-acetyl-CoA conversion  
834 pathways. Circles highlighted in green represent species that showed biomass growth in medium  
835 containing benzoate in the Fetzer study.

836 Figure 3 – Example of the interactive network visualization included on OrtSuite results. (A) The  
837 complete network with species colored by reaction (B) Species can be highlighted for simple  
838 identification (C). Tooltips on reaction link out the KEGG if the reaction identifier is given.

840 Table 1 - Species names, strain and abbreviation codes of the genomes used to validate OrtSuite (Supplementary data -  
841 Test\_genome\_set). The genomic potential, based on KEGG database, to completely encode all proteins involved in a BTA pathway is  
842 identified in the column “BTA pathway” (P1 – Anaerobic conversion of benzoate to acetyl-CoA 1; P2 – Anaerobic conversion of  
843 benzoate to acetyl-CoA 2; P3 – Aerobic conversion of benzoate to acetyl-CoA). \* indicates no literature was found connecting benzoate  
844 degradation and the respective species.

<b>Name and strain</b>	<b>Abbreviation code</b>	<b>KEGG id</b>	<b>BTA pathway</b>	<b>Accession number</b>	<b>Ref.</b>
<i>Acinetobacter defluvii</i> WCHA30	adv	T05474	P3	CP029389-CP029397	(Hu <i>et al</i> , 2017)
<i>Arabidopsis thaliana</i>	ath	T00041	-	GCF_000001735	*
<i>Azoarcus sp.</i> KH32C	aza	T02502	P2	AP012304, AP012305	(Junghare <i>et al</i> , 2015)
<i>Azoarcus sp.</i> DN11	azd	T05691	P2	CP021731	(Devanadera <i>et al</i> , 2019)
<i>Azoarcus sp.</i> CIB	azi	T04019	P2	CP011072	(Valderrama <i>et al</i> , 2012)
<i>Burkholderia cepacia</i> DDS 7H-2	bced	T03302	P3	CP007785-CP007787	(Jenul <i>et al</i> , 2018)
<i>Burkholderia vietnamiensis</i> G4	bvi	T00493	P3	CP000614-CP000621	(O’Sullivan <i>et al</i> , 2007)
<i>Cycloclasticus sp.</i> P1	cyq	T02265	P3	CP003230	(Wang <i>et al</i> , 2008)
<i>Cycloclasticus zancles</i> 78-ME	cza	T02780	P3	CP005996	(Messina <i>et al</i> , 2016)
<i>Desulfosporosinus orientis</i> DSM 765	dor	T01675	-	CP003108	(Robertson <i>et al</i> , 2000)
<i>Aromatoleum aromaticum</i> EbN1	eba	T00222	P2	CR555306-CR5553068	(Rabus <i>et al</i> , 2016)
<i>Latimeria chalumnae</i> (coelacanth)	lcm	T02913	-	GCF_000225785	*
<i>Magnetospirillum sp.</i> XM-1	magx	T04231	P2	LN997848-LN997849	(Meyer-Cifuentes <i>et al</i> , 2017)
<i>Paraburkholderia aromaticivorans</i> BN5	parb	T05169	P3	CP022989-CP022996	(Lee <i>et al</i> , 2019)
<i>Rhodococcus ruber</i> P14	rrz	T05142	P3	CP024315	(Peng <i>et al</i> , 2018)
<i>Sulfuritalea hydrogenivorans</i> sk43H	shd	T03591	P2	AP012547	(Sperfeld <i>et al</i> , 2019)

---

<i>Staphylococcus sciuri</i> FDAARGOS 285	sscu	T05176	-	CP022046-CP022047	(Mrozik & Labuzek, 2002)
<i>Thauera sp.</i> MZ1T	tmz	T00804	P2, P3	CP001281-CP001282	(Suvorova & Gelfand, 2019)

---



846 Table 2 - OrtSuite workflow runtime and clustering performance. The total runtime of each  
847 OrtSuite step when analyzing the genomic potential of species in Test\_genome\_set dataset in three  
848 pathways (P1, P2 and P3) for the conversion of benzoate to acetyl-CoA (BTA). Steps were  
849 performed with default parameters on a laptop with 4 cores and 16 GB of RAM. Pair-wise  
850 precision and recall results of OrthoFinder using BLAST and DIAMOND as an alignment search  
851 tool. Clustering was performed on the Test\_genome\_set dataset plus the mutated genomes.

<b>OrtSuite step</b>	<b>Runtime</b>
ORAdb construction and Generation of GPR_rules	2h47m
Generation of protein ortholog clusters	54m
Functional annotation of sequences in ortholog clusters	6m
Defining putative microbial interactions	3m
Total	3h50m
Precision (BLAST)	0.63
Recall (BLAST)	0.77
Precision (DIAMOND)	0.64
Recall (DIAMOND)	0.85

852

# ORA database generation and Gene-Protein-Reaction (GPR) rules

Reaction / pathway identifiers [txt]:

- EC
- KO (KEGG orthology ID)
- RID (KEGG reaction ID)
- KEGG pathway map

Manual editing of sequences, reactions and GPRs

**User-defined ORAdb**

# Ortholog clustering

- Protein sequence data [FASTA]
- Species id

Orthofinder

Cluster of orthologs

# Functional annotation

Relaxed search

DIAMOND BLAST

Reduced cluster of orthologs

Restrictive search

DIAMOND BLAST

Annotated sequences

# Putative interactions and exploration

Interactive network visualization [R / HTML]

Functional potential [CSV]

Putative species interactions[CSV]

**Step 1**

**Step 2**

**Step 3**

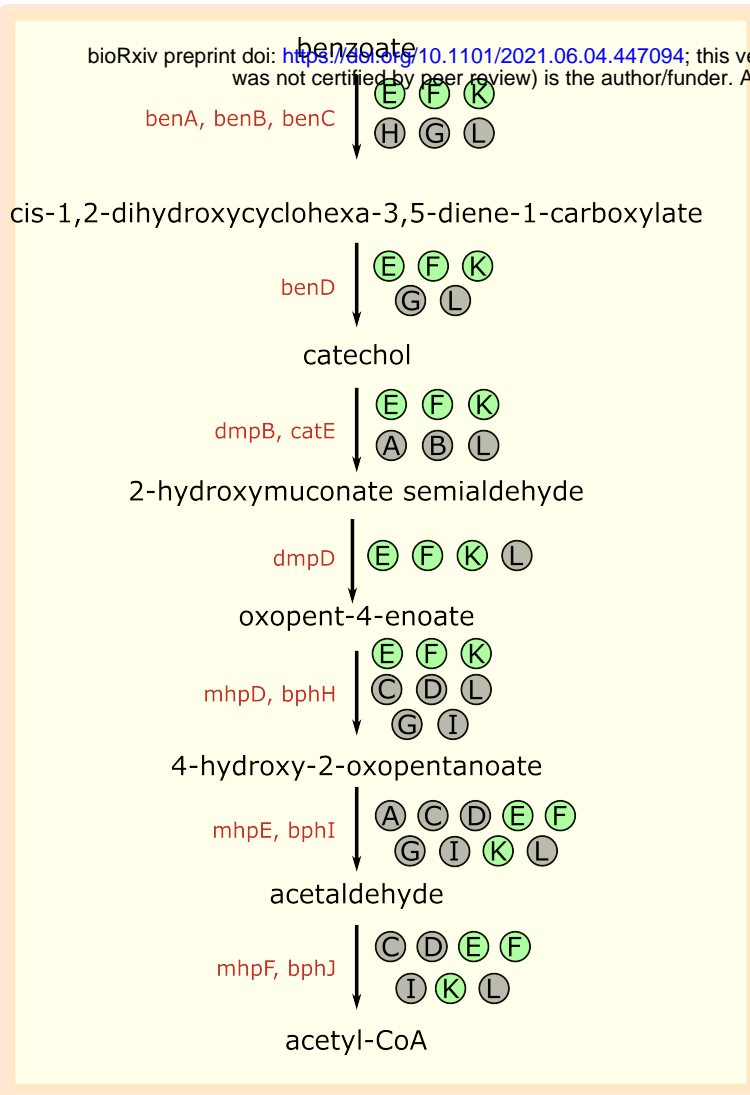
Optional step

Input [format]

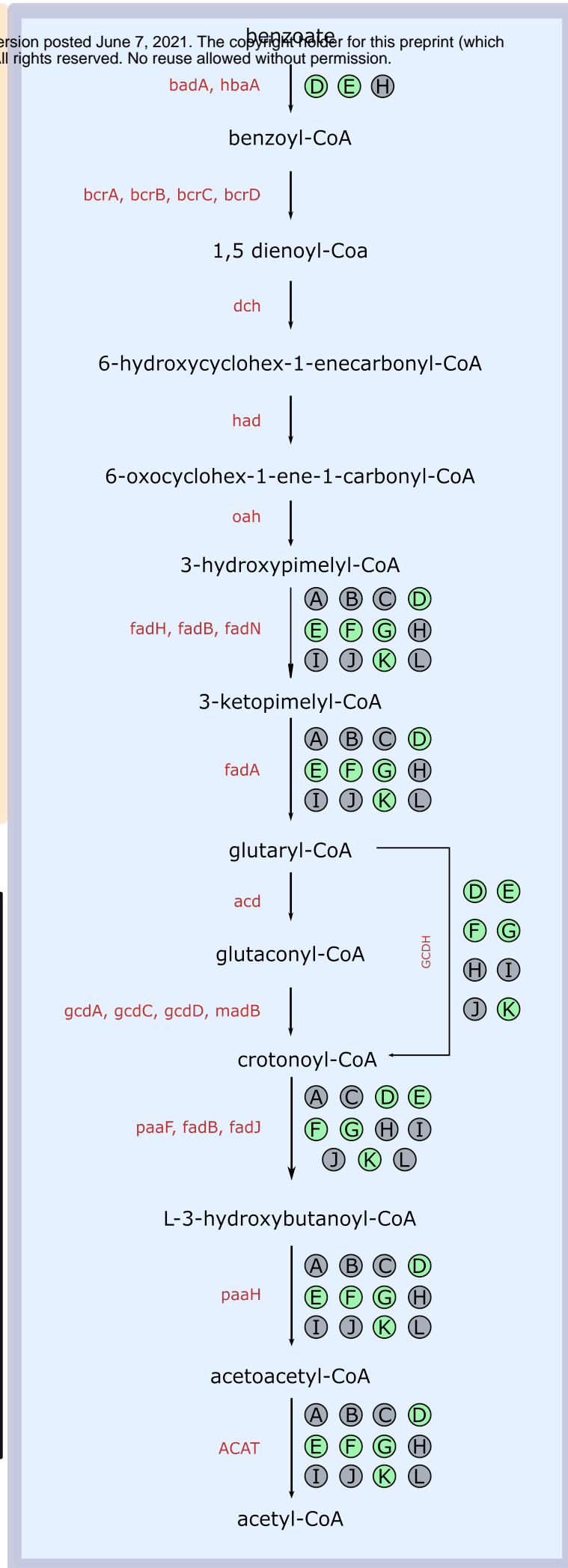
Key tools

Output [format]

## Aerobic



## Anaerobic



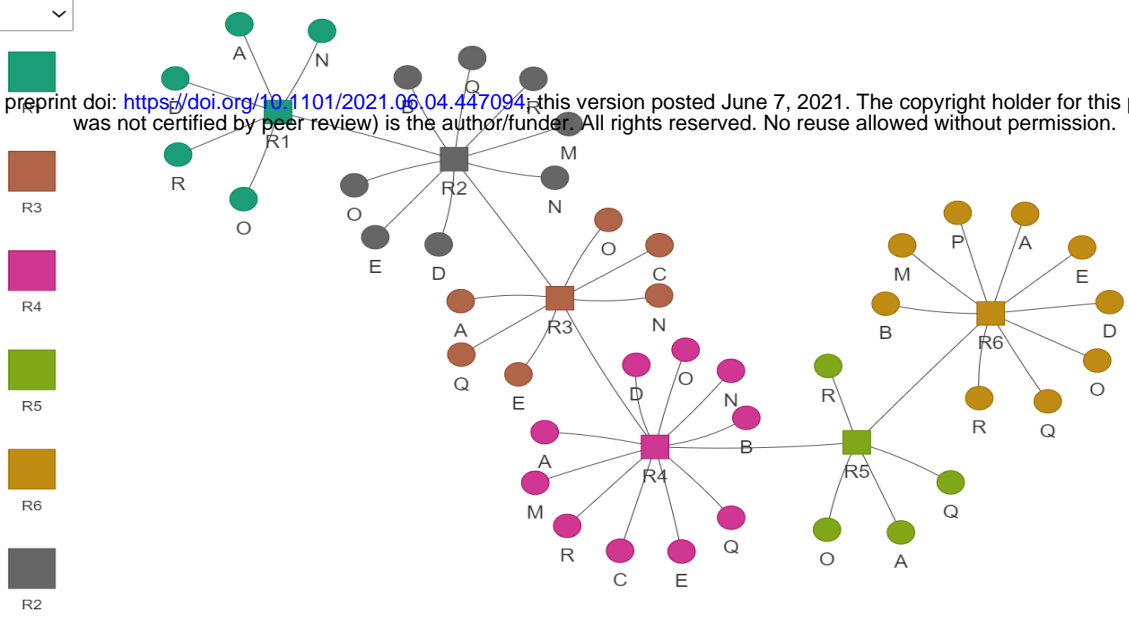
### Legends

- (A) *Bacillus subtilis*
- (B) *Paenibacillus polymyxa*
- (C) *Brevibacillus brevis*
- (D) *Comamonas testosteroni*
- (E) *Cupriavidus necator*
- (F) *Pseudomonas putida*
- (G) *Pseudomonas fluorescens*
- (H) *Variovorax paradoxus*
- (I) *Rhodococcus sp.*
- (J) *Acidovorax facilis*
- (K) *Rhodococcus ruber*
- (L) *Sphingobium yanoikuyae*

**Species with biomass growth in medium containing benzoate (Fetzer et al., 2015)**

Select by label

bioRxiv preprint doi: <https://doi.org/10.1101/2021.06.04.447094>; this version posted June 7, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



Select by label

- M
- A
- B
- C
- D
- E
- M**
- N
- O
- P
- Q
- R
- R1
- R2
- R3
- R4
- R5
- R6

