

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

**João Pedro Saraiva**

Helmholtz-Zentrum für Umweltforschung UFZ

**Alexandre Bartholomäus**

GFZ German Research Centre for Geosciences

**René Kallies**

Helmholtz-Zentrum für Umweltforschung UFZ

**Marta Gomes**

"Universidade do Minho"

**Marcos Bicalho**

Helmholtz-Zentrum für Umweltforschung UFZ Abteilung Umweltmikrobiologie

**Jonas Coelho Kasmanas**

Helmholtz-Zentrum für Umweltforschung UFZ Abteilung Umweltmikrobiologie

**Carsten Vogt**

Helmholtz-Zentrum für Umweltforschung UFZ

**Antonis Chatzinotas**

Helmholtz-Zentrum für Umweltforschung UFZ

**Peter Stadler**

Universität Leipzig

**Oscar Dias**

"Universidade do Minho"

**Ulisses Rocha** (✉ [ulisses.rocha@ufz.de](mailto:ulisses.rocha@ufz.de))

Helmholtz-Zentrum für Umweltforschung UFZ Abteilung Umweltmikrobiologie <https://orcid.org/0000-0001-6972-6692>

---

## Software article

**Keywords:** functional annotation, microbial interactions, microbial modelling, orthologs, partial genome-scale models

**Posted Date:** August 31st, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-52281/v2>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

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 Jonas Coelho Kasmanas<sup>1,4,7</sup>, Carsten Vogt<sup>1</sup>, Antonis Chatzinotas<sup>1,5,6</sup>, Peter Stadler<sup>7,8,9,10,11</sup>, Oscar  
5 Dias<sup>3</sup>, Ulisses Nunes da 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 Mathematics and Computer Sciences, University of Sao Paulo, Sao Carlos, Brazil

12 <sup>5</sup> Institute of Biology, Leipzig University, Leipzig, Germany

13 <sup>6</sup> German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Leipzig,  
14 Germany

15 <sup>7</sup> Bioinformatics Group, Department of Computer Science, Interdisciplinary Center for  
16 Bioinformatics, and Competence Center for Scalable Data Services and Solutions  
17 Dresden/Leipzig, University of Leipzig, Leipzig, Germany

18 <sup>8</sup> Max Planck Institute for Mathematics in the Sciences, Leipzig, Germany

19 <sup>9</sup> Institute for Theoretical Chemistry, University of Vienna, Wien, Austria

20 <sup>10</sup> Facultad de Ciencias, Universidad Nacional de Colombia, Bogotá, Colombia

21 <sup>11</sup> Santa Fe Institute, Santa Fe, U.S.A.

22  
23 \*Correspondence: [ulisses.rocha@ufz.de](mailto:ulisses.rocha@ufz.de)  
24

25 Running title: Mining interactions with OrtSuite

26  
27 Summary blurb: Predicting synergistic species interactions using the genomic potential of  
28 microbial communities.

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

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

47

## 48 Introduction

49 In environments where microorganisms play a crucial role, the microbial community functional  
50 potential encompasses the building blocks for all possible interspecies interactions (Maestre *et al*,  
51 2012; Mulder *et al*, 2001). For example, in environments rich in methane, microbial communities  
52 are dominated by species with genes encoding proteins involved in methanogenesis (Lyu *et al*,  
53 2018). Soil microbes, especially those in the rhizosphere, are genetically adapted to support plants  
54 in the resistance against pathogens and tolerance to stress (Mendes *et al*, 2018). In this context,  
55 natural ecosystems are populated by an enormous number of microbes (Locey & Lennon, 2016).  
56 For example, soil environments can contain more than  $10^{10}$  organisms per gram of soil  
57 heterogeneously distributed making a global search for interspecies interactions unfeasible  
58 (Raynaud & Nunan, 2014). The exponential increase in high-throughput sequencing data and the  
59 development of computational sciences and bioinformatics pipelines have advanced our  
60 understanding of microbial community composition and distribution in complex ecosystems (Roh  
61 *et al*, 2010). This knowledge increased our ability to reconstruct and functionally characterize  
62 genomes in complex communities, for example, by recovering metagenome-assembled genomes  
63 (MAGs) (Parks *et al*, 2017; Pasolli *et al*, 2019; Tully *et al*, 2018). While several tools have been  
64 developed to improve the reconstruction of MAGs, the same cannot be said for predicting  
65 interspecies interactions (Morin *et al*, 2018). Studies by Parks (Parks *et al*, 2017) and Tully (Tully  
66 *et al*, 2018), while advancing the reconstruction of MAGs, did not perform any functional  
67 characterization or prediction of interspecies interactions. Pasolli and collaborators (Pasolli *et al*,  
68 2019) performed functional annotation of representative species in their study by employing  
69 several tools such as EggNOG (Huerta-Cepas *et al*, 2017), KEGG (Kanehisa *et al*, 2004) and  
70 DIAMOND (Buchfink *et al*, 2015). However, the sheer number of representative genomes (4930)  
71 and the lack of focus on specific ecosystem processes make predicting interspecies interactions  
72 challenging.

73 Furthermore, the challenge of predicting interspecies interactions increases because of the  
74 multitude of potential interactions between species in microbial communities and between  
75 microbes and their hosts (e.g., plants, animals and microeukaryotes) (Slade *et al*, 2017). An  
76 integrated pipeline for annotation and visualization of metagenomes (MetaErg) (Dong & Strous,  
77 2019) attempted to address some of the challenges in metagenome annotation such as the inference  
78 of biological functions and integration of expression data. MetaErg performs comprehensive  
79 annotation and visualization of MAGs by integrating data from multiple sources such as Pfam  
80 (Mistry *et al*, 2021), KEGG (Kanehisa *et al*, 2004) and FOAM (Prestat *et al*, 2014). However,  
81 MetaErg's full genome annotation requires elevated processing times and computational resources  
82 due to its untargeted approach. In addition, there is a lack of a user-friendly tool to explore the  
83 results tables and graphs to extract pathway-specific information tied to each MAG and thus infer  
84 potential species interactions based on their functional profiles.

85 Genome-based modeling approaches have routinely been used to study single organisms as well  
86 as microbial communities (Gottstein *et al*, 2016). For example, constraint-based models are highly  
87 employed in studying and predicting metabolic networks (Heirendt *et al*, 2019). These models are  
88 generated upon the premise that any given function is feasible as long as the protein-encoding gene

89 is present. Although species may lack the genetic potential to perform all functions necessary to  
90 survive in a given ecosystem, outside laboratory conditions, microbes do not exist in isolation and  
91 may benefit from their interaction with other species. By assessing the genomic content of  
92 individual species, we are able to identify groups of microbes whose combined content may  
93 account for complete ecosystem functioning. However, generating full genome metabolic  
94 networks for each microbial community species is time-consuming as they require information not  
95 easily obtained for each community member, such as biomass composition and nutritional  
96 requirements.

97 In order to decrease complexity and facilitate analysis, it is possible to limit the search of  
98 interactions to groups of organisms (e.g., microbe-microbe or host-microbe) or specific ecosystem  
99 processes (e.g., nitrification or deadwood decomposition). A network-based tool for predicting  
100 metabolic capacities of microbial communities and interspecies interactions (NetMet) was recently  
101 developed (Tal *et al*, 2020). This tool only requires a list of species-specific enzyme identifiers  
102 and a list of compounds required for a given environment. However, besides the necessity of  
103 previous annotation of genomes, NetMet does not consider the rules that govern each reaction (e.g.  
104 protein complexes). Accurate annotation of gene function from sequencing data is essential to  
105 predict ecosystem processes potentially performed by microbial communities, particularly in cases  
106 where an ecosystem process is performed by the synergy of two or more species. Simple methods  
107 for the annotation of genomes rely, for instance, on the search for homologous sequences.  
108 Computational tools such as BLAST (Altschul *et al*, 1990) and DIAMOND (Buchfink *et al*, 2015)  
109 compare nucleotide or protein sequences to those present in databases. These approaches allow  
110 inferring the function of uncharacterized sequences from their homologous pairs whose function  
111 is already known. The degree of confidence in the assignment of biological function is increased  
112 if this has been validated by, for example, experimental data. Approaches based on orthology are  
113 increasingly used for genome-wide functional annotation (Huerta-Cepas *et al*, 2017). Orthologs  
114 are homologous sequences that descend from the same ancestor separated after a speciation event  
115 retaining the same function (Koonin, 2005). OrthoMCL (Li *et al*, 2003), CD-HIT (Li & Godzik,  
116 2006) and OrthoFinder (Emms & Kelly, 2015, 2019) are just a few tools that identify homologous  
117 relationships between sequences using orthology. OrthoFinder is more accurate than several other  
118 orthogroup inference methods since it considers gene length in detecting ortholog groups by  
119 introducing a score transformation step (Emms & Kelly, 2015). However, OrthoFinder, due to its  
120 all-versus-all sequence alignment approach, requires intensive computational resources resulting  
121 in long-running times when using large data sets for clustering. Because of the enormous number  
122 of potential combinations, limiting the scope of research to specific ecosystem processes may  
123 reduce the computational and resource costs associated with integrating ortholog clustering tools  
124 and functional annotation strategies. Still, having a pipeline that performs targeted annotation of  
125 genomes and genomic-based prediction of putative synergistic species interactions can assist  
126 researchers in the discovery of key players in any metabolic process. Furthermore, the  
127 identification of potential species interactions can lead to the design of synthetic microbial  
128 communities with a wide range of applications such as in bioremediation (Sharma & Shukla,  
129 2020), energy production (Jiang *et al*, 2020) and human health (Clark *et al*, 2021).

130 In this study, we developed OrtSuite; a workflow that can: (i) perform accurate ortholog based  
131 functional annotation, (ii) reveal putative microbial synergistic interactions, and (iii) digest and  
132 present results for pathway and community driven biological questions. These different features  
133 can be achieved with the use of three bash commands in a reasonable computational time. This  
134 research question / hypothesis-targeted approach integrates a user-defined, Ortholog–Reaction  
135 Association database (ORAdb) with up-to-date ortholog clustering tools. OrtSuite allows the  
136 search for putative microbial interactions by calculating the combined genomic potential of  
137 individual species in specific user-defined ecosystem processes. OrtSuite also provides a visual  
138 representation of the species’ genetic potential mapped to each of the reactions defined by the user.  
139 We evaluate this workflow using a clearly defined set of reactions involved in the well-described  
140 benzoate-to-Acetyl-CoA (BTA) conversion. Further, we used this workflow to functionally  
141 characterize a set of known benzoate degraders. OrtSuite’s ability to identify putative interspecies  
142 interactions was evaluated on species whose potential interactions have been previously predicted  
143 under controlled conditions (Fetzer *et al*, 2015).

144

## 145 **Results**

146 One of the motivations to develop OrtSuite was to facilitate the targeted analysis of microbial  
147 communities’ genomic potential, including the prediction of putative synergistic interspecies  
148 interactions. To simplify combining targeted functional annotation with the prediction of species  
149 interactions, we developed OrtSuite to integrate ortholog clustering tools (Emms & Kelly, 2019)  
150 with sequence alignment programs (Buchfink *et al*, 2015).

151

### 152 **OrtSuite is a flexible and user-friendly pipeline**

153 Three simple-to-use scripts were created to collectively perform all tasks associated with OrtSuite  
154 and provide a user-friendly execution. Users would only be required to provide a list of identifiers  
155 related to the ecosystem process of interest and the FASTA files (in protein format) of the species  
156 for which they intend to predict interactions. Next, the users only need to execute three simple  
157 bash commands that cover database generation, functional annotation, and species interactions.

158 Briefly, OrtSuite performs the following processing steps (Figure 1) (for further details see  
159 Methods): Step 1 – In this step, the script *DB\_construction.sh* takes the list of identifiers provided  
160 by the user and automatically downloads the protein sequences that will populate ORAdb; Step 2  
161 - In this step, the script *DB\_construction.sh* takes the list of KO identifiers obtained during Step 1  
162 and downloads the Gene-Protein-Reaction (GPR) rules from KEGG Modules; Step 3 – In this step,  
163 the function *orthofinder* performs the clustering of orthologs; Step 4 – In this step, the script  
164 *annotate\_and\_predict.sh* takes as input the FASTA files containing the Open Reading Frames  
165 (ORFs) of the genomes of interest and performs functional annotation (aligning them against the  
166 sequences in ORAdb); Step 5 – In this step, the script *annotate\_and\_predict.sh* performs the  
167 prediction of putative synergistic interspecies interactions (Figure 1) using the output file  
168 “*Reactions\_mapped\_to\_species.csv*” generated during Step 4. While not necessary, additional

169 control is given to the user with the option to establish thresholds in the minimum e-values (during  
170 sequence alignment of sequences in ortholog clusters to ORAdb). Other constraints include  
171 restricting the number of putative microbial interactions based on the presence of transporters and  
172 subsets of reactions to be performed by individual species (Supplementary data – Table S1). Data  
173 in public repositories continues to be added or updated. Thus, manual inspection of the files in the  
174 ORAdb and GPR rules, although not mandatory, is strongly advised.

175 Users can choose from two alternatives to install OrtSuite. They may use a docker image for  
176 personal computers or conda packages (recommended for installation for High-Performance  
177 Computers). We created a user-friendly git repository (<https://github.com/mdsufz/OrtSuite>) that  
178 provides users with a user-friendly guide covering the installation and the three scripts used to run  
179 our pipeline and the generated outputs.

180

### 181 **Computing time of OrtSuite stages**

182 We evaluated the runtime of each OrtSuite step on a set of genomes whose genomic potential in  
183 converting benzoate to acetyl-CoA was known (Table 1). OrtSuite was executed on a laptop with  
184 4 cores and 16 Gigabytes of RAM. We ran all OrtSuite steps on default settings, and recorded the  
185 total runtime of each step (Table 2). The entire workflow was completed in 3 h 50 min, and the  
186 longest step was the construction of the ORAdb, which consisted of 2 h and 47 min which  
187 involved. The user can modify the number of cores used during functional annotation, further  
188 decreasing run times.

189

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

191 Point mutations can have a drastic effect on the functional profile of microbes by altering the  
192 expected amino acid composition. Thus, we evaluated the impact of point mutations during the  
193 clustering of orthologs using OrthoFinder (Emms & Kelly, 2019). OrthoFinder allows users to  
194 choose between DIAMOND (Buchfink *et al*, 2015), BLAST (Altschul *et al*, 1990) and MMSeqs2  
195 (Steinegger & Söding, 2017, 2) as sequence aligners. DIAMOND and BLAST are the most  
196 commonly used sequence aligners. Therefore, we evaluated the clustering of orthologs these two  
197 tools. Nevertheless, the user may opt for MMSeqs2 as the sequence aligner when using OrtSuite.  
198 To test which of the two sequence aligners (DIAMOND or BLAST) yielded the best results, we  
199 performed ortholog clustering of a dataset consisting of the original target genomes and a set of  
200 artificially mutated genomes (Supplementary data - Test\_genome\_set) using both aligners. The  
201 results showed a 0.01 difference between BLAST and DIAMOND precision (Table 2). However,  
202 DIAMOND showed a 9.5% higher recall than that observed for BLAST what suggests DIAMOND  
203 may have higher sensitivity in the clustering of sequences with the same function. All artificially  
204 mutated sequences (even those with mutation rates of 25%) were clustered together with their non-  
205 mutated ortholog. In parallel, we also performed sequence alignment using NCBI's BLASTp  
206 (Madden, 2003) between the protein sequences of the DNA-mutated and un-mutated genes. E-  
207 values of sequence alignments in all species ranged from 0 to  $5e^{-180}$  and percentage of identity

208 from 61.32 to 98.84% (Supplementary data – Table S2). For validation of the OrtSuite workflow,  
209 clustering of protein orthologs was repeated using only the original unmutated 18 genomes and  
210 the default aligner (DIAMOND). We also generated a complete overview of the results generated  
211 during the clustering of orthologs (e.g., number of genes in ortholog clusters, number of  
212 unassigned genes, and number of ortholog clusters) (Supplementary data - Table S3).

213

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

215 The third step of OrtSuite consists of performing cluster annotation in a two-stage process. In the  
216 first, only 50% of sequences are used to align the sequences from ORAdb. Those with a minimum  
217 e-value proceed to the second stage, where all sequences contained in this cluster will be aligned.  
218 At the end, annotation of clusters will take into consideration additional parameters such as bit  
219 scores. To evaluate the thresholds used in the annotation of ortholog clusters, we used one relaxed  
220 (0.001) and four restrictive ( $1e^{-4}$ ,  $1e^{-6}$ ,  $1e^{-9}$  and  $1e^{-16}$ ) e-value cutoffs. An overview of the results  
221 (e.g., number of clusters containing orthologs from ORAdb, number of ortholog clusters with  
222 annotated sequences) is shown in (Supplementary data – Table S4). The performance of OrtSuite  
223 in the functional annotation of the genomes in the *Test\_genome\_set* is shown in (Supplementary  
224 data – Table S5). On average, 96% of the annotations assigned by KEGG were also identified by  
225 OrtSuite. The complete list of functional annotation results using the different e-value cutoffs is  
226 available in the Supplementary data - Table S6, Table S7, Table S8 and Table S9. Similarly, we  
227 used different e-value cutoffs for testing the mapping of species with the genetic potential for each  
228 reaction (considering the GPR rules) (Supplementary data – Table S10, Table S11, Table S12 and  
229 Table S13). The four different e-value cutoffs used during the restrictive search stage yielded  
230 similar results in terms of annotation. However, the largest decrease in the number of ortholog  
231 clusters that transits from the relaxed search to the restrictive occurs when using an e-value cutoff  
232 of  $1e^{-16}$  (Supplementary data – Table S4). The difference in computing time between lower and  
233 higher e-value thresholds was negligible (< 2 min). Other annotation tools, such as NCBI's BLAST  
234 tool (Altschul *et al*, 1990), BlastKOALA (Kanehisa *et al*, 2016) and Prokka (Seemann, 2014), can  
235 annotate full genomes, the latter at a relatively fast pace. On average, full genome annotation of  
236 our genomes in the *Test\_genome\_set* dataset using Prokka required 12 mins on a standard laptop  
237 with 16 Gigabytes of RAM and four CPUs to complete. BlastKOALA required approximately 3  
238 hours to annotate a single genome. However, the use of these tools resulted in longer runtimes or  
239 in additional manual processing of files generated from full genome annotations for filtering  
240 pathways of interest.

241

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

243 To test OrtSuite's ability to identify species with the genetic potential to perform a pathway  
244 individually, we defined sets of reactions used in three alternative pathways for converting  
245 benzoate to acetyl-CoA (Supplementary data – Table S14). Next, we compared the results to the  
246 species' known genomic content in each alternative pathway (Supplementary data – Table S15).  
247 OrtSuite matched KEGG's predictions in species' ability to perform each alternative benzoate

248 degradation pathway in all but two species - *Azoarcus* sp. DN11 and *Thauera* sp. MZ1T.  
249 Furthermore, OrtSuite identified five species capable of performing conversion pathways not  
250 contemplated in KEGG. *Azoarcus* sp. KH32C, *Aromatoleum aromaticum* EbN1,  
251 *Magnetospirillum* sp. XM-1 and *Sulfuritalea hydrogenivorans* sk43H have the genetic potential to  
252 perform both pathways involving the anaerobic conversion of benzoate to acetyl-CoA while  
253 *Azoarcus* sp. CIB has to genetic potential to achieve all alternative pathways (except when using  
254 an e-value cutoff of  $1e^{-16}$ ). No genes in *Thauera* sp. MZ1T involved in the conversion of crotonyl-  
255 CoA to 3-Hydroxybutanoyl-CoA (R03026) were identified by OrtSuite; this enzyme is essential  
256 for the anaerobic conversion of benzoate to acetyl-CoA. OrtSuite's performance yielded similar  
257 results between all tested e-value cutoffs. However, we observed a higher drop in the number of  
258 ortholog clusters whose sequences are all annotated with the same function when using an e-value  
259 cutoff of  $1e^{-16}$ . Thus, we set the default e-value for the restrictive search to  $1e^{-9}$ .

260

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

262 In this study, we tested the ability of OrtSuite in identifying interspecies interactions involved in  
263 the conversion of benzoate to acetyl-CoA where experimental data were available. We assessed  
264 the prediction of synergistic interspecies interactions on a set of sequenced isolates  
265 (Supplementary data - Fetzer\_genome\_set.zip). In a previous study (Fetzer *et al*, 2015), the  
266 potential of these isolates to grow in benzoate were analyzed individually and in combination  
267 under three different environmental conditions. These conditions were: low substrate  
268 concentration (1g/L benzoate); high substrate concentration (6g/L benzoate); and, high substrate  
269 concentration with additional osmotic stress (6g/L benzoate supplemented with 15g /L of NaCl).  
270 In that study, Fetzer et al investigated if the presence or absence of a particular species positively  
271 or negatively affected biomass production. Since under specific conditions the presence of a  
272 degrader alone was not sufficient for biomass production, they further analyzed if potential species  
273 interactions could be of relevance. Briefly, Fetzer and collaborators defined minimal communities  
274 for all environmental conditions. Next, they tested whether the inclusion of other species in a  
275 community stimulated biomass production. When co-cultures produced biomass, the authors  
276 suggested the species in those communities had the potential to synergistically metabolize  
277 benzoate (Fetzer *et al*, 2015). Using OrtSuite, we aimed to identify which potential species  
278 interactions predicted by Fetzer and collaborators could result from their combined genetic  
279 potential.

280 Our dataset contained 69,193 protein sequences distributed across the 12 species, resulting in 59  
281 Megabytes of data. More than 84% of all genes were placed in 9,533 ortholog clusters. In addition,  
282 541 clusters were composed of sequences obtained from all 12 species (Supplementary data -  
283 Table S16). OrtSuite's annotation stage resulted in 326 ortholog clusters with annotated sequences  
284 from ORAdb (Supplementary data - Table S17). The mapping of KOs to each species in the  
285 *Fetzer\_genome\_set* is available as supplementary data (Table S18). The genomic potential of each  
286 species for aerobic and anaerobic benzoate metabolizing pathways is shown in Figure 2. The  
287 complete mapping of reactions to each species is available in the supplementary data (Table S19).  
288 Based on the 326 ortholog clusters and the Gene-Protein-Reaction (GPR) rules (Supplementary

289 data - Table S20), five species (*Cupriavidus necator* JMP134, *Pseudomonas putida* ATCC17514,  
290 *Rhodococcus sp.* Isolate UFZ (Umweltforschung Zentrum), *Rhodococcus ruber* BU3 and  
291 *Sphingobium yanoikuyae* DSM6900) contained all protein-encoding genes required to perform  
292 aerobic conversion of benzoate to acetyl-CoA. In the Fetzer study, *Rhodococcus sp.* Isolate UFZ  
293 and *S. yanoikuyae* did not show growth in a medium containing benzoate. The incomplete  
294 functional potential of *C. testosteroni* ATCC 17713 and *P. putida* ATCC17514 to perform aerobic  
295 conversion of benzoate to acetyl-CoA is at odds with their reported growth as monocultures in the  
296 presence of benzoate as shown in the Fetzer study. The number of species with the genetic potential  
297 for each reaction involved in the aerobic benzoate degradation pathway (P3) is shown in  
298 (Supplementary data - Table S21). All species with the complete genomic potential to perform a  
299 complete pathway were excluded when calculating interspecies interactions since they do not  
300 require the presence of others. However, species identified by OrtSuite with the complete  
301 functional potential to perform each defined pathway were also included to compare to the results  
302 in the Fetzer study presented above. A total of 2382 combinations of species interactions were  
303 obtained whose combined genetic potential covered all reactions. The complete list of potentially  
304 interacting species is available in the supplementary data (Table S22).

305 In the anaerobic degradation pathways (P1 and P2) no species presented the genomic content to  
306 encode proteins involved in the conversion of benzoyl-CoA to Cyclohexa-1,5-diene-1-carboxyl-  
307 CoA (R02451) (Supplementary data - Table S23). This reaction requires the presence of a protein  
308 complex either composed of four subunits (K04112, K04113, K04114, K04115) or composed of  
309 two subunits (K19515, K19516). The genomes of the 12 species studied contained all subunits in  
310 either protein complex. Therefore, no species interactions were identified that would allow the  
311 complete anaerobic conversion of benzoate to acetyl-CoA. In the low substrate environment,  
312 OrtSuite identified 826 of 830 (99.5%) species combinations showing growth. In the high substrate  
313 environment, OrtSuite predicted 644 of 646 (99.7%). In the high substrate+salt stress environment,  
314 OrtSuite predicted all 271 (100%) combinations of species exhibiting growth (Supplementary data  
315 - Table S24).

## 316 Discussion

317 We designed OrtSuite to allow hypothesis-driven and user-friendly exploration of microbial  
318 interactions. Our team achieved this by integrating up-to-date clustering tools with faster sequence  
319 alignment methods and limiting the scope to user-defined ecosystem processes or metabolic  
320 functions. Using only three bash commands required to run the complete workflow, OrtSuite is a  
321 user-friendly tool capable of running in a customary computer (four cores and 16GB of RAM)  
322 with even faster runtimes when using high-performance computing.

323 The clustering of orthologs by OrthoFinder using DIAMOND (Buchfink *et al*, 2015) showed  
324 higher sensitivity and lower runtime compared to BLAST (Altschul *et al*, 1990) which has also  
325 been shown by Hernández-Salmerón and Moreno-Hagelsieb (Hernández-Salmerón & Moreno-  
326 Hagelsieb, 2020). Furthermore, low e-values and medium to high identity percentages in the  
327 sequence alignments between mutated and original genes indicated that the mutated genes still  
328 share enough sequence similarity to the original protein sequence. These results suggest that  
329 mutation rates of up to 25% of single DNA base pairs will not have an observable effect on the

330 clustering of orthologs. OrthoFinder's algorithm removes the gene length bias from the sequence  
331 alignment process, which may also explain why mutated genes were clustered with the original.  
332 Although it has been suggested that most genetic variations are neutral, changes in single base  
333 pairs can have a drastic effect on protein function (e.g., depending on the location of the mutation)  
334 (Ng & Henikoff, 2006). To this purpose, experimental functional studies can be used to validate  
335 previously unannotated orthologs. Furthermore, this study case does not consider the distribution  
336 of mutations across species and gene families, which can also have different effects on the  
337 clustering of orthologs (Khanal *et al*, 2015). Therefore, future studies increasing the rates of DNA  
338 base-pair substitutions and other types of mutations and experiments targeting protein function in  
339 ortholog clusters are needed.

340 Next, we aimed to improve and facilitate functional annotation and prediction of synergistic  
341 microbial interactions. Exploring the great amount of data generated from full genome annotation  
342 of individual species from complex microbial communities is daunting. This is evident in a study  
343 by Singleton and collaborators (Singleton *et al*, 2021) where the connection between structure and  
344 function required the analysis of metagenomics data, 16S and molecular techniques such as  
345 fluorescent in situ hybridization and Raman spectroscopy. When we looked at functional  
346 annotation alone, two challenges arose. First, performing all vs. all sequence alignments in  
347 complex communities is resource-consuming (time and computational power). Second, manual  
348 inspection of each annotated genome for target genes or pathways is required. Identifying  
349 interspecies interactions based on the microbe's complete genomic potential is also challenging.  
350 For example, ecologists increasingly employ network approaches, but selecting the most  
351 appropriate approach is not always straightforward and easy to implement (Delmas *et al*, 2019).  
352 OrtSuite overcomes these challenges by first performing cluster annotation in a two-stage process  
353 and limited to a user-defined set of functions, decreasing the number of sequence alignments  
354 necessary. The user-defined database coupled with the scripts for automated identification of  
355 interspecies interactions contained in OrtSuite decreases the time required to generate the data and  
356 facilitates its interpretation by the user. Additionally, OrtSuite generates a graphical representation  
357 of the network enabling the use of the whole microbial community (Figure 3)  
358 ([https://github.com/mdsufz/OrtSuite/blob/master/network\\_example.png](https://github.com/mdsufz/OrtSuite/blob/master/network_example.png)).

359 OrtSuite not only confirmed all but two of KEGG's predictions in species' ability to perform each  
360 alternative benzoate degradation pathway used in this study but also identified five species capable  
361 of performing conversion pathways not contemplated in KEGG. On average, an additional 18.3  
362 KO identifiers were mapped to genes not previously annotated in the species used in this study.  
363 Using e-value and bit score as the filtering criteria rather than sequence identity, employed by  
364 KEGG, may explain the increase in functionally annotated genes. For example, the alignment of a  
365 sequence of *A. defluvii* (adv: AWL30228.1) to the sequences in ORAdb annotated as K04105  
366 (conversion of benzoate to benzoyl-CoA) showed high bit-scores (200.7) and low e-values ( $2e^{-54}$ )  
367 but the identity percentage did not exceed 28.6%. The use of e-values and bit scores to infer  
368 function has been reviewed by Pearson (Pearson, 2013). Pearson suggests that e-values and bit  
369 scores are more sensitive and reliable than identity percentages in finding homology since they  
370 consider the evolutionary distance of aligned sequences, the sequence lengths and the scoring  
371 matrix.

372 To test the prediction of putative synergistic microbial interactions, we used data from an  
373 independent study performed by Fetzer and collaborators (Fetzer *et al*, 2015); hereafter Fetzer  
374 study. In the Fetzer study, five species showed biomass growth (estimated by optical density at  
375 590nm wavelength) in a medium containing benzoate. We evaluated whether these species  
376 possessed the complete genomic content to encode all proteins required for each benzoate to  
377 acetyl-CoA conversion pathway. The remaining seven species could not grow as monocultures in  
378 media with benzoate as the sole carbon source. Therefore, we evaluated whether the lack of growth  
379 was explained by the absence of essential protein-encoding genes involved in converting benzoate  
380 to acetyl-CoA. The Fetzer study also showed that, under specific nutrient and stress conditions,  
381 total biomass production was influenced by the presence of non-degrading species. Thus, we  
382 evaluated whether putative species interactions identified by OrtSuite fit the results obtained by in  
383 the Fetzer study. OrtSuite confirmed the functional potential for aerobic conversion of benzoate to  
384 acetyl-CoA in three of the five species whose growth in monocultures was observed during their  
385 research. In the Fetzer's study, *S. yanoikuyae* (accession number GCA\_903797735.1) and  
386 *Rhodococcus sp.* (accession number GCA\_903819475.1) were not able to grow as monoculture  
387 in the presence of benzoate. However, OrtSuite predicted that both possessed the functional  
388 potential to aerobically convert benzoate to acetyl-CoA. In their study, growth was considered  
389 when optical densities (OD) were above 0.094. The OD measured for *S. yanoikuyae* was 0.0916  
390 in a medium containing 1g/L of benzoate. The annotation of genes with the ability to perform the  
391 complete aerobic conversion of benzoate to acetyl-CoA combined with a small difference in OD  
392 to the minimum threshold suggests that *S. yanoikuyae* can grow, albeit slowly, on low benzoate  
393 containing medium. In the case of *Rhodococcus sp.* Isolate UFZ, the OD was never measured  
394 above 0.022 what might indicate another slow-growing species. Another possible explanation is  
395 that although these two species possess the genes necessary for aerobic benzoate degradation, they  
396 are not active. In Fetzer's study, the observed growth of *Comamonas testosteroni* ATCC11996  
397 and *Pseudomonas fluorescens* DSM6290 in the low benzoate environment was not explained by  
398 OrtSuite. To note, benzoate conversion intermediates were not determined in the Fetzer  
399 experiment. Hence, these two species may utilize reactions or pathways that were not included in  
400 the benzoate degradation pathways used in our study. Despite the presence of benzoate degraders,  
401 another possible explanation as to the unobserved growth in Fetzer's study for certain experimental  
402 conditions is the lack of tolerance of these species to high benzoate concentrations. For example,  
403 *C. necator* growth was stimulated at low benzoic acid concentrations but inhibited at high  
404 concentrations (Wang *et al*, 2014). In addition, the set of genes used in our study did not consider  
405 the presence of stress-related factors. To assess these effects, stress-resistance associated genes  
406 and reactions such as those involved in medium acidification (Kitko *et al*, 2009) could be added  
407 as constraints. Similar results were obtained when using a high substrate and salt stress medium.  
408 Under these conditions, the presence of benzoate degraders alone may not be sufficient to achieve  
409 growth. Benzoate degradation has been shown to decrease in hyperosmotic environments (Bazire  
410 *et al*, 2007). Therefore, additional constraints such as genes that confer resistance to environmental  
411 stressors or adverse conditions sodium chloride (NaCl) could be included in identifying  
412 interspecies interactions under different or changing environmental conditions.

413 No single species or combination of species possessed the complete genomic potential to  
414 anaerobically convert benzoate to acetyl-CoA via the two proposed pathways (P1 and P2). Since

415 all growth experiments were conducted in aerobic conditions, the species in question may only use  
416 benzoate as a carbon source in aerobic environments. To fully explore all the species potential to  
417 convert benzoate, additional degradation pathways could be checked in the future using a multi-  
418 omics approach. For example, OrtSuite users could potentially integrate the use of  
419 (meta)transcriptomic data during the prediction of interspecies interactions by excluding species  
420 showing no gene expression of the selected pathways. However, the analysis and integration of  
421 (meta)transcriptomic data is not trivial and would add more levels of complexity to consider (e.g.,  
422 expression of a gene can be high but protein be inactive) and is out of OrtSuite's scope.  
423 Furthermore, the only constraints added were related to the reactions that composed each pathway.  
424 Additional constraints can be included in future studies, such as potential mandatory transport-  
425 associated reactions, to increase confidence in the proposed interspecies interactions. Also, species  
426 interactions can be manually excluded if, for instance, antibiotic-producing species are known to  
427 inhibit the growth of others. OrtSuite confirmed that most interspecies interactions (> 99%)  
428 identified by Fetzer and collaborators were possible due to their combined metabolic potential to  
429 aerobically degrade benzoate to acetyl-CoA but not under anoxic conditions.

430 In this study, we ran OrtSuite on a dataset composed of 18 genomes (Table 1). To determine if this  
431 range would be within the number of genomes in regular microbiome studies, we calculated the  
432 average number of MAGs from different studies focusing on their recovery. A study performed  
433 by Parks and collaborators (Parks *et al*, 2017) analyzed sequencing data from 149 projects. Most  
434 projects (91%) consisted of less than 20 samples. On average, they recovered 5.3 metagenome-  
435 assembled genomes (MAGs) per metagenome. Work performed by Pasolli and collaborators  
436 (Pasolli *et al*, 2019) on microbial diversity in the human microbiome recovered, on average, 16  
437 MAGs per metagenomic library. From the 46 studies used in their work, 30 consisted of less than  
438 200 samples. Another study by Tully and collaborators focusing on marine environments (Tully  
439 *et al*, 2018) recovered 2631 MAGs from 234 samples (average of 11 MAGs per sample). Our  
440 analysis demonstrates that the average number of MAGs recovered from a metagenomic library  
441 ranges from five to 16. Therefore, by using a regular laptop, users can perform targeted functional  
442 annotation and interspecies interactions predictions using OrtSuite in average-sized metagenomes.

443 In summary, OrtSuite allows hypothesis-driven exploration of potential interactions between  
444 microbial genomes by limiting the search universe to a user-defined set of ecosystem processes.  
445 This is achieved by rapidly assessing the genetic potential of a microbial community for a given  
446 set of reactions considering the relationships between genes and proteins. The two-step annotation  
447 of clusters of orthologs with a personalized ORAdb decreases the overall number of sequence  
448 alignments that need to be computed. User-specified constraints, such as the presence of  
449 transporter genes, further reduce the search space for putative microbial interactions. Users have  
450 substantial control over several steps of OrtSuite: from manual curation of ORAdb, custom  
451 sequence similarity cutoffs to the addition of constraints for inference of putative microbial  
452 interactions. The reduction of the search space of synergistic interactions by OrtSuite will also  
453 allow more comprehensive and computationally demanding tasks to be performed. Such as  
454 (Community) Flux Balance Analysis, which depend heavily on genome-scale metabolic models  
455 (Thommes *et al*, 2019; Ravikrishnan & Raman, 2021). As long as links between genes, proteins  
456 and reactions exist, the flexibility and easy usage of OrtSuite allow its application to the study of

457 any given ecosystem process. Nevertheless, assessing the functional potential of microbes is just  
458 the first step in deciphering synergistic microbial interactions. Linking the functional potential of  
459 microbial communities to transcriptomic or proteomic data will improve predictions and provide  
460 further insights into other types of microbial interactions.

461

## 462 **Materials and Methods**

### 463 **OrtSuite workflow**

464 The OrtSuite workflow consists of three main tasks performed using three bash commands (Figure  
465 1). The first task consists of generating a user-defined ortholog-reaction associated database  
466 (ORAdb) and collecting the gene-protein-reaction (GPR) rules. This task takes as input a list of  
467 KEGG identifiers which will be used to download all protein sequences associated with a set of  
468 reactions/pathway of interest. Next, all gene-protein-rules (GPRs) associated with each reaction  
469 will be downloaded from KEGG Modules. In the second task, OrtSuite employs OrthoFinder  
470 (Emms & Kelly, 2015) to generate ortholog clusters. This task takes as input a folder with the  
471 location of the genomic sequences. The third task consists of the functional annotation of species,  
472 identification of putative synergistic interspecies interactions, and generation of visual  
473 representations of the results.

474

### 475 **OrtSuite task 1 (green box, Figure 1) – User defined Ortholog-Reaction Association database** 476 **(ORAdb) and Gene-Protein-Reaction (GPR) rules file**

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

- 481 ● a location of the project folder where all results will be stored
- 482 ● a text file with a list of KEGG identifiers (one identifier per line)
- 483 ● the full path to the OrtSuite installation folder

484 The list of identifiers can be KEGG reactions (RID) (e.g. R11353, R02451), enzyme commission  
485 (EC) numbers (e.g. 1.3.7.8, 4.1.1.103) or KEGG ortholog identifiers (e.g. K07539, K20941). This  
486 file is used by OrtSuite to automatically retrieve the KEGG Ortholog identifiers (KO) (in case the  
487 identifiers provided are not KO identifiers) and to download all their associated protein sequences  
488 (Kanehisa *et al*, 2004). OrtSuite makes use of the python library *grequests*, which allows multiple  
489 queries in KEGG and subsequently decreases the time required for retrieving the ortholog  
490 associated sequences. The user-defined ORAdb will be composed of KO-specific sequence files  
491 in FASTA format associated with all reactions/enzymes of interest. Users also can manually add  
492 or edit the sets of reactions and the associated protein sequences in the ORAdb. This feature is  
493 particularly important since many reactions related to ecosystem processes are constantly being  
494 discovered and updated and might not be included in the latest version of KEGG. In addition,  
495 during the execution of the *DB\_construction.sh* OrtSuite performs the automated download of the  
496 gene-protein-reaction (GPR) rules from KEGG Modules. This feature is vital since enzymes can  
497 catalyze many reactions with a single (i.e., one protein) or multiple subunits (i.e., protein  
498 complexes). We advise users to manually curate the final table to guarantee accurate results despite  
499 the automated process. An example of the final GPR table is shown in the Supplementary data  
500 (Table S20).

501

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

503 The second task of OrtSuite, takes a set of protein sequences and generates clusters of orthologs.  
504 This set of protein sequences can originate from single isolates or from the complete set of protein  
505 sequences recovered from metagenomes or metagenome-assembled genomes. Indeed, using  
506 protein sequences from isolates, metagenome-assembled genomes, and co-culture experiments  
507 will benefit significantly from OrtSuite's reduction of the universe of potential microbial  
508 interactions based on the user-defined ORAdb. Orthology considers that phylogenetically distinct  
509 species can share functional similarities based on a common ancestor (Gabaldón & Koonin, 2013).  
510 Potentially, genes with similar functions will be grouped together. To perform this task, the  
511 OrtSuite pipeline uses OrthoFinder (Emms & Kelly, 2015). Three sequence aligners are available  
512 in OrthoFinder – DIAMOND (Buchfink *et al*, 2015), BLAST (Altschul *et al*, 1990) and MMSeqs2  
513 (Steinegger & Söding, 2017). DIAMOND (v0.9.22) is used by default due to its improved trade-  
514 off between execution time and sensitivity (Emms & Kelly, 2019). This task is performed by  
515 running the command *orthofinder* located in the installation folder of OrthoFinder. This command  
516 takes as input the full path to the folder containing the protein sequences to be clustered and the  
517 full path to the folder where results are to be stored.

518

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

520 The third task of OrtSuite consists of the assignment of functions to protein sequences contained  
521 in the ortholog clusters. Functional annotation of these clusters consists of a two-step process  
522 termed relaxed and restrictive search, respectively. The goal of the relaxed search is to decrease  
523 the number of alignments required to assign functions to sequences in the ortholog clusters. Here,  
524 50% of the sequences from each cluster are randomly selected and aligned to all sequences  
525 associated with each reaction present in the ORAdb. Only the e-value is considered during this  
526 stage. Ortholog clusters where e-values meet a user-defined threshold to sequences in the ORAdb  
527 proceed to the restrictive search. The default e-value was set to 0.001, as the main objective of the  
528 relaxed search is to capture as many sequences for annotation as possible while avoiding an  
529 excessive number of sequence alignments. In the restrictive search, all sequences in the  
530 transitioned ortholog clusters are aligned to all the sequences in the reaction set(s) present in the  
531 ORAdb to which they had a hit during the relaxed search. Again, the query sequence is only  
532 assigned to the function of a reference sequence if the e-value is below a determined threshold  
533 (default  $1e^{-9}$ ). Next, an additional filter is applied based on annotation bit score values (default 50).  
534 Although we established default values for the relaxed and restrictive search and bit score, the user  
535 can define the thresholds for all individual parameters.

536 The identification of putative interactions between species is based on all combinations of bacterial  
537 isolates with the genomic content to perform the user-defined pathway defined in the ORAdb. The  
538 input for this task consists of: (1) a binary table generated at the end of the functional annotation,  
539 which indicates the presence or absence of sequences annotated to each reaction in the ORAdb in  
540 each species (e.g., Supplementary Table S10); (2) a set of Gene-Protein-Reaction (GPR) rules for

541 all reactions considered (e.g., Supplementary data - Table S20); and (3) a user-defined tab-  
542 delimited file where the sets of reactions for complete pathways, subsets of reactions required to  
543 be performed by single species and transporter-associated genes (e.g., Supplementary data – Table  
544 S1) are described. Manual filtering can be performed to further reduce the vast amount of putative  
545 microbial interactions and increase confidence in the results. For example, results can be queried  
546 for known cross-feeding relationships between species or interactions that remove toxic  
547 compounds. Also, putative interactions can be removed if they are not biologically feasible. The  
548 user also may have an interest in assessing subsets of microbial interactions using specific criteria.  
549 Therefore, additional constraints can be applied to the putative microbial interactions, further  
550 reducing the search space. These include the degree of completeness of a pathway, the number of  
551 reactions expected to be performed by a single species or the presence or absence of transporter  
552 genes. Additionally, graphical network visualization is also produced during this step (Figure 3).  
553 The graphical network visualization is implemented in R using the packages visNetwork (v2.0.9),  
554 reshape2 (v1.4.3), and RColorBrewers (v1.1-2) but also requires the pandoc linux library.  
555 Graphical visualization was implemented with R v3.6 but also tested with v4.0. The visualization  
556 creates a HTML file that allows interactive network exploration and provides hyperlinks to KEGG  
557 if available.

558 All tasks - functional annotation, prediction of putative microbial interactions, and generation of  
559 graphical visualizations - are performed by running the script *annotate\_and\_predict.sh* included  
560 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  
561 predictions of individual species and combinations of species with the genetic potential to perform  
562 each defined pathway is stored in text files located in a folder termed "interactions".

563

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

565 We selected three alternative pathways involved in the conversion of benzoate to acetyl-CoA  
566 (BTA) to test the functional annotation and prediction of putative synergistic microbial interactions  
567 using OrtSuite (Supplementary data - Table S14). Two pathways consisted of benzoate's anaerobic  
568 degradation to acetyl-CoA via benzoyl-CoA differing only in the reactions required for  
569 transformation of glutaryl-CoA to crotonyl-CoA (hereafter, respectively, P1 and P2). P1 first  
570 converts glutaryl-CoA to glutaconyl-CoA and then to crotonoyl-CoA while P2 directly converts  
571 glutaryl-CoA to crotonoyl-CoA. One pathway consisted in the aerobic degradation of benzoate via  
572 catechol (hereafter P3). The complete number of reactions, enzymes, KO identifiers and KO-  
573 associated sequences in each alternative pathway is shown in the supplementary data  
574 (Supplementary data - Table S25).

575

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

577 To assess the performance of OrtSuite, we selected the transformation of benzoate to acetyl-CoA  
578 as a model pathway and a set of previously characterized species known to be involved in this  
579 pathway (Table 1). This set of species was divided in two groups. The first group contained

580 sequenced genomes of species whose ability to convert benzoate to acetyl-CoA has been  
581 demonstrated by KEGG (Kanehisa *et al*, 2004) and were selected as positive controls. These  
582 species were classified according to their genomic potential: complete, if all protein-encoding  
583 genes required for a BTA pathway were present in their genome or partial, if not all protein-  
584 encoding genes were present. The second group consisted of species that lacked all required  
585 protein-encoding genes and were selected as negative controls. In total, we selected 18 species as  
586 positive controls. Seven of them have the genetic potential to perform the alternative P2 pathway;  
587 eight have the genetic potential to perform alternative path P3 (positive controls); and none can  
588 completely perform the alternative path P1. To note that species *Thauera sp.* MZ1T has the genetic  
589 potential to perform P2 and P3 pathways. Four organisms were selected as negative controls. Using  
590 their genomes, we evaluated the performance of OrtSuite based on precision and recall rates for  
591 clustering of orthologs and the correct functional annotation of sequences. Also, a set of genomes  
592 from the species containing the genetic potential to degrade benzoate (*Burkholderia vietnamiensis*  
593 G4, *Azoarcus sp.* CIB and *Aromatoleum aromaticum* EbN1) were artificially mutated at the  
594 nucleotide level at different rates to determine how levels of point mutations in open reading  
595 frames (ORFs) affected clustering of ortholog groups.

596

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

598 In a study performed by Fetzer and collaborators (Fetzer *et al*, 2015), community biomass  
599 production of mono- and mixed-cultures was assessed in a medium containing benzoate. The  
600 authors used this data to infer potential species interactions. We processed this set of genomes with  
601 OrtSuite to determine the species' genetic potential to degrade benzoate, either individually or due  
602 to their interaction. Our results were compared to those obtained by Fetzer and collaborators and  
603 used to assess whether potential microbial interactions could be derived from their combined  
604 genetic potential.

605

### 606 **Evaluation of ortholog clustering**

607 We evaluated the clustering of orthologs by measuring the pairwise precision and recall. Clustering  
608 precision measures how many pairs of sequences associated with the same molecular function are  
609 grouped and is calculated by dividing the number of correctly clustered sequences by the total  
610 number of clustered sequences (Equation 1).

611

---

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

---

612

613 where, correctly clustered sequences refer to the pairs of sequences that share the same function  
614 and are clustered together and total number of clustered sequences refers to all pairs of sequences  
615 that are clustered together irrespective of sharing the same function.

616

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

620

---

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

---

621

622 where, correctly clustered sequences refer to the pairs of sequences that share the same function  
623 and are clustered together and total true sequence clusters refers to all pairs of sequences that have  
624 the same function.

625

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

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

634 *Burkholderia vietnamiensis* G4 was mutated on the gene K05783, *Azoarcus sp.* CIB on the gene  
635 K07537 and *Aromatoleum aromaticum* EbN1 on the gene K07538. Each gene was mutated at rates  
636 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  
637 original genome (e.g., *Burkholderia vietnamiensis* G4 strain K05783\_25, where “K05783” is the  
638 KEGG ortholog identifier and “25” is the rate of mutation). A total of 18 strains were generated  
639 (six in silico mutated strains per genome). The complete set of original and artificially mutated  
640 genomes is available in a compressed file (Supplementary data - Test\_genomes\_set.zip).

641

### 642 **Evaluation of functional annotation**

643 Functional annotation was evaluated based on the data collected from KEGG (Altschul *et al*,  
644 1990). Annotation performance is calculated by dividing the number of matching annotated  
645 sequences by the total number of annotations (Equation 3).

646

---

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

---

647

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

651

## 652 **Evaluation of microbial interaction predictions**

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

### 660 *Bacterial cultures and sequencing*

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

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

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

677

## 678 **Data Availability:**

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

- 680 ● The genomes used to test the workflow are available at National Centre for Biotechnology  
681 Information (<https://www.ncbi.nlm.nih.gov/>) under the accession identifiers [CP029389-](#)  
682 [CP029397](#), [GCF\\_000001735](#), [AP012304](#), [AP012305](#), [CP021731](#), [CP011072](#), [CP007785-](#)  
683 [CP007787](#), [CP000614-CP000621](#), [CP003230](#), [CP005996](#), [CP003108](#), [CR555306-](#)  
684 [CR5553068](#), [GCF\\_000225785](#), [LN997848-LN997849](#), [CP022989-CP022996](#), [CP024315](#),  
685 [AP012547](#), [CP022046-CP022047](#) and [CP001281-CP001282](#).
- 686 ● The genome assemblies used to predict interspecies interactions are available at National  
687 Centre for Biotechnology Information (<https://www.ncbi.nlm.nih.gov/>) with the study  
688 accession PRJEB38476: (<https://www.ncbi.nlm.nih.gov/bioproject/648592>).

- 689       • OrtSuite scripts: GitHub (<https://github.com/mdsufz/OrtSuite>).

690

691 **Acknowledgments:** We thank the early users of OrtSuite Sandra Silva, Felipe Côrrea for their  
692 help with debugging and for workflow suggestions. We also thank Diogo Lima and Emanuel  
693 Cunha for their assistance in the implementation of the script required to generate the Gene-  
694 Protein-Reaction (GPR) rules; and Nicole Steinbach for her work in the sequencing of the isolates  
695 used as the Fetzer test set. This work was funded by the Helmholtz Young Investigator grant VH-  
696 NG-1248 Micro ‘Big Data’.

697

698 **Author Contributions:** JS, OD, PS and UNR developed the concept of OrtSuite. JS, MG, AB and  
699 UNR developed the OrtSuite workflow. JCK developed the docker/conda installation packages.  
700 JS, MG, AB and UNR performed the benchmarks. CV provided information and data for defining  
701 benzoate to acetyl-CoA conversion pathways. RK sequenced bacterial isolates that were provided  
702 by AC. AB created the interactive network visualization module. JS and UNR wrote the  
703 manuscript. All authors read and commented on different versions of the manuscript and approved  
704 the final manuscript.

705

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

707

## 708 **References**

709 Altschul SF, Gish W, Miller W, Myers EW & Lipman DJ (1990) Basic local alignment search tool. *Journal of*  
710 *molecular biology* 215: 403–10

711 Bazire A, Diab F, Jebbar M & Haras D (2007) Influence of high salinity on biofilm formation and benzoate  
712 assimilation by *Pseudomonas aeruginosa*. *Journal of Industrial Microbiology and Biotechnology* 34: 5–  
713 8

714 Buchfink B, Xie C & Huson DH (2015) Fast and sensitive protein alignment using DIAMOND. *Nat Methods*  
715 12: 59–60

716 Clark RL, Connors BM, Stevenson DM, Hromada SE, Hamilton JJ, Amador-Noguez D & Venturelli OS  
717 (2021) Design of synthetic human gut microbiome assembly and butyrate production. *Nat Commun*  
718 12: 3254

719 Delmas E, Besson M, Brice M-H, Burkle LA, Riva GVD, Fortin M-J, Gravel D, Guimarães PR, Hembry DH,  
720 Newman EA, *et al* (2019) Analysing ecological networks of species interactions. *Biological Reviews*  
721 94: 16–36

722 Devanadera A, Vejarano F, Zhai Y, Suzuki-Minakuchi C, Ohtsubo Y, Tsuda M, Kasai Y, Takahata  
723 Y, Okada K & Nojiri H (2019) Complete Genome Sequence of an Anaerobic Benzene-Degrading  
724 Bacterium, *Azoarcus* sp. Strain DN11. *Microbiol Resour Announc* 8

725 Dong X & Strous M (2019) An Integrated Pipeline for Annotation and Visualization of Metagenomic  
726 Contigs. *Front Genet* 10

727 Emms DM & Kelly S (2015) OrthoFinder: solving fundamental biases in whole genome comparisons  
728 dramatically improves orthogroup inference accuracy. *Genome biology* 16: 157–157

729 Emms DM & Kelly S (2019) OrthoFinder: phylogenetic orthology inference for comparative genomics.  
730 *Genome Biology* 20: 238

731 Fetzer I, Johst K, Schäwe R, Banitz T, Harms H & Chatzinotas A (2015) The extent of functional redundancy  
732 changes as species' roles shift in different environments. *Proc Natl Acad Sci USA* 112: 14888–14893

733 Gabaldón T & Koonin EV (2013) Functional and evolutionary implications of gene orthology. *Nature*  
734 *Reviews Genetics* 14: 360–366

735 Gottstein W, Olivier BG, Bruggeman FJ & Teusink B (2016) Constraint-based stoichiometric modelling from  
736 single organisms to microbial communities. *Journal of The Royal Society Interface* 13: 20160627

737 Heirendt L, Arreckx S, Pfau T, Mendoza SN, Richelle A, Heinken A, Haraldsdóttir HS, Wachowiak J,  
738 Keating SM, Vlasov V, *et al* (2019) Creation and analysis of biochemical constraint-based models  
739 using the COBRA Toolbox v.3.0. *Nature Protocols* 14: 639–702

740 Hernández-Salmerón JE & Moreno-Hagelsieb G (2020) Progress in quickly finding orthologs as reciprocal  
741 best hits: comparing blast, last, diamond and MMseqs2. *BMC Genomics* 21: 741

742 Hu Y, Feng Y, Zhang X & Zong Z (2017) *Acinetobacter defluvii* sp. nov., recovered from hospital sewage.  
743 *International Journal of Systematic and Evolutionary Microbiology*, 67: 1709–1713

744 Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, von Mering C & Bork P (2017) Fast  
745 Genome-Wide Functional Annotation through Orthology Assignment by eggNOG-Mapper. *Mol*  
746 *Biol Evol* 34: 2115–2122

747 Jenul C, Sieber S, Daepfen C, Mathew A, Lardi M, Pessi G, Hoepfner D, Neuburger M, Linden A,  
748 Gademann K, *et al* (2018) Biosynthesis of fragin is controlled by a novel quorum sensing signal. *Nat*  
749 *Commun* 9: 1–13

750 Junghare M, Patil Y & Schink B (2015) Draft genome sequence of a nitrate-reducing, o-phthalate  
751 degrading bacterium, *Azoarcus* sp. strain PA01T. *Standards in Genomic Sciences* 10: 90

752 Jiang Y, Dong W, Xin F & Jiang M (2020) Designing Synthetic Microbial Consortia for Biofuel Production.  
753 *Trends Biotechnol* 38: 828–831

754 Kanehisa M, Goto S, Kawashima S, Okuno Y & Hattori M (2004) The KEGG resource for deciphering the  
755 genome. *Nucleic acids research* 32: D277–D280

756 Kanehisa M, Sato Y & Morishima K (2016) BlastKOALA and GhostKOALA: KEGG Tools for Functional  
757 Characterization of Genome and Metagenome Sequences. *J Mol Biol* 428: 726–731

758 Khanal A, Yu McLoughlin S, Kershner JP & Copley SD (2015) Differential Effects of a Mutation on the  
759 Normal and Promiscuous Activities of Orthologs: Implications for Natural and Directed Evolution.  
760 *Mol Biol Evol* 32: 100–108

761 Kitko RD, Cleeton RL, Armentrout EI, Lee GE, Noguchi K, Berkmen MB, Jones BD & Slonczewski JL (2009)  
762 Cytoplasmic Acidification and the Benzoate Transcriptome in *Bacillus subtilis*. *PLOS ONE* 4: e8255

763 Koonin EV (2005) Orthologs, Paralogs, and Evolutionary Genomics. *Annual Review of Genetics* 39: 309–338

764 Li L, Stoeckert CJ & Roos DS (2003) OrthoMCL: identification of ortholog groups for eukaryotic genomes.  
765 *Genome Res* 13: 2178–2189

766 Lee Y, Lee Y & Jeon CO (2019) Biodegradation of naphthalene, BTEX, and aliphatic hydrocarbons  
767 by *Paraburkholderia aromaticivorans* BN5 isolated from petroleum-contaminated soil. *Sci Rep* 9:  
768 860

769 Li W & Godzik A (2006) Cd-hit: a fast program for clustering and comparing large sets of protein or  
770 nucleotide sequences. *Bioinformatics* 22: 1658–1659

771 Locey KJ & Lennon JT (2016) Scaling laws predict global microbial diversity. *PNAS*: 201521291

772 Lynch M (2010) Evolution of the mutation rate. *Trends Genet* 26: 345–352

773 Lyu Z, Shao N, Akinyemi T & Whitman WB (2018) Methanogenesis. *Curr Biol* 28: R727–R732

774 Madden T (2003) The BLAST Sequence Analysis Tool National Center for Biotechnology Information (US)

775 Maestre FT, Castillo-Monroy AP, Bowker MA & Ochoa-Hueso R (2012) Species richness effects on  
776 ecosystem multifunctionality depend on evenness, composition and spatial pattern. *Journal of*  
777 *Ecology* 100: 317–330

778 Mendes LW, Raaijmakers JM, de Hollander M, Mendes R & Tsai SM (2018) Influence of resistance breeding  
779 in common bean on rhizosphere microbiome composition and function. *ISME J* 12: 212–224

780 Messina E, Denaro R, Crisafi F, Smedile F, Cappello S, Genovese M, Genovese L, Giuliano L, Russo  
781 D, Ferrer M, *et al* (2016) Genome sequence of obligate marine polycyclic aromatic hydrocarbons-  
782 degrading bacterium *Cycloclasticus* sp. 78-ME, isolated from petroleum deposits of the sunken  
783 tanker Amoco Milford Haven, Mediterranean Sea. *Marine Genomics* 25: 11–13

784 Meyer-Cifuentes I, Fiedler S, Müller JA, Kappelmeyer U, Mäusezahl I & Heipieper HJ (2017) Draft  
785 Genome Sequence of Magnetospirillum sp. Strain 15-1, a Denitrifying Toluene Degradator Isolated  
786 from a Planted Fixed-Bed Reactor. *Genome Announc* 5

787

788 Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL, Tosatto SCE, Paladin L,  
789 Raj S, Richardson LJ, *et al* (2021) Pfam: The protein families database in 2021. *Nucleic Acids Research*  
790 49: D412–D419

791 Morin M, Pierce EC & Dutton RJ (2018) Changes in the genetic requirements for microbial interactions with  
792 increasing community complexity. *eLife* 7: e37072

793 Mrozik A & Labuzek S (2002) A comparison of biodegradation of phenol and homologous  
794 compounds by *Pseudomonas vesicularis* and *Staphylococcus sciuri* strains. *Acta Microbiol Pol* 51:  
795 367–378

796 Mulder CPH, Uliassi DD & Doak DF (2001) Physical stress and diversity-productivity relationships: The  
797 role of positive interactions. *PNAS* 98: 6704–6708

798 Ng PC & Henikoff S (2006) Predicting the effects of amino acid substitutions on protein function. *Annu Rev*  
799 *Genomics Hum Genet* 7: 61–80

800 O’Sullivan LA, Weightman AJ, Jones TH, Marchbank AM, Tiedje JM & Mahenthiralingam E (2007)  
801 Identifying the genetic basis of ecologically and biotechnologically useful functions of the  
802 bacterium *Burkholderia vietnamiensis*. *Environmental Microbiology* 9: 1017–1034

803 Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, Hugenholtz P & Tyson GW  
804 (2017) Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of  
805 life. *Nat Microbiol* 2: 1533–1542

806 Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, Beghini F, Manghi P, Tett A, Ghensi P, *et*  
807 *al* (2019) Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes  
808 from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* 176: 649-662.e20

809 Pearson WR (2013) An Introduction to Sequence Similarity (“Homology”) Searching. *Curr Protoc*  
810 *Bioinformatics* 0 3

811 Peng T, Luo A, Kan J, Liang L, Huang T & Hu Z (2018) Identification of A Ring-Hydroxylating  
812 Dioxygenases Capable of Anthracene and Benz[a]anthracene Oxidization from *Rhodococcus* sp.  
813 P14. *MMB* 28: 183–189

814 Prestat E, David MM, Hultman J, Taş N, Lamendella R, Dvornik J, Mackelprang R, Myrold DD, Jumpponen  
815 A, Tringe SG, *et al* (2014) FOAM (Functional Ontology Assignments for Metagenomes): a Hidden  
816 Markov Model (HMM) database with environmental focus. *Nucleic Acids Res* 42: e145

817 Rabus R, Boll M, Heider J, Meckenstock RU, Buckel W, Einsle O, Ermler U, Golding BT, Gunsalus  
818 RP, Kroneck PMH, *et al* (2016) Anaerobic Microbial Degradation of Hydrocarbons: From  
819 Enzymatic Reactions to the Environment. *MMB* 26: 5–28

820 Ravikrishnan A & Raman K (2021) Unraveling microbial interactions in the gut microbiome. *bioRxiv*:  
821 2021.05.17.444446

822 Raynaud X & Nunan N (2014) Spatial Ecology of Bacteria at the Microscale in Soil. *PLOS ONE* 9: e87217

823 Robertson WJ, Franzmann PD & Mee BJ (2000) Spore-forming, Desulfosporosinus-like sulphate-  
824 reducing bacteria from a shallow aquifer contaminated with gasoline. *Journal of Applied*  
825 *Microbiology* 88: 248–259

826 Roh SW, Abell GCJ, Kim K-H, Nam Y-D & Bae J-W (2010) Comparing microarrays and next-generation  
827 sequencing technologies for microbial ecology research. *Trends Biotechnol* 28: 291–299

828 Seemann T (2014) Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30: 2068–2069

829 Sharma B & Shukla P (2020) Designing synthetic microbial communities for effectual bioremediation: A  
830 review. *Biocatalysis and Biotransformation* 38: 405–414

831 Singleton CM, Petriglieri F, Kristensen JM, Kirkegaard RH, Michaelsen TY, Andersen MH, Kondrotaitė Z,  
832 Karst SM, Dueholm MS, Nielsen PH, *et al* (2021) Connecting structure to function with the recovery  
833 of over 1000 high-quality metagenome-assembled genomes from activated sludge using long-read  
834 sequencing. *Nat Commun* 12: 2009

835 Slade EM, Kirwan L, Bell T, Philipson CD, Lewis OT & Roslin T (2017) The importance of species identity  
836 and interactions for multifunctionality depends on how ecosystem functions are valued. *Ecology*  
837 98: 2626–2639

838 Sperfeld M, Diekert G & Studenik S (2019) Anaerobic aromatic compound degradation in  
839 *Sulfuritalea hydrogenivorans* sk43H. *FEMS Microbiol Ecol* 95

840 Steinegger M & Söding J (2017) MMseqs2 enables sensitive protein sequence searching for the analysis of  
841 massive data sets. *Nat Biotechnol* 35: 1026–1028

842 Suvorova IA & Gelfand MS (2019) Comparative Genomic Analysis of the Regulation of Aromatic  
843 Metabolism in Betaproteobacteria. *Front Microbiol* 10

844 Tal O, Selvaraj G, Medina S, Ofaim S & Freilich S (2020) NetMet: A Network-Based Tool for Predicting  
845 Metabolic Capacities of Microbial Species and their Interactions. *Microorganisms* 8: 840

846 Thommes M, Wang T, Zhao Q, Paschalidis IC & Segrè D (2019) Designing Metabolic Division of Labor in  
847 Microbial Communities. *mSystems* 4

848 Tully BJ, Graham ED & Heidelberg JF (2018) The reconstruction of 2,631 draft metagenome-assembled  
849 genomes from the global oceans. *Scientific Data* 5: 170203

850 Valderrama JA, Durante-Rodríguez G, Blázquez B, García JL, Carmona M & Díaz E (2012) Bacterial  
851 Degradation of Benzoate: CROSS-REGULATION BETWEEN AEROBIC AND ANAEROBIC  
852 PATHWAYS. *J Biol Chem* 287: 10494–10508

853 Wang B, Lai Q, Cui Z, Tan T & Shao Z (2008) A pyrene-degrading consortium from deep-sea  
854 sediment of the West Pacific and its key member *Cycloclasticus* sp. P1. *Environmental Microbiology*

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

857

858

859 **Figure legends**

860 **Figure 1 - OrtSuite workflow.**

861 OrtSuite takes a text file containing a list of identifiers for each reaction in the pathway of interest  
862 supplied by the user to retrieve all protein sequences from KEGG Orthology and are stored in  
863 ORAdb. Subsequently, the same list of identifiers is used to obtain the Gene-Protein-Reaction  
864 (GPR) rules from KEGG Modules (**Task 1**). Protein sequences from samples supplied by the user  
865 are clustered using OrthoFinder (**Task 2**). In **Task 3**, the functional annotation, identification of  
866 putative synergistic species interactions and graphical visualization of the network are performed.  
867 The functional annotation consists of a two-stage process (relaxed and restrictive search). Relaxed  
868 search performs sequence alignments between 50% of randomly selected sequences from each  
869 generated cluster. Clusters whose representative sequences share a minimum E-value of 0.001 to  
870 sequences in the reaction set(s) in ORAdb continue to the restrictive search. Here, all sequences  
871 from the cluster are aligned to all sequences in the corresponding reaction set(s) to which they had  
872 a hit (default E-value =  $1e^{-9}$ ). Next, the annotated sequences are further filtered to those with a bit  
873 score greater than 50 and are used to identify putative microbial interactions based on their  
874 functional potential. Constraints can also be added to reduce the search space of microbial  
875 interactions (e.g., subsets of reactions required to be performed by single species, transport-related  
876 reactions). Additionally, an interactive network visualization of the results is produced and  
877 accessed via a HTML file.

878 **Figure 2 - Mapping of the Fetzer genome set to benzoate pathways.**

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

883 **Figure 3 – Example of the interactive network visualization included on OrtSuite results.**

884 **(A)** The complete network with species is colored by reaction. **(B)** Species can be highlighted for  
885 simple identification. **(C)** Tooltips on reaction link out the KEGG if the reaction identifier is  
886 given.

Figure 1 **Step 1: ORA database generation and**

**Step 2: 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



**Step 3: Ortholog clustering**

- Protein sequence data [FASTA]
- Species id

Orthofinder

Cluster of orthologs

**Step 4: Functional annotation**

Relaxed search

DIAMOND BLAST

Reduced cluster  
of orthologs

Restrictive search

DIAMOND BLAST

Annotated sequences

**Step 5: Putative interactions  
and exploration**

**Interactive  
network  
visualization**

**Functional  
potential [CSV]**

**Putative species  
interactions [CSV]**

**Task 1**

**Task 2**

**Task 3**

Optional step  
----->

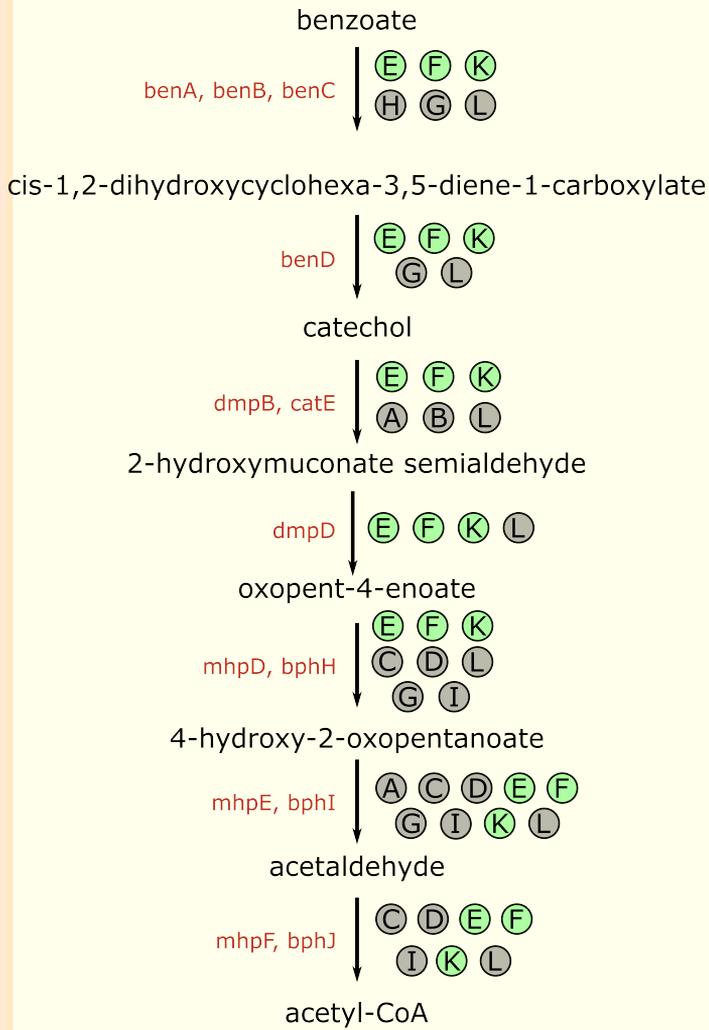
Input [format]

Key tools

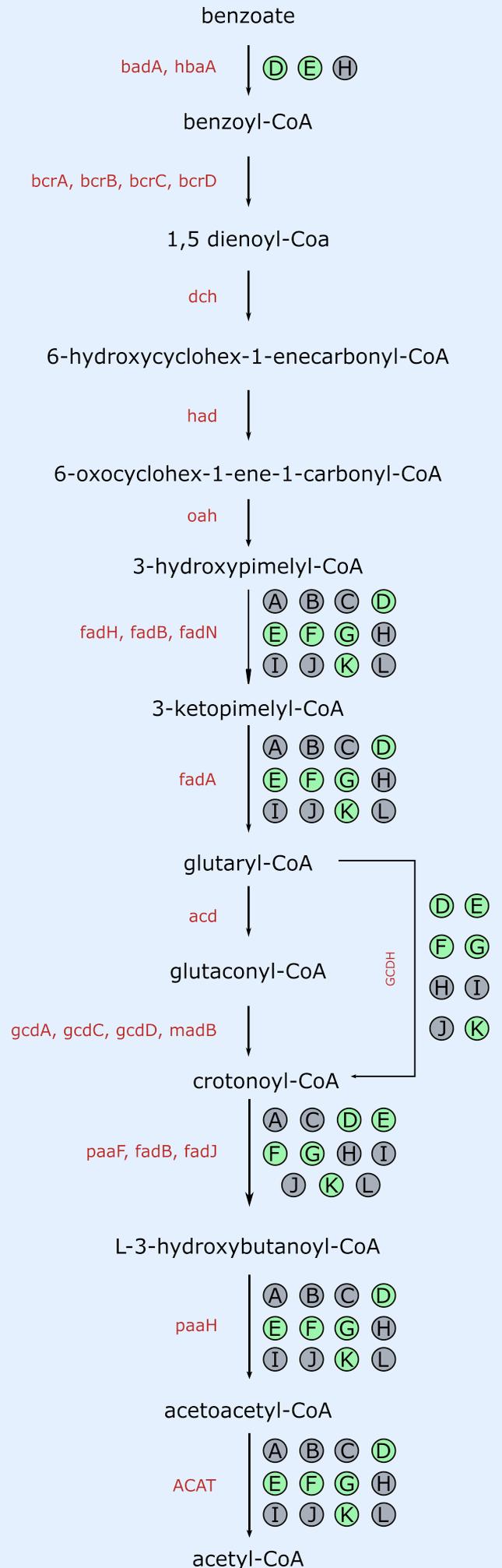
Output [format]

Figure 2

## Aerobic



## Anaerobic

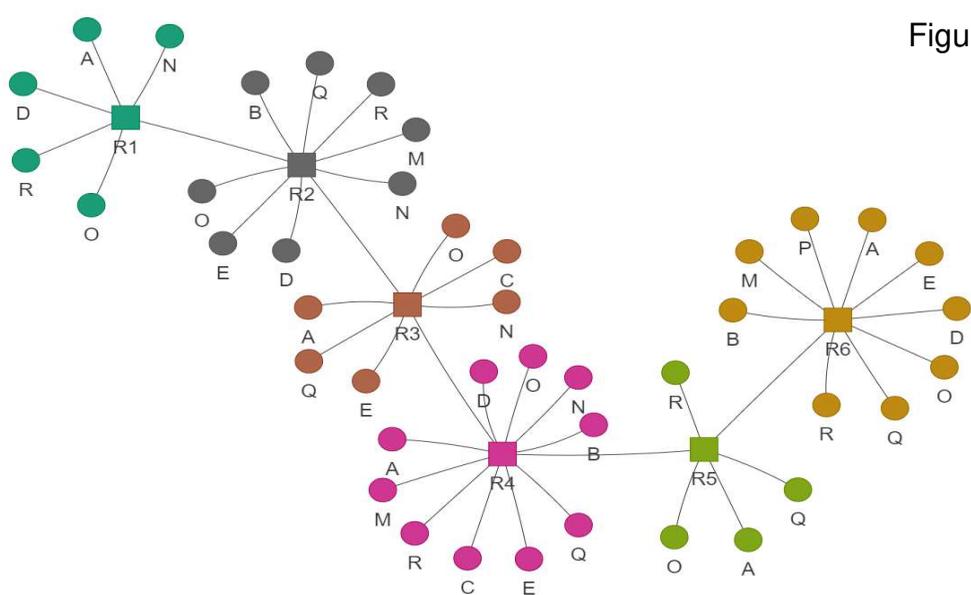
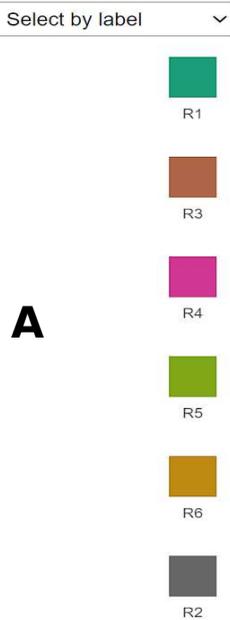


### Legends

- (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*

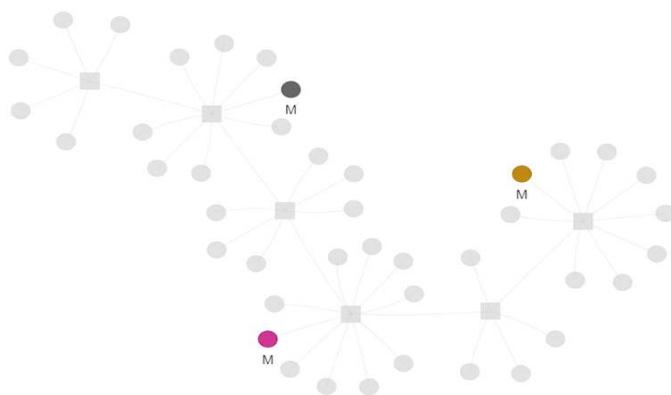
**Species with biomass growth in medium containing benzoate (Fetzer et al., 2015)**

Figure 3



**A**

**B**



**C**

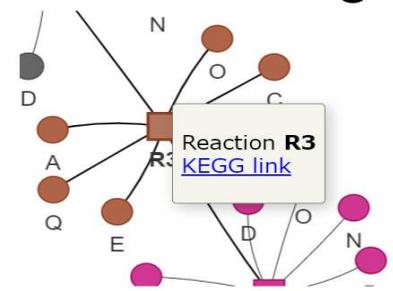


Table 1 - Species names, strain and abbreviation codes of the genomes used to validate OrtSuite (Supplementary data - Test\_genome\_set). The genomic potential, based on KEGG database, to completely encode all proteins involved in a BTA pathway is identified in the column “BTA pathway” (P1 – Anaerobic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

Name and strain	Abbreviation code	KEGG id	BTA pathway	Accession number	Ref.
<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)

---

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

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SaraivaetalSupplementarydataTableS3.xls](#)
- [SaraivaetalSupplementarydataTableS16.xls](#)
- [SaraivaetalSupplementarydataTableS20.xls](#)
- [SaraivaetalSupplementarydataTableS22.xls](#)
- [SaraivaetalSupplementarydataTableS24.xls](#)
- [SaraivaetalSupplementarydataTableS6.xls](#)
- [SaraivaetalSupplementarydataTableS7.xls](#)
- [SaraivaetalSupplementarydataTableS8.xls](#)
- [SaraivaetalSupplementarydataTableS9.xls](#)
- [SaraivaetalSupplementarydataTableS10.xls](#)
- [SaraivaetalSupplementarydataTableS11.xls](#)
- [SaraivaetalSupplementarydataTableS12.xls](#)
- [SaraivaetalSupplementarydataTableS13.xls](#)
- [SaraivaetalSupplementarydataTableS1.docx](#)
- [SaraivaetalSupplementarydataTableS2.docx](#)
- [SaraivaetalSupplementarydataTableS4.docx](#)
- [SaraivaetalSupplementarydataTableS5.docx](#)
- [SaraivaetalSupplementarydataTableS14.docx](#)
- [SaraivaetalSupplementarydataTableS15.docx](#)
- [SaraivaetalSupplementarydataTableS17.docx](#)
- [SaraivaetalSupplementarydataTableS18.docx](#)
- [SaraivaetalSupplementarydataTableS19.docx](#)
- [SaraivaetalSupplementarydataTableS21.docx](#)
- [SaraivaetalSupplementarydataTableS23.docx](#)
- [SaraivaetalSupplementarydataTableS25.docx](#)