

1 **OrtSuite – a flexible pipeline for annotation of ecosystem processes and prediction of putative**  
2 **microbial interactions**

3

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

6

7

8 <sup>1</sup>Department of Environmental Microbiology, Helmholtz Centre for Environmental Research –  
9 UFZ, Permoserstraße 15, 04318 Leipzig, Germany

10 <sup>2</sup>Centre of Biological Engineering, University of Minho, R. da Universidade, 4710-057, Braga,  
11 Portugal

12 <sup>3</sup>Institute of Biology, Leipzig University, Augustusplatz 10, 04109 Leipzig, Germany

13 <sup>4</sup>German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Deutscher Platz  
14 5e, 04103 Leipzig, Germany

15 <sup>5</sup>Bioinformatics Group, Department of Computer Science, Interdisciplinary Center for  
16 Bioinformatics, and Competence Center for Scalable Data Services and Solutions Dresden/Leipzig,  
17 University of Leipzig, PF 100920, 04009 Leipzig, Germany

18 <sup>6</sup>Max Planck Institute for Mathematics in the Sciences, Inselstraße 22, 04103 Leipzig, Germany

19 <sup>7</sup>Institute for Theoretical Chemistry, University of Vienna, Währinger Str. 17, 1090 Wien, Austria

20 <sup>8</sup>Facultad de Ciencias, Universidad Nacional de Colombia, Cra 45, Bogotá, Colombia

21 <sup>9</sup>Santa Fe Institute, 1399 Hyde Park Rd, Santa Fe, NM 87501, U.S.A.

22

23 Author email addresses: Joao Saraiva ([joao.saraiva@ufz.de](mailto:joao.saraiva@ufz.de)), Marta Gomes

24 ([martalopesgomes@hotmail.com](mailto:martalopesgomes@hotmail.com)), Rene Kallies ([rene.kallies@ufz.de](mailto:rene.kallies@ufz.de)), Carsten Vogt

25 ([carsten.vogt@ufz.de](mailto:carsten.vogt@ufz.de)), Antonis Chatzinotas ([antonis.chatzinotas@ufz.de](mailto:antonis.chatzinotas@ufz.de)), Peter Stadler

26 ([studla@bioinf.uni-leipzig.de](mailto:studla@bioinf.uni-leipzig.de)), Oscar Dias ([odias@ceb.uminho.pt](mailto:odias@ceb.uminho.pt)), Ulisses Nunes da Rocha

27 ([ulisses.rocha@ufz.de](mailto:ulisses.rocha@ufz.de))

28

29 \*Corresponding author: [ulisses.rocha@ufz.de](mailto:ulisses.rocha@ufz.de)

30

31

32

33 **Abstract**

34 **Background**

35 The exponential increase in high-throughput sequencing data and the development of computational  
36 sciences and bioinformatics pipelines has advanced our understanding of microbial community  
37 composition and functional and genetic diversity in complex ecosystems. Despite these advances,  
38 the identification of microbial interactions from genomic and metagenomics data remains a major  
39 bottleneck. To address this challenge, we present OrtSuite, a flexible workflow to predict putative  
40 microbial interactions based on genomic content.

41 **Results**

42 OrtSuite combines ortholog clustering strategies with genome annotation based on a user-defined  
43 set of functions allowing for hypothesis-driven data analysis such as assessing microbial  
44 interactions in specific ecosystems. OrtSuite allows users to install and run all workflow  
45 components and analyze the generated outputs using a simple pipeline consisting of 23 bash  
46 commands and one R command. Annotation is based on a two-stage process. First, only a subset of  
47 sequences from each ortholog cluster are aligned to all sequences in the Ortholog-Reaction  
48 Association database (ORAdb). Next, all sequences from clusters that meet a user-defined identity  
49 threshold are aligned to all sequence sets in ORAdb to which they had a hit. This approach results  
50 in a decrease in time needed for functional annotation. Further, OrtSuite identifies putative  
51 interspecies interactions based on their individual genomic content based on constraints given by  
52 the users. Additional control is afforded to the user at several stages of the workflow: 1) The  
53 construction of ORAdb only needs to be performed once for each specific process also allowing  
54 manual curation; 2) The identity and sequence similarity thresholds used during the annotation  
55 stage can be adjusted; and 3) Constraints related to pathway reaction composition and known  
56 species contributions to ecosystem processes can be defined.

57 **Conclusions**

58 OrtSuite is an easy to use workflow that allows for rapid functional annotation based on a user  
59 curated database. Further, this novel workflow allows the identification of interspecies interactions  
60 through user-defined constraints. Due to its low computational demands, for small datasets (e.g.  
61 maximum 100 genomes) OrtSuite can run on a personal computer. For larger datasets (> 100  
62 genomes), we suggest the use of computer clusters. OrtSuite is an open-source software available at  
63 <https://github.com/mdsufz/OrtSuite>.

64

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

67

## 68 **Background**

69 In environments where microorganisms play a key role, ecosystem processes are dependent on the  
70 functional potential of the associated microbiome and universe of possible interspecies interactions  
71 within the microbiome members [1,2]. For example, in environments rich in methane, microbial  
72 communities are dominated by species with genes encoding proteins involved in methanogenesis  
73 [3]. Soil microbes, especially those in the rhizosphere are genetically adapted to perform functions  
74 such as pathogen resistance and tolerance to stress [4]. Natural ecosystems are populated by an  
75 enormous number of microbes [5]. For example, soil environments can contain more than  $10^{10}$   
76 species per gram of soil making a global search for interspecies interactions unfeasible [6]. The  
77 exponential increase in high-throughput sequencing data and the development of computational  
78 sciences and bioinformatics pipelines has advanced our understanding of microbial community  
79 composition and distribution in complex ecosystems [7]. This knowledge increased our ability to  
80 reconstruct and functionally characterize genomes in complex communities, for example by the  
81 recovery of metagenome-assembled genomes (MAGs) [8–10]. While tools have been developed to  
82 improve the reconstruction of MAGs, the same cannot be said for predicting interspecies  
83 interactions [11]. The challenge of the latter increases due to the multitude of potential interactions

84 not only between species in microbial communities but also between microbes and their hosts (e.g.,  
85 plants, animals and microeukaryotes) [12]. In order to decrease complexity and facilitate analysis,  
86 the search of interactions can be limited to groups of organisms (e.g. microbe-microbe or host-  
87 microbe).

88 To improve our understanding on how microorganisms work within their communities and  
89 contribute to different ecosystem processes, it is essential to assess their functional capacity.  
90 Accurate annotation of gene function is essential to predict, from sequencing data, ecosystem  
91 processes potentially performed by microbial communities, particularly in cases where an  
92 ecosystem process is performed by the synergy of two or more species.

93 Simple methods for the annotation of genomes rely, for instance, on the search for  
94 homologous sequences. Computational tools such as BLAST [13] and DIAMOND [14] allow the  
95 comparison of nucleotide or protein sequences to those present in databases. These approaches  
96 allow inferring the function of uncharacterized sequences from their homologous pairs whose  
97 function is already known. The degree of confidence in the assignment of biological function is  
98 increased if this has been validated by, for example, experimental data. Orthologs are homologous  
99 sequences that descend from the same ancestor separated after a speciation event retaining the same  
100 function [15]. Approaches based on orthology are increasingly used for genome-wide functional  
101 annotation [16]. For example, OrthoMCL [17], CD-HIT [18] and OrthoFinder [19,20] are just a few  
102 tools that identify homologous relationships between sequences using orthology.

103 OrthoFinder has been shown to be more accurate than several other orthogroup inference  
104 methods since it takes into account gene length in the detection of ortholog groups by introducing a  
105 score transformation step [19]. However, OrthoFinder, due to its all-versus-all sequence alignment  
106 approach, requires intensive computational resources resulting in long running times when using  
107 large data sets for clustering. Nevertheless, integrating ortholog clustering tools with functional  
108 annotation strategies can improve our ability to functionally characterize microbial communities  
109 and reduce the time needed to obtain results.

110 Genome-based modelling approaches have routinely been used to study single organisms as  
111 well as microbial communities [21]. For example, constraint-based models are highly employed in  
112 the study and prediction of metabolic networks [22]. These models are generated upon the premise  
113 that any given function is feasible as long as the protein-encoding gene is present. Although species  
114 may lack the genetic potential to perform all functions necessary to survive in a given ecosystem, in  
115 nature microbes do not exist in isolation and usually benefit from their interaction with other  
116 species. By assessing the genomic content of individual species we are able to identify groups of  
117 microbes whose combined content may account for complete ecosystem functioning.

118 In this study, we developed OrtSuite; a pipeline to perform rapid and accurate functional  
119 annotation of various microbial species based on the detection of orthologs. This pipeline integrates  
120 the use of a user-defined database – Ortholog-Reaction Association database (ORAdb) – with up-  
121 to-date ortholog clustering tools. Additionally, OrtSuite allows the search for putative microbial  
122 interactions by calculating the combined genomic potential of individual species in specific  
123 ecosystem processes. We evaluate this pipeline using a clearly defined set of reactions involved in  
124 the well-described benzoate-to-Acetyl-CoA (BTA) conversion. Further, we used this pipeline to  
125 functionally characterize a set of known benzoate degraders. OrtSuite’s ability to identify putative  
126 interspecies interactions was evaluated on species whose interactions have been previously  
127 predicted under controlled conditions [23].

128

## 129 **Implementation**

130 The OrtSuite workflow, excluding installation, consists of four main tasks performed by the use of  
131 10 bash commands and one R command (Figure 1): (1) generation of a user defined ortholog-  
132 reaction associated database (ORAdb) consisting of protein sequences associated with a set of  
133 reactions/pathway of interest; (2) generation of ortholog clusters; (3) functional annotation of  
134 sequences in the generated ortholog clusters using ORAdb, consisting of a two-stage process termed  
135 relaxed and restrictive search; and, (4) identification of putative microbial interactions based on the

136 genomic content. To assess the performance of OrtSuite, we selected the transformation of benzoate  
137 to acetyl-CoA as a model pathway and a set of previously characterized species known to be  
138 involved in this pathway. This set of species was divided into two groups. The first group contained  
139 species known to possess all or a subset of the necessary protein encoding genes involved in the  
140 model pathway and were selected as positive controls. The second group consisted of species  
141 known to lack all required protein encoding genes and were selected as negative controls. Using  
142 their genomes, we evaluated the performance of OrtSuite based on precision and recall rates for  
143 both, clustering and functional annotation of sequences. Also, a set of genomes from the species  
144 containing the genetic potential to degrade benzoate were artificially mutated at the nucleotide level  
145 at different rates in order to determine how different open reading frames (ORFs) need to be in  
146 order to be clustered in a different ortholog group. Putative microbial interactions were assessed  
147 based on known species interactions and their potential to degrade benzoate individually or as a  
148 result of their interaction [22].

149

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

151 We selected three alternative pathways involved in the conversion of benzoate to acetyl-CoA  
152 (BTA) (Figure 2). Two pathways consisted in the anaerobic degradation of benzoate to acetyl-CoA  
153 via benzoyl-CoA differing only in the reactions required for transformation of glutaryl-CoA to  
154 crotonyl-CoA (hereafter, respectively, BTA\_P1 and BTA\_P2). BTA\_P1 first converts glutaryl-CoA  
155 to glutaconyl-CoA and then to crotonoyl-CoA while BTA\_P2 directly converts glutaryl-CoA to  
156 crotonoyl-CoA. One pathway consisted in the aerobic degradation of benzoate via catechol  
157 (hereafter BTA\_P3). The complete number of reactions, enzymes, KO identifiers and KO-  
158 associated sequences in each alternative pathway is shown in the supplementary data (Additional  
159 file 1: Table S1).

160

#### 161 *Selection of species with known genetic potential to convert benzoate to acetyl-CoA*

162 We designed OrtSuite to identify putative microbial interactions based on the combined genomic  
163 potential of species in a given set of enzymatic reactions. To test this feature, we used previously  
164 sequenced genomes of species whose ability to convert benzoate to acetyl-CoA has been  
165 demonstrated (Additional file 2: Table S2). These species were classified according to their  
166 genomic potential: *complete*, if all protein encoding genes required for a BTA pathway were present  
167 in their genome; *partial*, if not all protein encoding genes were present; and *absent*, if none of the  
168 required protein encoding genes were present in their genomes. In total, we selected 18 species.  
169 Seven of them have the genetic potential to perform the alternative BTA\_P2 pathway; eight have  
170 the genetic potential to perform alternative path 3 (positive controls); and, none able to completely  
171 perform the alternative path 1 (for BTA pathways, see Figure 2). As negative controls, four of the  
172 selected organisms were unable to perform completely any of the BTA pathways. Note that species  
173 *Thauera sp. MZ1T* has the genetic potential to perform BTA\_P2 and BTA\_P3 pathways.

174

175 *ORAdb – User defined ortholog-reaction association database (OrtSuite step 1)*

176 The ORAdb used for functional annotation consists of sets of protein sequences involved in the  
177 enzymatic reactions that compose a pathway/function of interest defined by the user. This database  
178 is generated by running the command *download\_kos* in OrtSuite. The user provides a text file with  
179 a list of identifiers (e.g. reaction (Additional file 3: *Example\_reaction\_list*), enzyme commission  
180 (EC) number (Additional file 4: *Example\_ec\_list*) or KEGG orthologs (Additional file 5:  
181 *Example\_ko\_list*). This file is used by OrtSuite to automatically retrieve the KEGG Ortholog  
182 identifiers (KO) and to download all their associated protein sequences [24]. OrtSuite makes use of  
183 the python library *grequests* which allows multiple queries in KEGG subsequently decreasing the  
184 time required for retrieving the ortholog associated sequence. The user-defined ORAdb will be  
185 composed of KO-specific sequence files in FASTA format associated with all reactions of interest.  
186 Users also have the opportunity to manually add or edit the sets of reactions and the associated  
187 protein sequences in the final ORAdb. This feature is of particular importance since many reactions

188 associated with ecosystem processes are constantly being discovered and updated and might not be  
189 included in the latest version of KEGG. In addition, the OrtSuite pipeline also allows the automated  
190 download of the gene-protein-reaction (GPR) rules from KEGG. Here, the user only needs to  
191 provide the location of the folder where the sequences of orthologs for each KO is stored. This  
192 feature is of particular importance since many reactions can be catalyzed by enzymes with a single  
193 (i.e., one protein) or multiple subunits (i.e., protein complexes).

194

#### 195 *Generation of protein ortholog clusters (OrtSuite step 2)*

196 The second stage of OrtSuite, takes a set of protein sequences and generates clusters of orthologs.  
197 This set of protein sequences can originate from single isolates or from the complete set of protein  
198 sequences recovered from metagenomes or metagenome-assembled genomes. Indeed, the use of  
199 protein sequences from isolates, metagenome-assembled genomes and co-culture experiments will  
200 benefit greatly from OrtSuite's reduction of the universe of potential microbial interactions based  
201 on a user defined ORAdb. Orthology considers that phylogenetically distinct species can share  
202 functional similarities based on a common ancestor [25]. Potentially, genes with equal function will  
203 be grouped together. To perform this task the OrtSuite pipeline uses OrthoFinder [19]. Two  
204 sequence aligners are available in OrthoFinder – DIAMOND [14] and BLAST [13]. DIAMOND is  
205 used by default due to its improved trade-off between execution time and sensitivity [20].

206

#### 207 *Functional annotation of ortholog clusters (OrtSuite step 3)*

208 The third stage of OrtSuite consists in the assignment of functions to protein sequences contained in  
209 the ortholog clusters. Functional annotation of these clusters consists of a two-step process termed  
210 relaxed and restrictive search, respectively. The goal of the relaxed search is to decrease the number  
211 of alignments required to assign functions to ortholog clusters. Here, 10% of the total number of  
212 sequences from each cluster are randomly selected and aligned to all sequences associated to each  
213 reaction present in the ORAdb. During this stage only sequence identity is evaluated. Ortholog

214 clusters where protein sequence identity meets a user-defined threshold to sequences in the ORAdb  
215 proceed to the restrictive search. The default value was set to 50% identity, as the main objective of  
216 the relaxed search is to capture as many sequences for annotation as possible whilst maintaining a  
217 reasonable number of sequence alignments. In the restrictive search, all sequences in the  
218 transitioned ortholog clusters are aligned to all the sequences in the reaction set(s) present in the  
219 ORAdb to which they had a hit during the relaxed search. The query sequence is only assigned to  
220 the function of a reference sequence if the combined score of sequence identity, positive matches,  
221 query sequence coverage and reference sequence coverage meets a user-defined threshold  
222 (Equation 1). The default value was set to 90% as an excessive restrictive search score can lead to a  
223 high number of sequences without annotation. Although we established default values for the  
224 relaxed and restrictive search, the user has the option to define the thresholds for all individual  
225 parameters. An explanation of how the different parameters might affect the results can be found in  
226 <https://github.com/mdsufz/OrtSuite/README.txt>.

227

228 *Equation 1*

$$229 \quad score = \frac{\%identity + \%positive\ matches + \%query\ seq.\ coverage + \%reference\ seq.\ coverage}{4}$$

230

231 where, percentage of identity (*%identity*) measures the exact alignment of amino acids between two  
232 sequences in a range of alignment. Here, sequences of interest with unknown function are provided  
233 by the user and aligned to the sequences in the user defined ORA database. The default identity  
234 value is set to 95% to capture as many perfect amino acid alignments as possible without being too  
235 stringent. Many mutations consistently occur in related protein groups that do not necessarily alter  
236 the function of a gene. These mutations alter the sequence composition and can result in low  
237 identity scores. The use of scoring matrices [26] is an example of measuring similarity between two

238 or more sequences when sequences are not identical. The percentage of positive matches (*%positive*  
239 *matches*) using scoring matrices reflects the number of sequences that, although not identical, are  
240 similar to each other in a range of alignment. The default value is set to 99% only returning  
241 sequences where a single amino acid is different in the whole alignment. The percentage of the  
242 query sequence that is covered (*%query seq.coverage*) refers to the percentage of the query  
243 sequence that overlaps with the reference sequence in the alignment. The default value is set 90%  
244 so that the entire query sequence is not required during the alignment. The percentage of reference  
245 sequence that is covered (*%reference seq.coverage*) refers to the percentage of the sequence in the  
246 ORA database that overlaps with the query sequence in the alignment. The default value is set 90%  
247 so that the entire reference sequence is not required during the alignment.

248

#### 249 *Defining putative microbial interactions (OrtSuite step 4)*

250 The identification of putative interactions between species is based on all combinations of bacterial  
251 isolates with the genomic content to perform the user-defined pathway defined in the ORAdb. The  
252 input for this task consists of: (1) a binary table generated at the end of the restrictive search for  
253 functional annotation which indicates the presence or absence of sequences annotated to each  
254 reaction in the ORA database in each species (Additional file 6: Table S3); (2) a set of Gene-  
255 Protein-Reaction (GPR) rules for all reactions considered (Additional file 7: Table S4); and (3) a  
256 user-defined tab-delimited file where the sets of reactions for complete pathways, subsets of  
257 reactions required to be performed by single species and transporter-associated genes (Additional  
258 file 8: Table S5) are described. The script *combinations.sh* included in OrtSuite  
259 (<https://github.com/mdsufz/OrtSuite/blob/master/combinations.sh>) is used to calculate all  
260 combinations of species capable of performing a set of reactions of interest to the user and results  
261 are stored in a text file.

262 To further reduce the vast amount of putative microbial interactions and to increase  
263 confidence in the results manual filtering can be performed to reflect available knowledge (e.g.

264 known cross-feeding relationship between species) or the likelihood of biologically feasible species  
265 interactions. The user also may have interest in assessing subsets of microbial interactions using  
266 specific criteria. Therefore, additional constraints can be applied to the list of putative microbial  
267 interactions to further reduce the search space. These include the degree of completeness of a  
268 pathway, the number of reactions expected to be performed by a single species, the presence or  
269 absence of transporter genes and the different number of species required to perform a complete  
270 pathway (e.g. only combinations of two or three species). The script *user\_constraints.sh* included in  
271 OrtSuite ([https://github.com/mdsufz/OrtSuite/blob/master/user\\_constraints.sh](https://github.com/mdsufz/OrtSuite/blob/master/user_constraints.sh)) allows the inclusion  
272 of these constraints.

273

#### 274 *Evaluation of ortholog clustering*

275 The clustering of orthologs was evaluated by measuring the pairwise precision and recall.  
276 Clustering precision measures how many pairs of sequences associated with the same molecular  
277 function are grouped together and is calculated by dividing the number of correctly clustered  
278 sequences by the total number of clustered sequences (Equation 2).

279

#### 280 *Equation 2*

$$281 \quad \textit{clustering precision} = \frac{\textit{correctly clustered sequences}}{\textit{total number of clustered sequences}}$$

282

283 where, *correctly clustered sequences* refers to the pairs of sequences that share the same function  
284 and are clustered together and *total number of clustered sequences* refers to all pairs of sequences  
285 that are clustered together irrespective of sharing the same function.

286 Clustering recall measures how many pairs of sequences with the same molecular function are not  
287 clustered together. Recall is calculated by dividing the number of correctly clustered sequences by  
288 the total true sequence clusters (Equation 3).

289

290 *Equation 3*

$$291 \quad \text{clustering recall} = \frac{\text{correctly clustered sequences}}{\text{total true sequence clusters}}$$

292

293 where, *correctly clustered sequences* refers to the pairs of sequences that share the same function  
294 and are clustered together and *total true sequence clusters* refers to all pairs of sequences that have  
295 the same function.

296

297 *Evaluation of sequence aligner used for clustering of orthologs*

298 Changes of a single DNA base can result in the production of a different amino acid which might  
299 result in a different protein. To determine the impact of mutations on the clustering of orthologs a  
300 single gene from three species was artificially mutated at different rates. These mutations were  
301 introduced in the nucleotide sequences of each gene. Only substitutions were considered since these  
302 are the most commonly studied [27] and none of the mutations were allowed to occur on the first  
303 and last codon. When, during the mutation, new stop or/and start codons were introduced, the  
304 translation was made for all the possible proteins and the largest was selected.

305 *Burkholderia vietnamiensis* G4 was mutated on the gene K05783, *Azoarcus sp.* CIB on the  
306 gene K07537 and *Aromatoleum aromaticum* EbN1 on the gene K07538. Each gene was mutated at  
307 rates 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  
308 original genome (e.g., *Burkholderia vietnamiensis* G4 strain K05783\_25, where “K05783” is the  
309 KEGG ortholog identifier and “25” is the rate of mutation). A total of 18 strains were generated (six  
310 *in silico* mutated strains per genome). The complete set of original and artificially mutated genomes

311 is available in a compressed file (Additional file 9: Set\_A\_genomes.zip) and at  
312 [https://github.com/mdsufz/OrtSuite/upload/master/Set\\_A\\_genomes.zip](https://github.com/mdsufz/OrtSuite/upload/master/Set_A_genomes.zip).

313

#### 314 *Evaluation of functional annotation*

315 Similar to clustering of orthologs, performance of functional annotation was determined based on  
316 the correct assignment of function to sequences in the ortholog clusters. *Annotation precision* is  
317 calculated by dividing the number of correctly annotated sequences by the total number of  
318 annotated sequences (Equation 4).

319

#### 320 *Equation 4*

$$321 \quad \textit{Annotation precision} = \frac{\textit{correctly annotated sequences}}{\textit{total number of annotated sequences}}$$

322

323 where, *correctly annotated sequences* refers to the number of sequences annotated with the correct  
324 function and *total number of annotated sequences* refers to the total number of sequences with  
325 annotation.

326 Annotation recall is calculated by dividing the number of correctly annotated sequences by  
327 total true annotated sequences (Equation 5).

328

#### 329 *Equation 5*

$$330 \quad \textit{Annotation recall} = \frac{\textit{correctly annotated sequences}}{\textit{total true sequence annotations}}$$

331

332 where, *correctly annotated sequences* refers to the number of sequences annotated with the correct  
333 function and *total true sequence annotations* refers to the total number of sequences annotated with  
334 the correct function and sequences which were wrongfully not assigned a function.

335

### 336 *Evaluation of microbial interaction predictions*

337 We evaluated the prediction of microbial interactions by selecting species with the potential to  
338 interact in the transformation of benzoate to acetyl-CoA identified by Fetzer and collaborators [23]  
339 (hereafter defined as Test\_genome\_set). These isolates did not have their genome sequenced prior  
340 this study. Therefore, bacterial cryo-cultures of the different isolates were revived on LB agar  
341 plates. Single colonies were picked and grown overnight in 2 ml LB medium at 37°C. The cells  
342 were pelleted by centrifugation. Cells were lysed and genomic DNA was extracted using a  
343 Nucleospin Tissue Kit (Machery and Nagel). Approximately 150 to 1000 ng of DNA were used for  
344 fragmentation (insert size: 300 – 700 bp) and sequencing libraries were prepared following the NEB  
345 Ultra II FS Kit protocol (New England Biolabs). Libraries were quantified using a JetSeq Library  
346 Quantification Lo-ROX Kit (Bioline) and quality-checked by Bioanalyzer (Agilent). These libraries  
347 were sequenced on an Illumina MiSeq Instrument with a final concentration of 8 pM using the v3  
348 600 cycles chemistry and 5% PhiX. The complete set of species and their genome sequences are  
349 available as compressed file (Additional file 10: Test\_genome\_set). According to Fetzer and  
350 collaborators [23], five of the twelve species (*C. testosteroni* ATCC 11996, *C. necator* JMP 134, *P.*  
351 *putida* ATCC 17514, *P. fluorescens* DSM 6290 and *R. ruber* BU3) were shown to grow as  
352 monocultures using benzoate as sole carbon and energy source (Additional file 11: Table S6). After  
353 their genomes were sequenced and assembled, we evaluated whether these species possessed the  
354 complete genomic content to encode all proteins required for each benzoate to acetyl-CoA  
355 conversion pathway. The remaining seven species were not able to grow as monocultures in media  
356 with benzoate as sole carbon source. Here, we evaluated whether the lack of growth was confirmed  
357 by lack of essential protein-encoding genes involved in conversion of benzoate to acetyl-CoA.

358 Additionally, Fetzter and collaborators [23] also observed that some of the species not able to grow  
359 as monocultures with benzoate as sole carbon and energy source could grow in medium with  
360 benzoate when being in a community with other species. Finally, their study also showed that,  
361 under nutrient specific conditions, total growth yields improved when communities contained  
362 species able to metabolize benzoate and several non-degrading species. Therefore, we evaluated  
363 whether putative species interactions identified by OrtSuite fit the results obtained by Fetzter et al.  
364 For the complete set of species combinations and benzoate degradation capabilities and effects  
365 identified by Fetzter and collaborators, see [23] (Additional file 12: Table S7).

366

## 367 **Results and discussion**

### 368 *Ortsuite is a flexible and user friendly pipeline*

369 The OrtSuite installation produces an environment on the user's computer with recently developed  
370 tools used for ortholog clustering [20] and sequence alignment [14]. Furthermore, OrtSuite contains  
371 several scripts that were developed to download sequence information, in bulk, and generate Gene-  
372 Protein-Reaction (GPR) rules from KEGG repositories in an automated manner by only requiring  
373 users to provide lists with identifiers (see <https://github.com/mdsufz/OrtSuite/>). Additional control  
374 is given to the user such as establishing thresholds in the minimum sequence identity (during  
375 sequence alignment of sequences in ortholog clusters to ORAdb) to restrict the number of putative  
376 microbial interactions based on constraints such as the presence of transporters and subsets of  
377 reactions to be performed by individual species (Additional file 8: Table S5). Since data in public  
378 repositories are frequently being added or updated and to include personal knowledge the user can  
379 manually curate the files in the ORAdb and GPR rules.

380 The entire workflow (excluding installation) consists of a total of 10 bash and 1 R easy to  
381 follow commands from the construction of the Ortholog-Reaction Association database (ORAdb),  
382 which will be used for the functional annotation of sequences and identification of potential  
383 microbial interactions.

384

385 *Computing time of OrtSuite stages*

386 The runtime of each OrtSuite step was evaluated on a set of genomes whose genomic content was  
387 known (Additional file 2: Table S2). The same set of genomes is used in the OrtSuite's GitHub  
388 tutorial page ([https://github.com/mdsufz/OrtSuite/blob/master/OrtSuite\\_tutorial.md](https://github.com/mdsufz/OrtSuite/blob/master/OrtSuite_tutorial.md)) with a total of  
389 75.5 Megabytes of data. OrtSuite was used to analyze this data on a laptop with 4 cores and 16  
390 Gigabytes of RAM. All OrtSuite steps were run on default settings, and the total runtime of each  
391 step was recorded (Additional file 13: Table S8). The total workflow was completed in 3 h 50 min  
392 and the longest single step runtime consisted of 2 h and 47 min which involved the construction of  
393 the ORAdb.

394 In this study, we ran OrtSuite on a dataset comprising 18 genomes and we postulate that up  
395 to 100 genomes can be analyzed in a customary laptop. To determine if this range would be within  
396 the number of genomes in regular microbiome studies we calculated the average number of  
397 metagenome-assembled genomes (MAGs) from metastudies where MAGs were recovered. A study  
398 performed by Parks and collaborators [8] analyzed sequencing data from 149 projects. Most  
399 projects (91%) consisted of less than 20 samples. On average, they recovered 5.3 metagenome-  
400 assembled genomes (MAGs) per metagenome. Work performed by Pasolli and collaborators [9] on  
401 microbial diversity in the human microbiome recovered, on average, 16 MAGs per metagenomic  
402 library. From the 46 studies used in their work, 30 consisted of less than 200 samples. Another  
403 study by Tully and collaborators focusing on marine environments [10] recovered 2631 MAGs  
404 from 234 samples (average of 11 MAGs per sample). Our analysis demonstrates that the average  
405 number of MAGs recovered from a metagenome currently range from five to 16. Therefore,  
406 performing targeted functional annotation and interspecies interactions predictions using OrtSuite in  
407 average sized metagenome samples is still feasible using a customary laptop.

408

409

410 *Higher recall rates during clustering of orthologs with DIAMOND*

411 OrthoFinder's [20] clustering of orthologs allows users to choose between DIAMOND [14] and

412 BLAST [13] as sequence aligners. To test which sequence aligner yielded the best results we

413 performed ortholog clustering of a dataset consisting of the original target genomes as well as a set

414 of artificially mutated genomes (Additional file 9: Set\_A\_genomes) using both aligners. The results

415 obtained (Table 1) do not show a striking difference between both aligners. However, an increase of

416 eight percent in the recall rate (which measures how many pairs of sequences with the same

417 molecular function are not clustered together) was observed when using DIAMOND.

418 *Table 1 – Pairwise precision and recall rates of clustering of orthologs using OrthoFinder*

	BLAST	DIAMOND
Precision	0.63	0.64
Recall	0.77	0.85

419 *Pair-wise precision and recall results of OrthoFinder using BLAST and DIAMOND as an alignment search tool. Clustering was performed on the*

420 *dataset Set\_A\_genomes plus the mutated genomes.*

421

422 All artificially mutated sequences (even those with mutation rates of 25%) were clustered

423 together with their non-mutated ortholog. In parallel, we also performed sequence alignment using

424 NCBI's BLASTp [28] between the protein sequences of the DNA-mutated and un-mutated genes

425 (Additional file 14: Table S9). E-values of sequence alignments in all species ranged from 0 to 5E-

426 180 and percentage of identity from 61.32 to 98.84% which indicates that the mutated genes still

427 shared enough sequence similarity to the original protein sequence. The results obtained during the

428 clustering of orthologs combined with the protein sequence alignments suggest that mutation rates

429 of up to 25% of single DNA base pairs will not have an observable effect on the clustering of

430 orthologs. OrthoFinder's algorithm removes the gene length bias from the sequence alignment

431 process, which may also explain why mutated genes were clustered with the original. Although it

432 has been suggested that most genetic variations are neutral, changes in single base pairs can have a

433 drastic effect on protein function (e.g. depending on the location where the mutation occurs) [29].  
434 To this purpose, experimental functional studies can be used to validate previously unannotated  
435 orthologs. Nevertheless, future studies increasing the rates of DNA base pair substitutions and other  
436 types of mutations as well as experiments targeting protein function in ortholog clusters are  
437 warranted.

438 For the subsequent steps of the OrtSuite workflow, clustering of protein orthologs was  
439 repeated using only the original genomes and the default aligner. A complete overview of the  
440 results generated during the clustering of orthologs (e.g. number of genes in ortholog clusters,  
441 number of unassigned genes and number of ortholog clusters) was also obtained (Additional file 15:  
442 Table S10 and at [https://github.com/mdsufz/OrtSuite/blob/master/examples/  
443 /Comparative\\_Genomics\\_Statistics/Statistics\\_Overall\\_test\\_set.tsv](https://github.com/mdsufz/OrtSuite/blob/master/examples/Comparative_Genomics_Statistics/Statistics_Overall_test_set.tsv).

444

#### 445 *Precision and recall rates of annotation improve with relaxed identity thresholds*

446 The third step of OrtSuite consists of performing cluster annotation in a two-stage process. In the  
447 first, only 10% of sequences are used in the alignment to the sequences ORAdb. Those with a  
448 minimum identity threshold proceed to the second stage where all sequences contained in the  
449 proceeded cluster are used. At the end, annotation of clusters will take into consideration additional  
450 parameters such as percentage of positive matches. To evaluate the thresholds used in each stage of  
451 the annotation of ortholog clusters we used two different relaxed (40 and 70%) and restrictive (90  
452 and 95%) search identity cutoffs. An overview of the results (e.g. number of clusters containing  
453 orthologs from ORAdb, number of ortholog clusters with annotated sequences) is also generated  
454 (Additional file 16: Table S11). Precision rates for functional annotation were not strikingly  
455 different when switching from 40 to 70% relaxed identity cutoffs and 90 to 95% restrictive search  
456 cutoffs (Table 2). However, recall rates were, on average 15 to 25 % higher when using a 90%  
457 identity cutoff during the restrictive search compared to 95%. Furthermore, the difference in  
458 computing time between lower and higher identity thresholds was negligible (< 2 min).

459 Other annotation tools, such as NCBI's BLAST tool [13] and Prokka [30], can annotate full  
 460 genomes, the latter at a relatively fast pace. On average, full genome annotation of our genomes in  
 461 the Set\_A\_genomes dataset using Prokka required 12 mins on a customary laptop with 16  
 462 Gigabytes of RAM and four CPUs to complete. However, exploring the greater amount of data  
 463 generated from full genome annotation of individual species from complex microbial communities  
 464 is a daunting task. Studies aiming to address specific ecosystem processes or metabolic functions  
 465 will require a manual inspection of each annotated genome to confirm the presence of the  
 466 associated genes. The same drawbacks are observed when attempting to identify interspecies  
 467 interactions based on the microbe's genomic potential. OrtSuite overcomes these challenges by  
 468 limiting the genome annotation to user-defined functions of interest. The user-defined database  
 469 coupled with the scripts for automated identification of interspecies interactions contained in  
 470 OrtSuite decreases the time required to generate the data and facilitates its interpretation by the user.  
 471

472 *Table 2 - Precision and recall rates of annotation for the Set\_A\_genomes dataset.*

Pathway	BPA_P1				BTA_P2				BTA_P3			
	40		70		40		70		40		70	
Relaxed search identity cutoff (%)	40		70		40		70		40		70	
Restricted search identity cutoff (%)	90	95	90	95	90	95	90	95	90	95	90	95
Precision	0.89	0.92	0.89	0.90	0.90	0.93	0.89	0.92	1	1	1	1
Recall	0.66	0.42	0.57	0.33	0.67	0.44	0.59	0.37	0.89	0.75	0.91	0.59

473 *Precision and recall rates of annotation using different combinations of identity and score cut-offs. The relaxed search identity cut-offs were set to 40*  
 474 *and 70%. The restrictive search score cut-offs were set to 90 and 95%. Annotation was performed on the genomes from the Set\_A\_genomes dataset*  
 475 *for all three benzoate to acetyl-CoA conversion pathways (P1\_BTA, P2\_BTA and P3\_BTA) used in this work.*  
 476

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

478 Ortsuite's ortholog cluster annotation approach (i.e. based on user-defined functions or pathways)  
 479 not only decreases the time needed to obtain the functional potential of individual species but also  
 480 facilitates analysis of the results for the whole community. To test Ortsuite's ability to identify

481 species with the genetic potential to perform a pathway individually we defined sets of reactions  
482 that are used in three alternative pathways for the conversion of benzoate to acetyl-CoA (Additional  
483 file 8: Table S5). Next we compared the results to the species' known genomic content in each  
484 alternative pathway (shown in Additional file 2: Table S2). This reported species' individual  
485 genomic potential to convert benzoate to acetyl-CoA in every alternative pathway was confirmed  
486 by OrtSuite. Using OrtSuite we identified seven species capable of performing one of the anaerobic  
487 conversion pathways (BTA\_P2). A total of eight species were identified with the ability to perform  
488 the aerobic conversion of benzoate to acetyl-CoA (BTA\_P3). No single species was identified by  
489 OrtSuite with the ability to convert benzoate to acetyl-CoA via the alternative anaerobic process  
490 (BTA\_P1). For these 12 species, cluster annotation using ORAdb required only 6 min using a  
491 laptop with 16 Gigabytes of RAM and 4 CPUs (Additional file 13: Table S8). In this study, we only  
492 defined three pathways for analysis but the user can increase this number without a substantial  
493 increase in running times. Additionally, the results are tailored to the user-defined constraints (see  
494 example, Additional file 8: Table S5) facilitating analysis.

495

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

497 We designed OrtSuite to allow hypothesis-driven exploration of microbial interactions by limiting  
498 the search by user-defined ecosystem processes or metabolic functions. In this study, we tested  
499 these OrtSuite functions by identifying interspecies interactions involved in the aerobic and  
500 anaerobic conversion of benzoate to acetyl-CoA where experimental data were available. OrtSuite's  
501 ability to predict interspecies interactions was assessed on a set of sequenced isolates  
502 (Test\_genome\_set) whose growth as mono- and mixed-cultures was analyzed previously [23] under  
503 three different environmental conditions (low substrate: 1g/L benzoate, high substrate: 6g/L  
504 benzoate and high substrate+salt stress: 6g/L benzoate supplemented with 15g /L of NaCl). This  
505 dataset contained 69,193 protein sequences distributed across the 12 species resulting in a total of  
506 59 Megabytes of data. More than 84% of all genes were placed in 9,533 ortholog clusters. In

507 addition, 541 clusters were composed of sequences obtained from all 12 species (Additional file 15:  
508 Table S10.tsv). OrtSuite's annotation stage resulted in 44 ortholog clusters with annotated  
509 sequences from ORAdb (Additional file 16: Table S11.csv). The mapping of KOs to each species in  
510 the Test\_genome\_set is available as supplementary data (Additional file 17: Table S12). The  
511 genomic potential of each species for aerobic and anaerobic benzoate metabolizing pathways is  
512 shown in Figure 3. The complete mapping of reactions to each species is available in the  
513 supplementary data (Additional file 18: Table S13). Based on the 44 ortholog clusters and the Gene-  
514 Protein-Reaction (GPR) rules, four species (*Cupriavidus necator* JMP134, *Pseudomonas putida*  
515 ATCC17514, *Rhodococcus ruber* BU3 and *Sphingobium yanoikuyae* DSM6900) contained all  
516 protein-encoding genes required to perform aerobic conversion of benzoate to acetyl-CoA. Except  
517 for *S. yanoikuyae* all others were in agreement with Fetzer's observation of their ability to grow in  
518 medium containing 1g/L of benzoate as the carbon source. In Fetzer's study, in a medium  
519 containing 1g/L of benzoate, growth was considered when optical densities (OD) were above 0.094.  
520 The OD measured for *S. yanoikuyae* was 0.0916. The annotation of genes with the ability to  
521 perform the complete aerobic conversion of benzoate to acetyl-CoA combined with the small  
522 difference in OD to the minimum threshold suggests that *S. yanoikuyae* indeed can grow on low  
523 benzoate containing medium but at lower growth rates. The observed growth of *Comamonas*  
524 *testosteroni* ATCC11996 and *Pseudomonas fluorescens* DSM6290 in the low benzoate  
525 environment was not correlated with their genomic potential according to OrtSuite in any benzoate-  
526 to-acetyl-CoA conversion pathway. It is likely that these two species utilize reactions or pathways  
527 that were not included in the benzoate degradation pathways used in our study.

528         Based on the genomic content, any combination of species where one of them possesses the  
529 ability to perform a complete pathway will always result in a potential interspecies interaction. To  
530 account for this scenario all species with the complete genomic potential to perform a complete  
531 pathway are excluded when calculating interspecies interactions. In the aerobic benzoate  
532 degradation pathway (P3\_BTA) the protein-encoding genes that are involved in the conversion of

533 2-Hydroxy-2,4-pentadienoate to 2-Hydroxymuconate semialdehyde (R02604) were lacking in all  
534 remaining species and therefore, no species interactions were obtained. Similarly, in the anaerobic  
535 degradation pathways (P1\_BTA and P2\_BTA) no species presented the genomic content to encode  
536 proteins involved in the reactions R02451, R03028, R05579, R05581, R05594 and R05597.  
537 Therefore, no species interactions were identified that would allow the complete anaerobic  
538 conversion of benzoate to acetyl-CoA. Next, we considered the interactions where species with  
539 complete conversion capabilities are present. In the low substrate environment, OrtSuite identified  
540 776 of 830 (93.4%) species combinations showing growth. In the high substrate environment,  
541 OrtSuite predicted 614 of 646 (95%). In the high substrate+salt stress environment, OrtSuite  
542 predicted all 271 (100%) combinations of species exhibiting growth (Additional file 19: Table S14).  
543 According to the same study in the high benzoate concentrations, 237 combinations of species did  
544 not exhibit growth ( $OD < 0.2545$ ) even in the presence of known benzoate degraders which  
545 occurred in 160 combinations. From this group, OrtSuite identified 185 interspecies interactions  
546 that had the potential to perform aerobic conversion of benzoate. A possible explanation as to why  
547 growth in Fetzer's study was not observed, despite the presence of benzoate degraders, is the lack of  
548 tolerance of these species to high benzoate concentrations. For example, *C. necator* growth was  
549 shown to be stimulated at low benzoic acid concentrations but inhibited at high concentrations [31].  
550 To assess these effects, additional stress-resistance associated genes and reactions such as those  
551 involved in medium acidification [32] could be added as constraints. Similar results were obtained  
552 when using a high substrate+salt stress medium. Growth of species combinations was not observed  
553 despite the presence of benzoate degraders. Benzoate degradation has been shown to decrease in  
554 hyperosmotic environments [33] therefore, additional constraints such as genes that confer  
555 resistance to environmental stressors or adverse conditions could be included during the  
556 identification of interspecies interactions under different or changing environmental conditions.

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

559 growth experiments were conducted in aerobic conditions, it is possible that the species in question  
560 are only capable of using benzoate as a carbon source in aerobic environments. To fully explore all  
561 the species potential to convert benzoate, additional degradation pathways could be tested in the  
562 future in a multi-omics context. Furthermore, the only constraints added were related to the  
563 reactions that composed each pathway. Additional constraints can be included in future studies,  
564 such as potential mandatory transport-associated reactions, to increase confidence in the proposed  
565 interspecies interactions. In summary, OrtSuite confirmed that most interspecies interactions  
566 (96.1%) identified by Fetzer and collaborators were possible via the aerobic conversion of benzoate  
567 to acetyl-CoA but were not possible at anoxic conditions.

568

## 569 **Conclusions**

570 The OrtSuite workflow allows hypothesis-driven exploration of potential interactions between  
571 microbial genomes by limiting the search universe to a user-defined set of ecosystem processes.  
572 This is achieved by rapidly assessing the genetic potential of a microbial community for a given set  
573 of reactions. The two-step annotation of clusters of orthologs with a personalized ORAdb decreases  
574 the overall number of sequence alignments that need to be computed. User-specified constraints,  
575 such as the presence of transporter genes, further reduces the search space for putative microbial  
576 interactions. Users have substantial control over several steps of OrtSuite: from manual curation of  
577 ORAdb, custom sequence similarity cutoffs to the addition of constraints for inference of putative  
578 microbial interactions. As long as links between genes and reactions exist, the flexibility and easy  
579 usage of OrtSuite allow its application to the study of any given ecosystem process.

580

## 581 **Declaration Sections**

582

## 583 **Ethics approval and consent to participate**

584 Not applicable

585 **Consent for publication**

586 Not applicable

587

588 **Availability of data and materials**

589 The genomes used to test the workflow are available at National Centre for Biotechnology  
590 Information (<https://www.ncbi.nlm.nih.gov/>) under the accession identifiers CP029389-CP029397,  
591 GCF\_000001735, AP012304, AP012305, CP021731, CP011072, CP007785-CP007787,  
592 CP000614-CP000621, CP003230, CP005996, CP003108, CR555306-CR5553068,  
593 GCF\_000225785, LN997848-LN997849, CP022989-CP022996, CP024315, AP012547,  
594 CP022046-CP022047 and CP001281-CP001282. The genome assemblies used to predict  
595 interspecies interactions are available at the European Nucleotide Archive (ENA -  
596 <https://www.ebi.ac.uk/ena>) with the study accession PRJEB38476. All analysis results and scripts  
597 used to generate them are available at <https://github.com/mdsufz/OrtSuite>.

598

599 **Competing interests**

600 The authors declare that they have no competing interests.

601

602 **Funding**

603 This work was funded by the Helmholtz Young Investigator grant VH-NG-1248 Micro ‘Big Data’.

604

605 **Author’s contributions**

606 JS, OD, PS and UNR developed the concept of OrtSuite. JS, MG and UNR developed the OrtSuite  
607 workflow. JS, MG and UNR performed the benchmarks. CV provided information and data for  
608 defining benzoate to acetyl-CoA conversion pathways. RK sequenced bacterial isolates provided by  
609 AC. JS and UNR wrote the manuscript. All authors read and approved the manuscript.

610

611 **Acknowledgments**

612 We thank the early users of OrtSuite Sandra Silva and Felipe Côrrea for their help with beta testing  
613 and for workflow suggestions. We also thank Diogo Lima for his assistance in the implementation  
614 of the script required to generate the Gene-Protein-Reaction (GPR) rules; and Nicole Steinbach for  
615 her work in the sequencing of the isolates used as the test set.

616

617 **Availability and requirements**

618 Project name: OrtSuite

619 Project home page: <https://github.com/mdsufz/OrtSuite>

620 Operating system: Linux64

621 Programming languages: Shell; Python3.6, R

622 Other requirements: Java, BLAST+ [34], BLAST [13], DIAMOND [14], OrthoFinder [20], mcl

623 [35]

624 License: MIT

625

626 **List of abbreviations**

627 BTA: Benzoate-to-Acetyl-CoA

628 CPU: Central Processing Unit

629 EC: Enzyme Commission

630 GPR: Gene-Protein-Reaction

631 KEGG: Kyoto Encyclopedia of Genes and Genomes

632 KO: KEGG ortholog

633 OD: Optic Density

634 ORAdb: Ortholog-Reaction Association database

635 ORF: Open Reading Frame

636 RAM: Random Access Memory

637 **References**

- 638 1. Maestre FT, Castillo-Monroy AP, Bowker MA, Ochoa-Hueso R. Species richness effects on  
639 ecosystem multifunctionality depend on evenness, composition and spatial pattern: Community  
640 attributes and multifunctionality. *J Ecol.* 2012;100:317–30.
- 641 2. Mulder CPH, Uliassi DD, Doak DF. Physical stress and diversity-productivity relationships: The  
642 role of positive interactions. *Proc Natl Acad Sci.* 2001;98:6704–8.
- 643 3. Lyu Z, Shao N, Akinyemi T, Whitman WB. Methanogenesis. *Curr Biol.* 2018;28:R727–32.
- 644 4. Liu F, Hewezi T, Lebeis SL, Pantalone V, Grewal PS, Staton ME. Soil indigenous microbiome  
645 and plant genotypes cooperatively modify soybean rhizosphere microbiome assembly. *BMC*  
646 *Microbiol.* 2019;19:201.
- 647 5. Locey KJ, Lennon JT. Scaling laws predict global microbial diversity. *Proc Natl Acad Sci.*  
648 2016;201521291.
- 649 6. Raynaud X, Nunan N. Spatial Ecology of Bacteria at the Microscale in Soil. *PLoS ONE.* 2014
- 650 7. Roh SW, Abell GCJ, Kim K-H, Nam Y-D, Bae J-W. Comparing microarrays and next-generation  
651 sequencing technologies for microbial ecology research. *Trends Biotechnol.* 2010;28:291–9.
- 652 8. Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, et al. Recovery of  
653 nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nat Microbiol.*  
654 2017;1.
- 655 9. Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, et al. Extensive Unexplored  
656 Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning  
657 Age, Geography, and Lifestyle. *Cell.* 2019;176:649-662.e20.
- 658 10. Tully BJ, Graham ED, Heidelberg JF. The reconstruction of 2,631 draft metagenome-assembled  
659 genomes from the global oceans. *Sci Data.* 2018;5:170203.
- 660 11. Morin M, Pierce EC, Dutton RJ. Changes in the genetic requirements for microbial interactions  
661 with increasing community complexity. *eLife.* 2018;7:e37072.
- 662 12. Slade EM, Kirwan L, Bell T, Philipson CD, Lewis OT, Roslin T. The importance of species  
663 identity and interactions for multifunctionality depends on how ecosystem functions are valued.  
664 *Ecology.* 2017;98:2626–39.
- 665 13. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J*  
666 *Mol Biol.* 1990;215:403–10.
- 667 14. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat*  
668 *Methods.* 2015;12:59–60.
- 669 15. Koonin EV. Orthologs, Paralogs, and Evolutionary Genomics. *Annu Rev Genet.* 2005;39:309–  
670 38.
- 671 16. Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, von Mering C, et al. Fast  
672 Genome-Wide Functional Annotation through Orthology Assignment by eggNOG-Mapper. *Mol*  
673 *Biol Evol.* 2017;34:2115–22.

- 674 17. Li L, Stoeckert CJ, Roos DS. OrthoMCL: Identification of Ortholog Groups for Eukaryotic  
675 Genomes. *Genome Res.* 2003;13:2178–89.
- 676 18. Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein or  
677 nucleotide sequences. *Bioinforma Oxf Engl.* 2006;22:1658–9.
- 678 19. Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome comparisons  
679 dramatically improves orthogroup inference accuracy. *Genome Biol.* 2015;16:157.
- 680 20. Emms DM, Kelly S. OrthoFinder: phylogenetic orthology inference for comparative genomics.  
681 *Genome Biol.* 2019;20:238.
- 682 21. Gottstein W, Olivier BG, Bruggeman FJ, Teusink B. Constraint-based stoichiometric modelling  
683 from single organisms to microbial communities. *J R Soc Interface.* 2016;13:20160627.
- 684 22. Heirendt L, Arreckx S, Pfau T, Mendoza SN, Richelle A, Heinken A, et al. Creation and  
685 analysis of biochemical constraint-based models using the COBRA Toolbox v.3.0. *Nat Protoc.*  
686 2019;14:639–702.
- 687 23. Fetzer I, Johst K, Schäwe R, Banitz T, Harms H, Chatzinotas A. The extent of functional  
688 redundancy changes as species' roles shift in different environments. *Proc Natl Acad Sci U S A.*  
689 2015;112:14888–93.
- 690 24. Kanehisa M, Goto S, Kawashima S, Okuno Y, Hattori M. The KEGG resource for deciphering  
691 the genome. *Nucleic Acids Res.* 2004;32:D277–D280.
- 692 25. Gabaldón T, Koonin EV. Functional and evolutionary implications of gene orthology. *Nat Rev*  
693 *Genet.* 2013;14:360–6.
- 694 26. Pertsemlidis A, Fondon JW. Having a BLAST with bioinformatics (and avoiding  
695 BLASTphemy). *Genome Biol.* 2001;2:reviews2002.1.
- 696 27. Lynch M. Evolution of the mutation rate. *Trends Genet.* 2010;26:345–52.
- 697 28. Madden T. The BLAST Sequence Analysis Tool. National Center for Biotechnology  
698 Information (US); 2003.
- 699 29. Ng PC, Henikoff S. Predicting the Effects of Amino Acid Substitutions on Protein Function.  
700 *Annu Rev Genomics Hum Genet.* 2006;7:61–80.
- 701 30. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinforma Oxf Engl.*  
702 2014;30:2068–9.
- 703 31. Wang W, Yang S, Hunsinger GB, Pienkos PT, Johnson DK. Connecting lignin-degradation  
704 pathway with pre-treatment inhibitor sensitivity of *Cupriavidus necator*. *Front Microbiol.* 2014
- 705 32. Kitko RD, Cleeton RL, Armentrout EI, Lee GE, Noguchi K, Berkmen MB, et al. Cytoplasmic  
706 Acidification and the Benzoate Transcriptome in *Bacillus subtilis*. *PLoS ONE.* 2009.
- 707 33. Bazire A, Diab F, Jebbar M, Haras D. Influence of high salinity on biofilm formation and  
708 benzoate assimilation by *Pseudomonas aeruginosa*. *J Ind Microbiol Biotechnol.* 2007;34:5–8.
- 709 34. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:  
710 architecture and applications. *BMC Bioinformatics.* 2009;10:421.

711 35. van Dongen S. Graph Clustering by Flow Simulation. University of Utrecht; 2000.

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737 **Figures legends**

738

739 Figure 1. OrtSuite workflow. Protein sequences from samples supplied by the user are clustered  
740 using OrthoFinder **(a)**. OrtScraper takes a text file containing a list of identifiers for each reaction in  
741 the pathway of interest supplied by the user to retrieve all protein sequences from KEGG **(b)**.  
742 Sequences mapped to reactions are stored in ORAdb **(c)**. Functional annotation **(d)** consists of a  
743 two-stage process (relaxed and restrictive search). Relaxed search **(e)** performs sequence alignments  
744 between 10% of randomly selected sequences from each generated cluster. Clusters whose  
745 representative sequences share a minimum 50% identity to sequences in reaction set(s) in ORAdb  
746 transition to the restrictive search **(f)**. Here, all sequences from the cluster are aligned to all  
747 sequences in the corresponding reaction set(s) to which they had a hit. Finally, the annotated  
748 sequences are used to identify putative microbial interactions based on their functional potential **(g**  
749 **and h)**. Additional constraints can be added to reduce the search space of microbial interactions  
750 (rounded blue rectangles) (e.g. subsets of reactions required to be performed by single species,  
751 transport-related reactions required).

752

753 Figure 2. Benzoate to acetyl-CoA conversion pathways. Reactions in blue background illustrate the  
754 anaerobic degradation of benzoate to acetyl-CoA via benzoyl-CoA. Reactions in yellow background  
755 illustrate the aerobic degradation of benzoate via catechol.

756

757 Figure 3. Mapping of the genomic potential of each species from the Test\_genome\_set dataset to  
758 each reaction in aerobic (yellow) and anaerobic (blue) benzoate-to-acetyl-CoA conversion  
759 pathways.

760

761

762

763 **Additional files**

764

765 Additional file 1:

766 **Table S1.** Number of reactions, enzymes, KO groups and KO-associated sequences represented in  
767 each alternative benzoate to acetyl-CoA conversion pathway used (BTA\_P1, BTA\_P2 and  
768 BTA\_P3).

769

770 Additional file 2:

771 **Table S2.** Species names, strain and abbreviation codes used to validate OrtSuite. The genomic  
772 potential, based on KEGG database, to completely encode all proteins involved in a BTA pathway  
773 is identified in the column “BTA pathway” (P1 – Anaerobic conversion of benzoate to acetyl-CoA  
774 1; P2 – Anaerobic conversion of benzoate to acetyl-CoA 2; P3 – Aerobic conversion of benzoate to  
775 acetyl-CoA) . The column OrtSuite\_result contains which BTA pathway(s) were identified as being  
776 completely performed by each species used for validation.

777

778 Additional file 3:

779 **Example\_reaction\_list.** Example file of a list of KEGG reaction identifiers used to generate  
780 ORAdb.

781

782 Additional file 4:

783 **Example\_ec\_list.** Example file of list of enzyme commission (EC) numbers used to generate  
784 ORAdb.

785

786 Additional file 5:

787 **Example\_ko\_list.** Example file of list of KEGG ortholog identifiers used to generate ORAdb.

788

789 Additional file 6:

790 **Table S3.** Example of the binary table with the mapping of KO identifiers (rows) to the species of  
791 interest (columns) (1 – indicates the presence of genes in the species; 0 – indicates the absence of  
792 genes in the species).

793

794 Additional file 7:

795 **Table S4.** Gene-Protein-Reaction (GPR) rules and metadata associated with all reactions present in  
796 the Ortholog Reaction-Associated database (ORAdb).

797

798 Additional file 8:

799 **Table S5.** User-defined constraints: pathway name, complete set of reactions present in each  
800 pathway, sets of reactions required to be performed by single species (each subset is described  
801 between parenthesis) and transport reactions. Transporter column describes the transport reaction  
802 (e.g. R00750) that is coupled to the reaction in the pathway (e.g. R02601). Thus, species that  
803 perform the latter must also contain the genes associated with the transport reaction.

804

805 Additional file 9:

806 **Set\_A\_genomes.** Compressed folder containing all the genomes (FASTA format) of the test species  
807 used to test the workflow. Species abbreviation codes are used as file names.

808

809 Additional file 10:

810 **Test\_genome\_set.** Sequenced genomes of the species used in a study by Fetzer and collaborators  
811 [23].

812

813 Additional file 11:

814 **Table S6.** Species used in the Fetzer study [23] and their ability to grow as monocultures in  
815 benzoate-containing medium.

816

817 Additional file 12:

818 **Table S7.** Growth results of microbial communities composed of individuals and groups of species  
819 from the Fetzer study [23] on three different environments. Environment 1 contained a benzoate  
820 concentration of 1g/L. Environment 2 contained a benzoate concentration of 6g/L. Environment 3  
821 contained a benzoate concentration of 6g/L and 15 g/L of NaCl.

822

823 Additional file 13:

824 **Table S8.** OrtSuite workflow runtime. The total runtime of each OrtSuite step when analyzing the  
825 genomic potential of species in Set\_A\_genomes dataset in three pathways (P1, P2 and P3) for the  
826 conversion of benzoate to acetyl-CoA (BTA). Steps were performed with default parameters on a  
827 laptop with 4 cores and 16 GB of RAM.

828

829 Additional file 14:

830 **Table S9.** Sequence alignments of original and mutated sequences using BLAST [13].

831

832 Additional file 15:

833 **Table S10.** Statistics obtained during the clustering of protein orthologs using the  
834 Test\_genome\_set. Results include numbers of species, genes, clusters and genes per cluster.

835

836 Additional file 16:

837 **Table S11.** Overview of the number of clusters, sequences and KOs during the annotation of the  
838 Test\_genome\_set.

839

840 Additional file 17:

841 **Table S12.** Mapping of species annotated with KO identifiers from ORAdb.

842

843 Additional file 18:

844 **Table S13.** Potential of species in the Test\_genome\_set to perform reactions associated with

845 benzoate degradation based on ORAdb. Species identifiers: (A) *Bacillus subtilis* ATCC, (B)

846 *Paenibacillus polymyxa* ATCC 842, (C) *Brevibacillus brevis* ATCC 8246, (D) *Comamonas*

847 *testosteroni* ATCC 11996, (E) *Cupriavidus necator* JMP 134, (F) *Pseudomonas putida* ATCC

848 17514, (G) *Pseudomonas fluorescens* DSM 6290, (H) *Variovorax paradoxus* ATCC 17713, (I)

849 *Rhodococcus sp.* (isolate UFZ), (J) *Acidovorax facilis* (isolate UFZ), (K) *Rhodococcus ruber* BU3,

850 (L) *Sphingobium yanoikuyae* DSM 6900. (1 – species with the complete genomic potential to

851 perform the reaction; 0 – species without the complete genomic potential to perform the reaction).

852

853 Additional file 19:

854 **Table S14.** Growth of single species and in combination with others measured by Fetzer and

855 collaborators in three different media (low substrate: 1g/L benzoate, high substrate: 6g/L benzoate

856 and high substrate+salt stress: 6g/L benzoate supplemented with 15 g/L of NaCl). Species

857 identifiers: (A) *Bacillus subtilis* ATCC, (B) *Paenibacillus polymyxa* ATCC 842, (C) *Brevibacillus*

858 *brevis* ATCC 8246, (D) *Comamonas testosterone* ATCC 11996, (E) *Cupriavidus necator* JMP 134,

859 (F) *Pseudomonas putida* ATCC 17514, (G) *Pseudomonas fluorescens* DSM 6290, (H) *Variovorax*

860 *paradoxus* ATCC 17713, (I) *Rhodococcus sp.* (isolate UFZ), (J) *Acidovorax facilis* (isolate UFZ),

861 (K) *Rhodococcus ruber* BU3, (L) *Sphingobium yanoikuyae* DSM 6900. Growth was considered

862 when optical density (OD) was equal or above 0.094, 0.2545 and 0.0752 in environments with low

863 substrate, high substrate and high substrate+salt stress benzoate, respectively.

864