

# 1 **IMP: a pipeline for reproducible integrated** 2 **metagenomic and metatranscriptomic analyses**

3

4 Shaman Narayanasamy<sup>†1</sup>, Yohan Jarosz<sup>†1</sup>, Emilie E.L. Muller<sup>1</sup>, Cédric C. Laczny<sup>1°</sup>, Malte  
5 Herold<sup>1</sup>, Anne Kaysen<sup>1</sup>, Anna Heintz-Buschart<sup>1</sup>, Nicolás Pinel<sup>2°</sup>, Patrick May<sup>1</sup>, and Paul  
6 Wilmes<sup>1\*</sup>

7

8 \* Correspondence:

9 [paul.wilmes@uni.lu](mailto:paul.wilmes@uni.lu)

10 <sup>†</sup> Equal contributors

11

## 12 **Author details**

13 <sup>1</sup> Luxembourg Centre for Systems Biomedicine, 6 avenue du Swing, University of  
14 Luxembourg, L-4367 Belvaux, Luxembourg

15 <sup>2</sup> Institute of Systems Biology, 401 Terry Avenue North, WA 98109, Seattle, USA.

16 <sup>°</sup> Current affiliation: CCL - Saarland University, Building E2 1, 66123 Saarbrücken, Germany;

17 NP - Universidad EAFIT, Carrera 49 No 7 sur 50, Medellín, Colombia.

18

19

## 20 **Abstract**

21 We present IMP, an automated pipeline for reproducible integrated analyses of coupled  
22 metagenomic and metatranscriptomic data. IMP incorporates preprocessing, iterative co-  
23 assembly of metagenomic and metatranscriptomic data, analyses of microbial community  
24 structure and function as well as genomic signature-based visualizations. Complementary use  
25 of metagenomic and metatranscriptomic data improves assembly quality and enables the  
26 estimation of both population abundance and community activity while allowing the recovery  
27 and analysis of potentially important components, such as RNA viruses. IMP is containerized  
28 using Docker which ensures reproducibility. IMP is available at <http://r3lab.uni.lu/web/imp/>.

29

30 **Keywords:** multi-omics data integration; metagenomics; metatranscriptomics; microbial  
31 ecology; microbiome; reproducibility

## 32 **Background and motivation**

33 Microbial communities are ubiquitous in nature and govern important processes related to  
34 human health and biotechnology [1, 2]. A significant fraction of naturally occurring  
35 microorganisms elude detection and investigation using classical microbiological methods due  
36 to their unculturability under standard laboratory conditions [3]. The issue of unculturability is  
37 largely circumvented through the direct application of high-resolution and high-throughput  
38 molecular measurements, most notably metagenomics, to microbial community samples  
39 collected *in situ* [4–6]. Beyond metagenomics, there is also a clear need to obtain functional  
40 readouts in the form of additional layers of omics data from consortia. Moreover, there is a  
41 growing desire to integrate the resulting meta-omics data to more conclusively link genetic  
42 potential to actual phenotype *in situ* [6]. For this purpose, specialized wet-lab methods have  
43 been developed to ensure that the generated data fulfill the premise of systematic measurements  
44 [7], as subsampling has been shown to inflate intra- and inter-sample variation, thereby  
45 hampering subsequent data integration, individual biomolecular fractions, i.e. DNA, RNA,  
46 proteins and metabolites are derived from single, unique samples [7, 8]. Next-generation  
47 sequencing (NGS) of microbial community derived DNA and reverse transcribed RNA  
48 (cDNA) results in metagenomic (MG) and metatranscriptomic (MT) data, respectively.  
49 Additional layers of meta-omic data include the metaproteome (MP) and the (meta-  
50 )metabolome [4–6, 9]. Although standardized and reproducible wet-lab methods have been  
51 developed for integrated omics of microbial communities, corresponding dry-lab workflows  
52 have yet to be formalized.

53

54 Computational solutions for the analysis of MG and MT data can be broadly  
55 categorized into reference-dependent or reference-independent (*de novo*) methods [5].  
56 Reference-dependent methods are based on the direct alignment/mapping of sequencing reads

57 onto isolate genomes, gene catalogs or MG data. A major drawback of such methods is the  
58 large number of sequencing reads discarded during data analysis, due to their dissimilarity from  
59 genomes/genes within the reference databases. More specifically, based on analyses of MG  
60 data from the human gut microbiome, which is arguably the most resolved microbial ecosystem  
61 in terms of representative isolate genomes, 43% of organisms are not represented by isolate  
62 genomes [10], while 74%-81% of sequencing data is typically represented within an integrated  
63 gene catalog based on MG data [11], exemplifying a substantial loss of potentially useful  
64 information when using isolate genome reference databases alone. Conversely, reference-  
65 independent methodologies, such as approaches based on *de novo* assembly, enable the  
66 retrieval of previously uncharacterized genomes and/or potentially novel genes, providing an  
67 added advantage over reference-based methods [4, 5, 12]. Furthermore, it has been shown that  
68 the assembly of sequencing reads into longer contiguous sequences (contigs) greatly improves  
69 the taxonomic assignments and prediction of genes as opposed to direct identification from  
70 short sequencing reads [13, 14].

71

72         Given the advantages of reference-independent methods, a wide array of MG-specific  
73 assemblers such as IDBA-UD [15] and MEGAHIT [16] have been developed. Most MT  
74 analyses consist of reference genome- or metagenome-dependent workflows, similar to the  
75 reference-dependent workflows used for MG data [17, 18]. However, reference-independent  
76 approaches for metatranscriptomic data exploitation are also available either using specialized  
77 metatranscriptome assemblers (e.g. IDBA-MT [19]), metagenome assemblers (e.g. IDBA-UD  
78 [15]) or single-species transcriptome assemblers (e.g. Trinity [20]) [14]. In both cases, the  
79 available assemblers are able to handle the uneven sequencing depths of MG and MT data.  
80 Although dedicated assembly methods have been developed for MG and MT data, formalized  
81 pipelines allowing both data types to be used in an integrated way have yet to be developed.

82

83 Automated data processing and analysis pipelines have so far been mainly developed  
84 for MG data. These include pipelines such as MOCAT [21] and MetAMOS [12] which  
85 incorporate the entire process of MG data analysis, ranging from preprocessing of sequencing  
86 reads to *de novo* assembly and post-assembly analysis (read alignment, taxonomic  
87 classification, gene annotation, etc.). Both pipelines use SOAPdenovo [22] as the default *de*  
88 *nov* assembler, performing single-length *k*mer-based assemblies which usually result in  
89 fragmented (low contiguity) assemblies with lower coverage levels, compared to multiple-  
90 length *k*mer-based assemblers [23]. However, MetAMOS offers more flexibility by providing  
91 multiple options for MG assemblers, in addition to being an easily installable and user-  
92 amendable pipeline for standardized MG data analysis. *De novo* MG assemblies may be further  
93 leveraged by binning the data to resolve and retrieve population-level genomes including those  
94 from hitherto undescribed taxa [17, 24–28].

95

96 Multi-omic analyses have already allowed unprecedented insights into microbial  
97 community structure and function *in situ*, using different *in silico* analysis approaches.  
98 Franzosa *et al.* (2014) [29] have applied reference-based analyses of MG and MT data to study  
99 healthy human fecal microbial community samples. Conversely, Hultman *et al.* (2015) [30]  
100 performed a multi-omic survey of microbial communities in permafrost soils, leveraging  
101 coupled MG, MT and MP data. The MG data was subjected to *de novo* assembly, annotation  
102 and binning of the resulting contigs. The remaining MT and MP datasets were analyzed based  
103 on the MG assembly and additional reference databases. Bremges *et al.* (2015) [31] performed  
104 analysis on a production-scale bioreactor, using deep MG and MT data. *De novo* assembly and  
105 annotation was carried out for the MG data, after which MT data was used to identify active  
106 metabolic pathways of target organisms. Furthermore, the authors bundled all tools and

107 dependencies within a Docker container to promote reproducibility of their workflow. In  
108 contrast to these aforementioned studies, arguably the most integrated multi-omic study to date  
109 was performed by Muller, Pinel *et al.* (2014) [32], which involved a temporal survey of a  
110 oleaginous microbial community from a biological wastewater treatment plant using a  
111 combination of MG, MT and MP data. In their study, MG and MT data were integrated via a  
112 *de novo* co-assembly procedure, the output of which was subsequently binned and annotated.  
113 The resulting gene set was subsequently used to identify peptides from the MP data. In addition,  
114 patterns of intra- and inter-population expression and genomic variation (single nucleotide  
115 polymorphisms, SNPs) were resolved. Furthermore, Roume, Heintz-Buschart *et al.* (2015) [33]  
116 performed a comparative study of two samples from a biological wastewater treatment plant.  
117 In this work, MG and MT data were also integrated through co-assemblies and the MP data  
118 was analyzed in relation to these co-assemblies. The latter two studies performed integration  
119 by co-assembling MG and MT data. Although these studies clearly demonstrated the power of  
120 multi-omic analyses in facilitating unprecedented insights into community structure and  
121 function, standardized and reproducible dry-lab workflows for integrating and analyzing the  
122 multi-omic data have so far been unavailable. Such approaches are pertinent to compare results  
123 between different studies and systems of study.

124

125 While MG analysis provides information on the gene coding potential (functional  
126 potential) of a given community, complementary usage of MT data enables the study of  
127 transcriptional activity, which more faithfully represents potential community-wide  
128 phenotypes [4, 9, 34]. Due to the absence of tools/workflows to handle multi-omic datasets,  
129 most of the aforementioned studies utilize non-standardized, *ad hoc* analyses, mostly made up  
130 of custom workflows, creating a challenge in reproducing the analyses [12, 35–37]. Here, we  
131 present the Integrated Meta-omic Pipeline (IMP), an open source *de novo* assembly-based

132 pipeline to perform standardized, automated and reproducible large-scale integrated analysis  
133 of multi-omic (MG and MT) data, derived from a single microbial community.

134

## 135 **Overview of IMP**

136 IMP leverages Docker for reproducibility and deployment. The interfacing with Docker is  
137 facilitated through a user-friendly Python wrapper script. As such, Python and Docker are the  
138 only prerequisites for the pipeline, enabling an easy installation and execution process.  
139 Workflow implementation and automation is achieved using Snakemake [38, 39]. The IMP  
140 workflow can be broadly categorized into three major processes: i) preprocessing, ii) assembly  
141 and iii) analysis ([Fig. 1](#)). The modular design and open source features of IMP also allow for  
142 customization of the pipeline to suit specific user-defined analysis requirements. Detailed  
143 parameters for IMP processes are described in the section [Details and parameters of IMP](#) and  
144 examples of detailed workflow schematic is provided in [Additional file 1: HTML S1 and S2](#).

145

### 146 ***Preprocessing of paired-end reads***

147 The input to IMP consists of MG and MT (preferably depleted of ribosomal RNA - rRNA)  
148 paired-end reads in FASTQ format, each comprising a set of two files containing pair-1 and  
149 pair-2 reads, respectively. MG and MT reads are preprocessed independently of each other and  
150 this involves an initial quality control step (see [Fig. 1](#) and section [Trimming and quality](#)  
151 [filtering](#)) [40] followed by filtering of reads that are deemed undesired in downstream  
152 processes. For the analysis of MG and MT data from human microbiome studies, the quality  
153 control steps should be followed by the optional filtering of human genome derived sequences  
154 (see [Fig. 1](#) and section [Human genome derived sequence filtering](#)). Additionally, *in silico*  
155 rRNA sequence depletion is applied to the MT data (see [Fig. 1](#) and section [Ribosomal RNA](#)  
156 [filtering](#)). In summary, the preprocessing implemented in IMP involves the independent

157 removal of sequences from MG and MT datasets which are undesired due to technical and/or  
158 biological reasons.

159

### 160 *IMP-based iterative co-assembly*

161 IMP implements an iterative co-assembly procedure ([Fig. 1](#) and [Additional file 2: Figure S1](#)),  
162 which combines the benefits from iterative assemblies and co-assemblies of multi-omic (MG  
163 and MT) data [32]. Additional rounds of assembly involving the unused reads (unmappable)  
164 from previous assembly are herein referred to as “iterative assemblies”, while assemblies  
165 involving both MG and MT reads are hereafter referred to as “co-assemblies”.

166

167         Preprocessed MT reads are first assembled to generate an initial set of MT contigs  
168 ([Additional file 2: Figure S1](#)). MT reads that remain unmapped to the initial set of MT contigs  
169 undergo an additional round of assembly. The combined set of MT contigs (from both the  
170 aforementioned MT assemblies) are then used as input, together with preprocessed MG and  
171 MT reads, to perform an initial co-assembly. The MG and MT reads which remain unmapped  
172 to the resulting contigs are recruited for an additional co-assembly iteration. The resulting  
173 contigs are further refined by performing a contig-level assembly, which aligns highly similar  
174 contigs against each other. This procedure aims at reducing redundancy by collapsing shorter  
175 contigs into longer contigs and/or improving contiguity by extending contigs via overlapping  
176 contig ends, thereby producing the final set of contigs ([Additional file 2: Figure S1](#)). The IMP-  
177 based iterative co-assembly is a key feature that facilitates the integration of MG and MT data  
178 as it allows maximization of overall data usage. Finally, preprocessed MG and MT reads are  
179 mapped onto the final contigs and the resulting alignment information is used in various  
180 downstream analysis procedures ([Fig. 1](#)). Please refer to section [Details and parameters of IMP](#)  
181 for information about programs and parameters.

182

### 183 *Post-assembly analysis and output*

184 The set of contigs resulting from the IMP-based iterative co-assembly undergo quality  
185 assessment as well as taxonomic [41] and functional annotation [42]. Additionally, non-linear  
186 dimensionality reduction of genomic signatures (NLDR-GS) is performed using VizBin [24,  
187 43] which provides two-dimensional (2D) embeddings, enabling the visualization of the  
188 contigs as scatter plots in a 2D map format. Further analysis steps include, but are not limited  
189 to, calculations of the contig- and gene-level depths of coverage and the calling of genomic  
190 variants (using three variant callers, see section [Variant calling](#)). The information from these  
191 analyses are condensed and integrated into the VizBin-based maps to produce augmented  
192 visualizations. The visualizations and various summaries of the output are compiled into a  
193 HTML report (for examples of HTML report see [Additional file 1: HTML S1 and S2](#) and for  
194 details of output see section [Output](#)).

195

## 196 **Results and discussion**

197 We demonstrate the performance and output of IMP on three multi-omic datasets, each  
198 consisting of MG and MT paired-end reads (see section [Coupled metagenomic and](#)  
199 [metatranscriptomic datasets](#) for details). A simulated mock (SM) community of 73 bacterial  
200 genomes [14] was mainly used to benchmark the IMP-based iterative co-assemblies in  
201 comparison to standard assembly strategies. Additional benchmarking of the iterative co-  
202 assemblies were performed on published datasets from a human fecal (HF) sample [29] and a  
203 wastewater (WW) sludge microbial community [32]. The latter datasets were used to assess  
204 the output and features of IMP ranging from preprocessing to post-assembly analyses.

205

### 206 *Preprocessing and filtering*

207 The preprocessing and filtering of sequencing reads is essential for the removal of low quality  
208 bases/reads and potentially unwanted sequences, prior to assembly and analysis. Preprocessing  
209 of NGS data, prior to assembly, has been shown to increase the quality of *de novo* assemblies,  
210 despite decreased numbers of input reads [44]. The preprocessing of MG and MT reads is  
211 handled in a tailored manner within IMP.

212

213 The results of the IMP and MetAMOS preprocessing procedures are summarized in  
214 [Additional file 3](#): Table S1. The preprocessing of the HF included the optional filtering of  
215 human genome derived-sequences, while the same step was omitted for the WW data. Even  
216 though the preprocessing using IMP yields both paired- and single-end reads, the following  
217 section discusses only the paired-end reads as they make up a large fraction of the data and are  
218 generally more informative compared to single-end reads. However, IMP retains the use of  
219 single-end reads in all downstream processes, unlike most available methods, which discard  
220 such reads. Since MetAMOS assumes all input to be MG data, it cannot be directly compared  
221 against the preprocessing using IMP.

222

223 The final output of IMP's preprocessing and filtering procedure retained 69.2% (29.8%  
224 low quality and 1.0% human genome derived sequences) and 89.2% (10.8% low quality) of  
225 MG paired-end reads for HF and WW, respectively. Similarly, approximately 90.8% (6.4%  
226 low quality, 1.2% rRNA and 1.7% human genome derived sequences) and 55.3% (19.8% low  
227 quality and 24.9% rRNA) of MT paired-end reads were retained from HF and WW,  
228 respectively. The filtering of human genome derived sequences are due to technical and privacy  
229 reasons, whereas *in silico* rRNA filtering helps remove remaining rRNA reads, which are  
230 usually abundant in cDNA libraries, even after the application of wet-lab rRNA depletion  
231 procedures. In summary, IMP is designed to perform stringent and standardized preprocessing

232 of MG and MT data in a tailored way enabling efficient data usage in subsequent steps.

233

### 234 *Assessment of the iterative assembly approach*

235 *De novo* assemblies of MG or MT data usually result in a large fraction of reads that are  
236 unmappable to the produced contigs and therefore remain unused, resulting in suboptimal data  
237 usage. Previous studies (e.g. Muller, Pinel *et al.* 2014 [32]; Schürch *et al.* 2014 [45]; Reyes *et*  
238 *al.* 2015 [46]) assembled the set of unmappable reads iteratively to successfully obtain  
239 additional contigs from these additional rounds of assembly and which lead to an increase in  
240 number of predictable genes.

241

242 In order to evaluate the best iterative assembly approach for IMP, we attempted to  
243 determine the opportune number of assembly iterations in relation to assembly quality metrics.  
244 The evaluation involved performing multiple iterations of recruiting unmappable reads to the  
245 previously generated non-redundant assembly, followed by a *de novo* assembly of those  
246 unmapped reads. The assembly from a given iteration was then merged with the previous  
247 assembly to reduce the redundancy (refer to section [Iterative single-omic assemblies](#) for  
248 details). The evaluation of additional assembly iterations for MG data of SM, HF and WW are  
249 summarized in [Fig. 2](#), based on four different metrics. Overall, each iteration on each of the  
250 different datasets (SM, HF and WW) lead to an increase in total length of assembly and  
251 increased the overall number of mappable reads, but differed in the observed gain of contigs  
252 and genes ([Fig. 2](#); and [Additional file 3](#): Table S2). A similar trend is noticeable for iterative  
253 assemblies on MT data (see [Additional file 2](#): Figure S2 and [Additional file 3](#): Table S3). The  
254 observed trends may be explained by the fact that the complexity of the data typically  
255 confounds assemblies [44]. The exclusion of mappable reads in each iteration of assembly

256 reduces the complexity of the data, which in turn allows additional contigs to be assembled and  
257 results in a higher cumulative output [44].

258

259         Considering the relatively low increase in longer contigs and genes beyond the first  
260 assembly iteration ([Fig. 2](#), [Additional file 2](#): Figure S2 and [Additional file 3](#): Table S2 and S3)  
261 and the extended runtimes required to perform additional assembly iterations, an opportune  
262 single iteration assembly approach was implemented in the workflow of IMP. This approach,  
263 which balances the maximization of output yield with runtime, is implemented within the IMP-  
264 based co-assembly approach.

265

## 266 ***Benchmarking***

267 The iterative co-assemblies were benchmarked against single-omic MG and MT assemblies  
268 and co-assemblies obtained using the state-of-the-art MG data analysis pipeline, MetAMOS  
269 [12].

270

### 271 *Single-omic assemblies and multi-omic iterative co-assemblies*

272 Separate single-omic iterative assemblies on all datasets were generated using the preprocessed  
273 MG and MT data (see [Iterative single-omic assemblies](#)). The iterative co-assemblies were  
274 executed in IMP using the two available assembler options for the co-assembly step, i.e. the  
275 default IDBA-UD [15] (hereafter referred to as IMP) and the optional MEGAHIT assembler  
276 [16] (referred to as IMP-MEGAHIT). Both assemblers are regarded as state-of-the-art because  
277 they perform assemblies on multiple *kmer* sizes, while MEGAHIT was also chosen due to its  
278 superior speed and efficient memory usage [15, 16].

279

280           The comparison between single-omic assemblies and IMP-based co-assemblies are  
281 summarized in [Table 1](#). The IMP-based co-assemblies consistently returned larger number of  
282 contigs, increased total length of assembled contigs and a higher number of predicted genes  
283 (partial genes included) compared to single-omic assemblies for all datasets ([Table 1](#)). The  
284 apparent slight reduction in contiguity (N50 statistic) is due to the addition of shorter contigs  
285 likely stemming from the increased sequencing depth of the combined MG and MT datasets,  
286 which also increases the complexity of the assembly process. By using the reference genomes  
287 from the SM data as ground truth, an improved recovery of reference genome fractions is  
288 apparent for the IMP-based co-assemblies. Importantly, a significant increase in the number of  
289 mappable MG and MT reads was observed within all co-assemblies compared to the respective  
290 single-omic assemblies ([Table 1](#)) which suggests superior data usage using the IMP-based  
291 approach. For example, the IMP-based iterative co-assemblies resulted in a large fraction of  
292 reads being mappable back to the contigs derived from the HF sample (average of  
293 approximately 88.0 % and 96.3 % for the MG and MT reads, respectively; [Table 1](#)), which is  
294 substantially higher compared to the numbers reported in a previous report in which MG  
295 sequencing data was mapped to an integrated gene catalog, i.e. 74%-81% [11]. In summary,  
296 the complementary use of MG and MT data in the context of *de novo* co-assembly results in  
297 an increased yield of output, while enhancing overall data usage for subsequent analyses.

298

### 299 *Quality assessment of the IMP-based iterative co-assembly procedure*

300 Iterative co-assemblies using IMP (referred to as IMP and IMP-MEGAHIT, see section above)  
301 and co-assemblies from MetAMOS [12] on all datasets were compared against each other.  
302 MetAMOS was chosen due to its similar aim of providing an open source, reproducible and  
303 standardized *de novo* assembly-based platform for large-scale microbiome sequencing  
304 analyses [12]. Although MetAMOS was developed specifically for MG data analysis, it was

305 hereby extended to perform MG and MT co-assemblies by including both MG and MT read  
306 libraries as input (see section [Execution of pipelines](#)) using two available assembler options:  
307 SOAPdenovo [22] (hereafter referred to as MetAMOS) and IDBA-UD [15] (hereafter referred  
308 as MetAMOS-IDBA\_UD). The results of the comparison are summarized in [Fig. 3](#) (see  
309 [Additional file 2](#): Figure S3 and [Additional file 3](#): Table S4 for detailed comparison and  
310 results).

311

312         Based on the SM data ([Fig. 3A](#)), IMP and MetAMOS-IDBA\_UD performed similarly  
313 for most measures but the IMP-based assemblies producing slightly better contiguity (N50  
314 statistic) and lower levels of apparent misassembly. Despite MetAMOS producing the largest  
315 number of contigs  $\geq$  1kb, its comparatively low N50 statistic indicates a highly fragmented  
316 assembly, which is further reflected in the low number of predicted unique genes. Conversely,  
317 this fragmented assembly is accompanied by relatively low misassembly rate, reinforcing the  
318 notion that shorter contigs are less prone to misassemblies [44]. However, longer contigs ( $\geq$   
319 1kb) are a prerequisite for population-level genome reconstruction and subsequent multi-omic  
320 data interpretation. On the other hand, IMP-MEGAHIT generated the highest number of  
321 predicted unique genes, recovered the largest fraction of reference genomes, while yielding  
322 comparatively large number of contigs  $\geq$  1kb, relatively high N50 statistic and a low rate of  
323 misassembly. The assessment based on real datasets shows comparable performance between  
324 IMP, IMP-MEGAHIT and MetAMOS-IDBA\_UD ([Fig. 3B](#) and [C](#)). In general, MetAMOS  
325 produced highly fragmented assemblies for the real datasets, possibly due to the single *k*mer  
326 length ( $k = 31$ ) assemblies, which tend to produce relatively fragmented assemblies compared  
327 to multiple *k*mer length assemblers (as in IMP, IMP-MEGAHIT and MetAMOS-IDBA\_UD)  
328 [23]. In summary ([Fig. 3D](#)), IMP and MetAMOS-IDBA\_UD performed similarly for most  
329 metrics when the same assembly program (IDBA-UD) was used by both pipelines. However,

330 the IMP iterative co-assemblies were generated using a lower number of reads compared to  
331 MetAMOS-IDBA\_UD due to the more stringent preprocessing procedures in IMP, which in  
332 turn yielded better quality assemblies (Fig. 3D) which are a prerequisite for population-level  
333 genome reconstruction and multi-omic data interpretation.

334

### 335 *Summary output from IMP*

336 The workflow of IMP is unique such that it allows integrated MG and MT data handling.  
337 Although MetAMOS may be extended to perform co-assemblies of MG and MT data, it does  
338 not discriminate between the two data types in its pre- and post-assembly procedures which is  
339 important given the disparate nature of MG and MT datasets.

340

341 IMP generates several output files, as detailed in the [Output](#) section which allow both  
342 reference-dependent and –independent analyses of the data. Information from these output files  
343 are condensed and summarized using different static and dynamic visualization methods,  
344 which are compiled into an HTML report ([Additional file 1: HTML S1 and S2](#)). [Fig. 4](#) presents  
345 selected output available from the analysis of IMP on the HF data. The generated taxonomic  
346 overviews are based on the alignment of contigs to the most closely related prokaryotic  
347 genomes within the NCBI genome database and the fraction of potential reference genome  
348 bases covered ([Fig. 4A](#); [Additional file 1: HTML S1](#)) [41]. The abundances of the predicted  
349 genes (based on average depths of coverage) may be represented both at the MG and MT levels  
350 and thus enable the comparison of functional potential ([Fig. 4B](#)) and actual expression ([Fig.](#)  
351 [4C](#)) of various KEGG Ontology categories (for details, see Krona charts within [Additional file](#)  
352 [1: HTML S1](#)). IMP also provides augmented VizBin-based 2D maps [24, 43], with additional  
353 layers of information integrated onto them, for example, variant densities ([Fig. 4D](#)) and MT to  
354 MG depth of coverage ratios ([Fig. 4E](#)). These visualizations may aid users in highlighting

355 subsets of contigs based on certain characteristics of interest, i.e. population  
356 heterogeneity/homogeneity, low/high transcriptional activity, etc. Please refer to [Additional](#)  
357 [file 1](#): HTML S1 and S2 for further examples.

358

### 359 *Integrated omic data allows identification of key microbiome characteristics*

360 The integration of MG and MT data provides unique opportunities for uncovering community-  
361 or population-specific traits which cannot be resolved from MG data alone. Here we provide  
362 two examples of insights gained through the direct comparison of the MG and MT results  
363 provided by IMP.

364

#### 365 *Identification of RNA viruses*

366 To identify differences in the information content of MG and MT complements, the contigs  
367 generated from IMP were inspected with respect to coverage by MG and MT reads. In the two  
368 exemplary datasets, a large fraction of the contigs resulted from the composite assembly of MG  
369 and MT data, followed by contigs composed exclusively of MG data and a small proportion  
370 composed exclusively of MT data ([Additional file 3](#): Table S5). Longer contigs ( $\geq 1$  kb)  
371 composed exclusively of MT reads and annotated with known viral/bacteriophage genes were  
372 enriched and retained for further inspection [Table 2](#) (for the complete list contigs, see  
373 [Additional file 3](#): Table S6 and S7). A sequence alignment-based search against the NCBI nr  
374 nucleotide database revealed that the longer contigs represent almost complete genomes of  
375 RNA viruses (see [Additional file 3](#): Table S8 and S9). This demonstrates that the incorporation  
376 of MT data in the assembly enables the recovery of nearly complete RNA virus genomes,  
377 thereby allowing their detailed study.

378

#### 379 *Identification of populations with apparent high transcriptional activity*

380 To further demonstrate the unique analytical capabilities of IMP, we set out to identify  
381 populations with a high transcriptional activity in the HF sample. Average depth of coverage  
382 at the contig- and gene-level is a common measure used to evaluate the abundance of microbial  
383 populations within communities [25, 27, 32]. Integrative analysis of MG and MT data by IMP  
384 further extends this measure by calculation of average MT to MG depth of coverage ratios,  
385 which provides information on transcriptional activity and can be visualized using augmented  
386 VizBin maps.

387

388 One particular cluster of contigs within the augmented VizBin maps displayed high MT  
389 to MG depth of coverage ratios (see [Additional file 2: Figure S4](#)). The subset of contigs  
390 (simplified as subset) within the selected cluster aligned to the genome of the *Escherichia coli*  
391 P12B. For comparison, we also identified a subset which was highly abundant at the MG level  
392 (lower MT to MG ratio), which aligned to the genome of *Collinsella intestinalis* DSM 13280  
393 strain. Based on these observations, we highlighted these subsets of contigs to produce an  
394 augmented VizBin map ([Fig. 5A](#)). In these, *C. intestinalis* and *E. coli* subsets are mainly  
395 represented by clear peripheral clusters which exhibit consistent intra-cluster MT to MG depth  
396 of coverage ratios ([Fig. 5A](#)). The subsets were manually inspected in terms of their distribution  
397 of average MG and MT depths of coverage, comparing them against the corresponding  
398 distributions of all the contigs. The MG-based average depths of coverage of the contigs from  
399 the entire community exhibited a bell-shape like distribution, with a clear peak (mode). On the  
400 contrary, MT depths of coverage exhibited a spread distribution, with a relatively lower mean  
401 (compared to MG distribution) and no clear peak ([Fig. 5B](#)). The *C. intestinalis* subset displays  
402 similar distributions to that of the entire community, whereas the *E. coli* subset exhibits an  
403 unusually high MT-based depth of coverage, and a low MG-based depth of coverage ([Fig. 5B](#)).  
404 Further inspection of the individual omic datasets revealed that the *E. coli* subset was not

405 covered by the MG-based contigs, while 70 % of the *E. coli* genome was recoverable from the  
406 MT-based assembly (Fig.5C). In contrast, the *C. intestinalis* subset demonstrated comparable  
407 genomic recovery in all co-assemblies and the MG-only assembly. As noted by the authors of  
408 the original study by Franzosa *et al.* (2014) [29], the cDNA conversion protocol used to  
409 produce the MT data is known to introduce approximately 1-2 % of *E. coli* genomic DNA into  
410 the cDNA as contamination which is then reflected in the MT data. According to our analyses,  
411 0.12 % of MG reads and 1.95 % of MT reads from this sample could be mapped onto the *E.*  
412 *coli* contigs which is consistent with the numbers quoted by Franzosa *et al.* (2014) [29]. This  
413 fraction of reads is sufficient for *de novo* reconstruction of approximately 70% of the *E. coli*  
414 genome. The integrative analyses of MG and MT data within IMP enables users to  
415 conveniently highlight notable cases such as this, and to further investigate inconsistencies  
416 and/or interesting characteristics within the multi-omic data.

417

## 418 **Summary**

419 IMP was developed in order to leverage the advantages associated with integrating MG and  
420 MT data for studying microbial community structure and function *in situ* [4, 6]. Accordingly,  
421 we present a self-contained workflow which performs reproducible integrative analyses of  
422 coupled MG and MT data derived from single and unique microbial community samples. IMP  
423 encapsulates all processes including preprocessing, assembly, and analyses within an  
424 automated reproducible pipeline.

425

426 We implemented customized preprocessing and filtering procedures for MG and MT  
427 data due to the distinct nature of these different omic data types. We also evaluated the IMP-  
428 based iterative co-assembly procedure and found it to produce higher amount of output volume  
429 (a higher number of contigs and genes), thereby resulting in enhanced data usage (reflected in

430 a higher fraction of read which can be mapped back to the contigs). IMP provides the option  
431 for the use of two state-of-the-art assemblers, whereby the default assembler, IDBA-UD,  
432 produces highly contiguous assemblies, while MEGAHIT balances the number of contigs with  
433 favorable contiguity with a high number of predicted genes and a relatively low rate of  
434 misassemblies. High quality assemblies yield better quality taxonomic information and gene  
435 annotations. Consequently, the post-assembly analysis of MG and MT data enables users to  
436 evaluate their community of interest based on taxonomy, functional potential and functional  
437 expression. The integrated co-assembly also provides the opportunity for analyses not possible  
438 based on MG data alone, such as the detection of RNA viruses and the identification of  
439 transcriptionally active populations. The output of IMP is compatible with, and may be  
440 exported to interactive tools such as VizBin [43] and Anvi'o [17] for binning and further  
441 analyses. Furthermore, the output data (annotated gene sets) may be used for the analysis and  
442 integration of additional omic data, most notably MP data.

443

444 The use of the Docker promotes reproducibility and sharing such that researchers are  
445 able to tag specific versions of the pipeline used for a particular study, thus enabling their peers  
446 to precisely replicate bioinformatic analyses workflows with relative ease and with minimal  
447 impact on overall performance of the employed bioinformatic tools [36, 37]. Static websites  
448 will be created and associated with every new version of IMP (Docker image), such that users  
449 will be able to download and launch specific versions of the pipeline to reproduce the work of  
450 others. Finally, the open source nature of IMP encourages a community-based effort to  
451 contribute and further improve the pipeline. The common scripting languages Bash, Make and  
452 Python, which Snakemake is based on [38, 39], are widely used within the bioinformatic  
453 community, thus reducing the learning curve for further development, improvement and  
454 customization by other users. Finally, the combination of open development and

455 reproducibility should promote the general paradigm of reproducible research within the  
456 microbiome research community.

457

## 458 **Details and parameters of IMP**

459 A Python (ver 3) wrapper script was implemented for user-friendly execution of IMP via the  
460 command line. The full list of dependencies, parameters (see below) and documentation are  
461 available on the IMP website (<http://r3lab.uni.lu/web/imp/>).

462

### 463 ***Reproducibility***

464 IMP is based around a Docker container that runs the Ubuntu 14.04 operating system, with all  
465 relevant dependencies. Five mounting points are defined for the Docker container with the -v  
466 option: i) input directory, ii) output directory, iii) database directory, iv) code directory, and v)  
467 configuration file directory. Environment variables are defined using the -e parameter,  
468 including: i) paired MG data, ii) paired MT data, and iii) configuration file. The latest IMP  
469 Docker image will be downloaded automatically upon launching the command, but users may  
470 also launch specific versions based on tags or use modified/customized versions of their local  
471 code base.

472

### 473 ***Automation***

474 Automation of the workflow is achieved using Snakemake 3.4.2, a Python-based make  
475 language implemented specifically for building reproducible bioinformatic workflows and  
476 pipelines. It allows seamless integration of Python code and shell (bash) commands, using  
477 make scripting style. It also provides checkpoints to continue interrupted analyses and/or rerun  
478 steps if required [38, 39].

479

480 ***Trimming and quality filtering***

481 Trimmomatic 0.32 [40] is used to perform trimming and quality filtering of MG and MT  
482 Illumina paired-end reads, using the following parameters: ILLUMINACLIP:TruSeq3-  
483 PE.fa:2:30:10; LEADING:20; TRAILING:20; SLIDINGWINDOW:1:3; MAXINFO:40:0.5;  
484 MINLEN:40. The parameters may be tuned in the IMP config file. The output from this step  
485 includes retained paired-ends and single-ends (mate discarded) which are all used for  
486 downstream processes.

487

488 ***Ribosomal RNA filtering***

489 SortMeRNA 2.0 is used for filtering rRNA from the MT data. The process is applied on  
490 FASTQ files for both paired- and single-end reads generated from the previous preprocessing  
491 step. Paired-end FASTQ files are interleaved prior to running SortMeRNA. If one of the mates  
492 within the paired-end read are classified as an rRNA sequence, then the entire pair is filtered  
493 out. After running SortMeRNA, the interleaved paired-end output is split into two separate  
494 paired-end FASTQ files. The filtered sequences (without rRNA) are used for the downstream  
495 processes. All available databases provided within SortMeRNA are used for filtering and the  
496 maximum memory usage parameter is set to 4 Gb (option: -m 4000).

497

498 ***Read mapping***

499 The read mapping procedure is performed using bwa mem aligner [47] with settings: -v 1  
500 (verbose output level), -M (Picard compatibility) introducing an automated samtools header  
501 using the -R option [47]. Paired- and single-end reads are mapped separately, and the resulting  
502 alignments are merged (using samtools merge). Read mapping is performed at various steps in  
503 the workflow including: i) filtering human sequences (optional), ii) recruitment of unmapped

504 reads within the IMP-based iterative co-assembly, and iii) mapping of preprocessed MG and  
505 MT reads to the final contig set.

506

### 507 *Extracting unmapped reads*

508 The extraction of unmapped reads (paired- and single-end) begins by mapping reads to a given  
509 reference sequence (see section [Read mapping](#)). The resulting alignment file (BAM format) is  
510 used as input for the extraction of unmapped reads. A set of paired-end reads considered  
511 unmappable if both or either one of the mates do not map to the given reference. The unmapped  
512 reads are converted from BAM to FASTQ format using samtools and BEDtools 2.17.0 -  
513 bamToFastq utility [48]. Similarly, unmapped single-end reads are also extracted from the  
514 alignment information.

515

### 516 *Filtering host sequences*

517 The human sequence filtering is performed by mapping the both paired- and single-end reads  
518 (see section [Read mapping](#)) onto the human genome version 38  
519 (<http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/>), followed by extraction of  
520 unmapped reads (see section [Extracting unmapped reads](#) for details). This filtering step may  
521 be omitted from the workflow using the IMP user configuration file, while users may replace  
522 the human genome with other FASTA sequences from other hosts based on their screening  
523 requirements.

524

### 525 *Parameters of the IMP-based iterative co-assembly*

526 The IMP-based iterative co-assembly implements MEGAHIT 1.0.3 [16] as the MT assembler  
527 while IDBA-UD 1.1.1 [15] is used as the default co-assembler (MG + MT), with MEGAHIT  
528 [16] as an alternative option for the co-assembler. All *de novo* assemblies are performed on

529 kmers ranging from 25-mers to 99-mers, with an incremental step of four. Accordingly, the  
530 command line parameters for IDBA-UD are --mink 25 --maxk 99 --step 4 --similar 0.98 --pre-  
531 correction [15]. Similarly, the command line options for MEGAHIT are --k-min 25 --k-max  
532 99 --k-step 4, except for the MT assemblies which are performed with an additional --no-bubble  
533 option to prevent merging of bubbles within the assembly graph [16]. Furthermore, contigs  
534 generated from the MT assembly (“pseudo long-reads”) are used as input within the -l flag of  
535 IDBA-UD or -r flag of MEGAHIT [15, 16]. The parameters for the assemblies may be adjusted  
536 in the user configuration file. Kmer ranges may be customized in the IMP user configuration  
537 file.

538

539

#### 540 *Annotation and assembly quality assessment*

541 Prokka 1.11 [42] with the --metagenome setting is used to perform functional annotation. The  
542 default BLAST and HMM databased of Prokka are used for the functional annotation.

543

544 MetaQUAST 3.1 [41] is used to perform taxonomic annotation of contigs with the  
545 maximum number of reference genomes set to 75 (--max-ref-number 75). In addition,  
546 MetaQUAST provides various assembly statistics.

547

#### 548 *Depth of coverage*

549 Contig- and gene-wise depth of coverage values are calculated (per base) using BEDtools  
550 2.17.0 [48] and aggregated (by average) using awk, adapted from the CONCOCT code [27]

551 (script: `map-bowtie2-markduplicates.sh`, Github URL: <https://github.com/BinPro/CONCOCT>).

553

554 ***Variant calling***

555 The variant calling procedure is performed using the following tools: i) Samtools mpileup  
556 0.1.19 [49], ii) Freebayes 0.9.21 [50] and iii) Platypus 0.8.1 [51], each using their respective  
557 default settings. The input is the merged paired- and single-end read alignment (BAM) against  
558 the final assembly FASTA file (see section [Read mapping](#)). The output files from the three  
559 methods are indexed using tabix and compressed using gzip. No filtering is applied to the  
560 variant calls, so that users may access all the information and filter them according to their  
561 requirements. The output from samtools mpileup is used for the VizBin-based visualizations.

562

563 ***Non-linear dimensionality reduction of genomic signatures (NLDR-GS)***

564 VizBin [43] performs NLDR-GS onto contigs  $\geq$  1kb, using default settings, to obtain 2D  
565 embeddings. Parameters can be modified in the IMP config file.

566

567 ***Visualization and reporting***

568 IMP compiles the multiple summaries and visualizations into a HTML report. FASTQC [52]  
569 is used to visualize the quality and quantity of reads before and after preprocessing.  
570 MetaQUAST [41] is used to report assembly quality and taxonomic associations of contigs. A  
571 custom script is used to generate KEGG-based [53] functional Krona plots by running  
572 KronaTools [54] (script: `genes.to.kronaTable.py`, GitHub URL:  
573 <https://github.com/EnvGen/metagenomics-workshop>). Additionally, VizBin output is  
574 integrated with the information derived from the IMP analyses, using a custom R script for  
575 analysis and visualization of the augmented maps. The R workspace (Rdat) is saved such that  
576 users are able to access it for further analyses. All the steps executed within an IMP run  
577 including parameters and runtimes are summarized in the form of a workflow diagram and a  
578 log-file.

579

## 580 *Output*

581 The output generated by IMP includes a multitude of large files. Paired- and single-end FASTQ  
582 files of preprocessed MG and MT reads are provided such that the user may employ them for  
583 additional downstream analyses. The output of the IMP-based iterative co-assembly consists  
584 of a FASTA file, while the alignments/mapping of MG and MT preprocessed reads to the final  
585 co-assembly are also provided as a binary alignment format (BAM), such that users may use  
586 these for further processing. Predicted genes and their respective annotations are provided in  
587 the various formats produced by Prokka [42]. Assembly quality statistics and taxonomic  
588 annotations of contigs are provided as per the output of MetaQUAST [41]. Two-dimensional  
589 embeddings from the NLDR-GS are provided such that they can be exported to and further  
590 curated using VizBin [43] for human-augmented binning. Additionally, abundance and  
591 expression information is represented by contig- and gene-level average depth of coverage  
592 values. MG and MT genomic variant information (VCF format), including both SNPs and  
593 INDELs (insertions and deletions), is also provided.

594

595 The HTML report (see [Additional file 1: HTML S1 and S2](#)) compiles various  
596 summaries and visualizations including, i) augmented VizBin-based maps, ii) MG- and MT-  
597 level functional Krona charts [54], iii) run time information iv) detailed schematic of the steps  
598 carried out within the IMP run, v) list of parameters and commands, and vi) additional reports  
599 [FASTQC report [52], MetaQUAST report [41]]. Please refer to documentation of IMP for a  
600 detailed list and description of the output.

601

## 602 *Customization and further development*

603 Basic parameters (input, output, assembler, configuration file) may be changed via the IMP  
604 command line, while more advanced parameters may be changed by editing the user  
605 configuration file (JSON format). Finally, users may directly edit the code of IMP to implement  
606 extensive changes to the pipeline, if required. The “--current” flag, within the IMP command  
607 line can be used to execute customized (local) versions of the code base. Finally, the IMP  
608 launcher script provides the option (flag: --enter) to launch the Docker container interactively,  
609 for development and testing purposes (described on the IMP website and documentation).

610

## 611 **Data and analyses**

### 612 *Coupled metagenomic and metatranscriptomic datasets*

613 The simulated MT data was obtained from the original authors [14], upon request. A  
614 complementary metagenome was simulated using the same set of 73 bacterial genomes used  
615 for the aforementioned simulated MT [14]. Simulated reads were obtained using the NeSSM  
616 MG simulator (default settings) [55].

617

618 Real data analyzed for this article included coupled MG and MT data, i.e. both datasets  
619 were obtained from the same unique sample. The published human fecal data derived from the  
620 healthy individual “X310763260” [29] was obtained from the NCBI Sequence Read Archive  
621 (metagenome SRA: SRX247379, metatranscriptome SRA:SRX247335). The wastewater  
622 sludge data was obtained in-house, but is available on the NCBI SRA (metagenome SRA:  
623 SRX389533, metatranscriptome SRA: SRX389534) [32].

624

### 625 *Iterative single-omic assemblies*

626 In order to determine the opportune number of iterations within the IMP-based iterative co-  
627 assembly strategy within IMP, we first performed an initial assembly using IMP preprocessed

628 SM, HF and WW MG reads with IDBA-UD [15] together with cap3 [56], which was used to  
629 further collapse the contigs and reduce the redundancy of the assembly. This initial assembly  
630 was followed by a total of three assembly iterations. Each iteration was made up of four  
631 separate steps: i) extraction of reads unmappable to the previous assembly (using the procedure  
632 explained in [Extracting unmapped reads](#)), ii) assembly of unmapped reads using IDBA-UD  
633 [15], iii) merging/collapsing the contigs from the previous assembly using cap3 [56], and iv)  
634 evaluation of the merged assembly using MetaQUAST [41]. The assembly was evaluated in  
635 terms of the per-iteration increase in mappable reads, assembly length, numbers of contigs  $\geq$   
636 1 kb, and number of unique genes.

637

638 Similar iterative assemblies were also performed for MT data of SM, HF and WW using  
639 MEGAHIT [16] except, CD-HIT-EST [57] was used to collapse the contigs at  $\geq 95\%$  identity  
640 ( $-c 0.95$ ) while MetaGeneMark [58] was used to predict genes. The parameters and settings of  
641 the other programs were the same as those defined in [Details and parameters of IMP](#).

642

643 The contigs from the first iteration of both the MG and MT iterative assemblies were  
644 selected to represent the control single-omic (MG-only and MT-only) assemblies and were  
645 compared against co-assemblies.

646

### 647 *Execution of pipelines*

648 MetAMOS was executed on each dataset using: i) the default setting, with SOAPdenovo as  
649 assembler (using a length of 31-mers), and ii) a custom version with IDBA-UD as the  
650 assembler (option: `-a idba-ud`) using both MG and MT paired-end FASTQ reads as input. All  
651 computations using MetAMOS were set to use eight computing cores per run (option: `-p 8`).

652

653 Similarly, IMP was executed for each dataset using different assemblers for the co-assembly  
654 step: i) default setting using IDBA-UD, and ii) MEGAHIT (option: -a megahit). Additionally,  
655 the analysis of HF data included the preprocessing step of filtering human genome sequences,  
656 which was omitted for WW data. Illumina TruSeq2 adapter trimming was used for WW data  
657 preprocessing, since the information was available. Computation was performed using eight  
658 computing cores. The customized parameters were specified in the IMP configuration file (see  
659 [Additional file 1: HTML S1 and S2 for exact configurations](#)).

660

### 661 *Assembly assessment and comparison*

662 Assemblies were assessed and compared at the contig level (scaffolds not considered), using  
663 MetaQUAST [41]. The gene calling function (flag: -f) was utilized to obtain the number of  
664 genes which were predicted from the various assemblies. An additional parameter within  
665 MetaQUAST was used for ground truth assessment of the simulated mock (SM) community  
666 assemblies, by providing the list of 73 reference genomes (flag: -R). MetaQUAST was applied  
667 to compare: i) single-omic assemblies and multi-omic co-assemblies and ii) co-assemblies  
668 from different pipelines.

669

### 670 *Analysis of contigs assembled from MT data*

671 A list of contigs with no MG depth of coverage together with additional information on these  
672 (contig length, annotation, MT depth of coverage) was retrieved using the R workspace (Rdat)  
673 which is provided as part IMP output. The sequences of these contigs were extracted and  
674 subjected to a BLAST search on NCBI to determine their potential origin. Furthermore, contigs  
675 with length  $\geq 1$  kb, average depth of coverage  $\geq 20$  bases and containing genes encoding known  
676 virus/bacteriophage functions were extracted.

677

## 678 *Analysis of subsets of contigs*

679 Subsets of contigs were identified by visual inspection of augmented VizBin maps generated  
680 by IMP. Detailed inspection of contig-level MT to MG depth of coverage ratios was carried  
681 out using the R workspace provided as part of IMP output. The alignment information of  
682 contigs to isolate genomes provided by MetaQUAST [41] were used to highlight subsets of  
683 contigs aligning to genomes of *Escherichia coli* P12B strain (*E. coli*) and *Collinsella*  
684 *intestinalis* DSM 13280 (*C. intestinalis*).

685

686 MetaQUAST [41] was used to compare the three co-assemblies carried out on the HF  
687 dataset (using IMP, IMP-MEGAHIT and MetAMOS-IDBA\_UD) against the corresponding  
688 single-omic MG and MT assemblies (see section Iterative single-omic assemblies). For the HF  
689 data, corresponding reference genomes were extracted from the IMP output and were provided  
690 to MetaQUAST (flag: -R) as the reference genome set.

691

## 692 *Computational platforms*

693 IMP and MetAMOS were executed on a Dell R820 machine with 32 Intel(R) Xeon(R) CPU  
694 E5-4640 @ 2.40GHz physical computing cores (64 virtual), 1024 TB of DDR3 RAM (32 GB  
695 per core) with Debian 7 Wheezy as the operating system. Additional computations outside the  
696 scope of the pipelines (IMP and MetAMOS) were performed on the Gaia cluster of the  
697 University of Luxembourg HPC platform [59].

698

## 699 **Availability**

700 IMP software and code are available under the BSD-4-Clause license, on the LCSB R3 website:  
701 <http://r3lab.uni.lu/web/imp/>. Scripts and commands for additional analyses are available at:  
702 [https://git-r3lab.uni.lu/shaman.narayanasamy/IMP\\_article\\_analyses](https://git-r3lab.uni.lu/shaman.narayanasamy/IMP_article_analyses).

## 703 **Competing interests**

704 The authors declare that they have no competing interests

705

## 706 **Authors' contributions**

707 SN, NP, EELM, PM, PW conceived the analysis and designed the workflow. SN, YJ, MH,  
708 CCL developed the software, wrote the documentation and tested the software. YJ ensured  
709 reproducibility of the software. SN performed data analyses. EELM, PM, AHB, AK, NP, PW  
710 participated in discussions and tested the software. SN, EELM, AHB, PM, NP, AK, MH, PW  
711 wrote and edited the manuscript. PW designed and supported the project. All authors read and  
712 agreed on the final version of the manuscript.

713

## 714 **Abbreviations**

715 NGS: next-generation sequencing

716 Contigs: contiguous sequence(s)

717 cDNA: complementary-DNA

718 MG: Metagenomic

719 MT: Metatranscriptomic

720 MP: Metaproteome/metaproteomic

721 SNPs: single nucleotide polymorphisms

722 INDELS: insertions and deletions

723 rRNA: ribosomal RNA

724 IMP: Integrated Meta-omic Pipeline

725 SM: Simulated mock

726 HF: Human fecal

727 WW: Wastewater  
728 bp: base pair  
729 Kb: kilo base  
730 KEGG: Kyoto Encyclopedia of Genes and Genomes  
731 VCF: Variant call format  
732 SRA: Sequence read archive  
733 NCBI: National Center for Biotechnology Information

734

## 735 **Acknowledgements**

736 We would like to acknowledge John Larsson from SciLifeLab (Sweden) for kindly providing  
737 the KEGG-based functional Krona plot scripts. Albi Celaj from the University of Toronto is  
738 thanked for supplying the *in silico* simulated metatranscriptomic data and the corresponding  
739 reference genomes. The University of Luxembourg High Performance Computing (HPC)  
740 facility is duly thanked for providing and maintaining the computing platform. The  
741 Reproducible Research Results (R3) team of the Luxembourg Centre for Systems Biomedicine  
742 is acknowledged for support of the project and for promoting reproducible research. This work  
743 was supported by an ATTRACT program grant (A09/03), a European Union Joint  
744 Programming in Neurodegenerative Diseases grant (INTER/JPND/12/01) and a proof of  
745 concept grant (PoC/13/02) to PW, Aide à la Formation Recherche (AFR) grants to SN (PHD-  
746 2014-1/7934898), and a CORE junior (C15/SR/10404839) to EELM, all funded by the  
747 Luxembourg National Research Fund (FNR).

748

## 749 **References**

- 750 1. Turnbaugh PJ, Ley RE, Hamady M, Fraser-liggett C, Knight R, Gordon JI: **The human**  
751 **microbiome project: exploring the microbial part of ourselves in a changing world.**  
752 *Nature* 2007, **449**:804–810.
- 753 2. Rittmann BE: **Microbial ecology to manage processes in environmental biotechnology.**  
754 *Trends Biotechnol* 2006, **24**:261–266.
- 755 3. Stewart EJ: **Growing unculturable bacteria.** *J Bacteriol* 2012, **194**:4151–4160.
- 756 4. Narayanasamy S, Muller EEL, Sheik AR, Wilmes P: **Integrated omics for the**  
757 **identification of key functionalities in biological wastewater treatment microbial**  
758 **communities.** *Microb Biotechnol* 2015, **8**:363–368.
- 759 5. Segata N, Waldron L, Ballarini A, Narasimhan V, Jousson O, Huttenhower C, Boernigen  
760 D, Tickle TL, Morgan XC, Garrett WS, Huttenhower C: **Computational meta’omics for**  
761 **microbial community studies.** *Mol Syst Biol* 2013, **9**:666.
- 762 6. Muller EEL, Glaab E, May P, Vlassis N, Wilmes P: **Condensing the omics fog of**  
763 **microbial communities.** *Trends Microbiol* 2013, **21**:325–333.
- 764 7. Roume H, Muller EEL, Cordes T, Renaut J, Hiller K, Wilmes P: **A biomolecular isolation**  
765 **framework for eco-systems biology.** *ISME J* 2013, **7**:110–121.
- 766 8. Roume H, Heintz-Buschart A, Muller EEL, Wilmes P: **Sequential isolation of**  
767 **metabolites, RNA, DNA, and proteins from the same unique sample.** *Methods Enzymol*  
768 2013, **531**:219–236.
- 769 9. Solomon K V, Haitjema CH, Thompson DA, O’Malley MA: **Extracting data from the**  
770 **muck: deriving biological insight from complex microbial communities and non-model**  
771 **organisms with next generation sequencing.** *Curr Opin Biotechnol* 2014, **28C**:103–110.

- 772 10. Sunagawa S, Mende DR, Zeller G, Izquierdo-Carrasco F, Berger SA, Kultima JR, Coelho  
773 LP, Arumugam M, Tap J, Nielsen HB, Rasmussen S, Brunak S, Pedersen O, Guarner F, de  
774 Vos WM, Wang J, Li J, Doré J, Ehrlich SD, Stamatakis A, Bork P: **Metagenomic species**  
775 **profiling using universal phylogenetic marker genes.** *Nat Methods* 2013, **10**:1196–1199.
- 776 11. Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, Arumugam M, Kultima JR, Prifti E,  
777 Nielsen T, Juncker AS, Manichanh C, Chen B, Zhang W, Levenez F, Wang J, Xu X, Xiao L,  
778 Liang S, Zhang D, Zhang Z, Chen W, Zhao H, Al-Aama JY, Edris S, Yang H, Wang J,  
779 Hansen T, Nielsen HB, Brunak S, et al.: **An integrated catalog of reference genes in the**  
780 **human gut microbiome.** *Nat Biotechnol* 2014, **32**:834–841.
- 781 12. Treangen TJ, Koren S, Sommer DD, Liu B, Astrovska I, Ondov B, Darling AE,  
782 Phillippy AM, Pop M: **MetAMOS: a modular and open source metagenomic assembly**  
783 **and analysis pipeline.** *Genome Biol* 2013, **14**:R2.
- 784 13. Nalbantoglu OU, Way SF, Hinrichs SH, Sayood K: **RAIphy: phylogenetic classification**  
785 **of metagenomics samples using iterative refinement of relative abundance index**  
786 **profiles.** *BMC Bioinformatics* 2011, **12**:41.
- 787 14. Celaj A, Markle J, Danska J, Parkinson J: **Comparison of assembly algorithms for**  
788 **improving rate of metatranscriptomic functional annotation.** *Microbiome* 2014, **2**:39.
- 789 15. Peng Y, Leung HCM, Yiu SM, Chin FYL: **IDBA-UD: a de novo assembler for single-**  
790 **cell and metagenomic sequencing data with highly uneven depth.** *Bioinformatics* 2012,  
791 **28**:1420–1428.
- 792 16. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W: **MEGAHIT: an ultra-fast single-node**  
793 **solution for large and complex metagenomics assembly via succinct de Bruijn graph.**  
794 *Bioinformatics* 2015, **31**:1674–1676.
- 795 17. Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, Delmont TO:

- 796 **Anvi'o: an advanced analysis and visualization platform for 'omics data.** *PeerJ* 2015,  
797 **3:e1319.**
- 798 18. Leimena MM, Ramiro-Garcia J, Davids M, van den Bogert B, Smidt H, Smid EJ,  
799 Boekhorst J, Zoetendal EG, Schaap PJ, Kleerebezem M: **A comprehensive**  
800 **metatranscriptome analysis pipeline and its validation using human small intestine**  
801 **microbiota datasets.** *BMC Genomics* 2013, **14**:530.
- 802 19. Leung HCM, Yiu S-M, Parkinson J, Chin FYL: **IDBA-MT: de novo assembler for**  
803 **metatranscriptomic data generated from next-generation sequencing technology.** *J*  
804 *Comput Biol* 2013, **20**:540–550.
- 805 20. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,  
806 Raychowdhury R, Zeng Q, Chen Z, Mauceli E, Hacohen N, Gnirke A, Rhind N, Palma F,  
807 Birren BW, Nusbaum C, Lindblad-toh K, Friedman N, Regev A: **Full-length transcriptome**  
808 **assembly from RNA-Seq data without a reference genome.** *Nat Biotechnol* 2011, **29**:644–  
809 652.
- 810 21. Kultima JR, Sunagawa S, Li J, Chen W, Chen H, Mende DR, Arumugam M, Pan Q, Liu  
811 B, Qin J, Wang J, Bork P: **MOCAT: a metagenomics assembly and gene prediction**  
812 **toolkit.** *PLoS One* 2012, **7**:e47656.
- 813 22. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y, Tang J, Wu  
814 G, Zhang H, Shi Y, Liu Y, Yu C, Wang B, Lu Y, Han C, Cheung DW, Yiu S-M, Peng S,  
815 Xiaoqian Z, Liu G, Liao X, Li Y, Yang H, Wang J, Lam T-W, Wang J: **SOAPdenovo2: an**  
816 **empirically improved memory-efficient short-read de novo assembler.** *Gigascience* 2012,  
817 **1**:18.
- 818 23. Lai B, Wang F, Wang X, Duan L, Zhu H: **InteMAP: Integrated metagenomic assembly**  
819 **pipeline for NGS short reads.** *BMC Bioinformatics* 2015, **16**:244.

- 820 24. Laczny CC, Pinel N, Vlassis N, Wilmes P: **Alignment-free visualization of**  
821 **metagenomic data by nonlinear dimension reduction.** *Sci Rep* 2014, **4**:4516.
- 822 25. Albertsen M, Hugenholtz P, Skarszewski A, Nielsen KL, Tyson GW, Nielsen PH:  
823 **Genome sequences of rare, uncultured bacteria obtained by differential coverage**  
824 **binning of multiple metagenomes.** *Nat Biotechnol* 2013, **31**:533–538.
- 825 26. Nielsen HB, Almeida M, Juncker AS, Rasmussen S, Li J, Sunagawa S, Plichta DR,  
826 Gautier L, Pedersen AG, Le Chatelier E, Pelletier E, Bonde I, Nielsen T, Manichanh C,  
827 Arumugam M, Batto J-M, Quintanilha Dos Santos MB, Blom N, Borruel N, Burgdorf KS,  
828 Boumezbear F, Casellas F, Doré J, Dworzynski P, Guarner F, Hansen T, Hildebrand F, Kaas  
829 RS, Kennedy S, Kristiansen K, et al.: **Identification and assembly of genomes and genetic**  
830 **elements in complex metagenomic samples without using reference genomes.** *Nat*  
831 *Biotechnol* 2014, **32**:822–828.
- 832 27. Alneberg J, Bjarnason BS, de Bruijn I, Schirmer M, Quick J, Ijaz UZ, Lahti L, Loman  
833 NJ, Andersson AF, Quince C: **Binning metagenomic contigs by coverage and**  
834 **composition.** *Nat Methods* 2014, **11**:1144–1146.
- 835 28. Kang DD, Froula J, Egan R, Wang Z: **MetaBAT, an efficient tool for accurately**  
836 **reconstructing single genomes from complex microbial communities.** *PeerJ* 2015,  
837 **3**:e1165.
- 838 29. Franzosa E a, Morgan XC, Segata N, Waldron L, Reyes J, Earl AM, Giannoukos G,  
839 Boylan MR, Ciulla D, Gevers D, Izard J, Garrett WS, Chan AT, Huttenhower C: **Relating**  
840 **the metatranscriptome and metagenome of the human gut.** *Proc Natl Acad Sci U S A*  
841 2014, **111**:E2329–E2338.
- 842 30. Hultman J, Waldrop MP, Mackelprang R, David MM, Mcfarland J, Blazewicz SJ, Harden  
843 J, Turetsky MR, Mcguire AD, Shah MB, Verberkmoes NC, Lee LH: **Multi-omics of**

- 844 **permafrost, active layer and thermokarst bog soil microbiomes.** *Nature* 2015, **521**:208–  
845 2112.
- 846 31. Bremges A, Maus I, Belmann P, Eikmeyer F, Winkler A, Albersmeier A, Pühler A,  
847 Schlüter A, Sczyrba A: **Deeply sequenced metagenome and metatranscriptome of a**  
848 **biogas-producing microbial community from an agricultural production-scale biogas**  
849 **plant.** *Gigascience* 2015, **4**:33.
- 850 32. Muller EEL, Pinel N, Laczny CC, Hoopman MR, Narayanasamy S, Lebrun LA, Roume  
851 H, Lin J, May P, Hicks ND, Heintz-Buschart A, Wampach L, Liu CM, Price LB, Gillece JD,  
852 Guignard C, Schupp JM, Vlassis N, Baliga NS, Moritz RL, Keim PS, Wilmes P:  
853 **Community integrated omics links the dominance of a microbial generalist to fine-tuned**  
854 **resource usage.** *Nat Commun* 2014, **5**:5603.
- 855 33. Roume H, Heintz-Buschart A, Muller EEL, May P, Satagopam VP, Laczny CC,  
856 Narayanasamy S, Lebrun LA, Hoopmann MR, Schupp JM, Gillece JD, Hicks ND,  
857 Engelthaler DM, Sauter T, Keim PS, Moritz RL, Wilmes P: **Comparative integrated omics:**  
858 **identification of key functionalities in microbial community-wide metabolic networks.**  
859 *npj Biofilms Microbiomes* 2015, **1**:15007.
- 860 34. Vanwonterghem I, Jensen PD, Ho DP, Batstone DJ, Tyson GW: **Linking microbial**  
861 **community structure, interactions and function in anaerobic digesters using new**  
862 **molecular techniques.** *Curr Opin Biotechnol* 2014, **27**:55–64.
- 863 35. Kenall A, Edmunds S, Goodman L, Bal L, Flintoft L, Shanahan DR, Shipley T: **Better**  
864 **reporting for better research: a checklist for reproducibility.** *BMC Neurosci* 2015, **16**:44.
- 865 36. Belmann P, Dröge J, Bremges A, McHardy AC, Sczyrba A, Barton MD: **Bioboxes:**  
866 **standardised containers for interchangeable bioinformatics software.** *Gigascience* 2015,  
867 **4**:47.

- 868 37. Di Tommaso P, Palumbo E, Chatzou M, Prieto P, Heuer ML, Notredame C: **The impact**  
869 **of Docker containers on the performance of genomic pipelines.** *PeerJ* 2015, **3**:e1273.
- 870 38. Köster J, Rahmann S: **Snakemake-a scalable bioinformatics workflow engine.**  
871 *Bioinformatics* 2012, **28**:2520–2522.
- 872 39. Koster J: **Reproducibility in next-generation sequencing analysis.**  
873 <http://dx.doi.org/10.17877/DE290R-7242> (2014). Accessed 05 Feb 2015.
- 874 40. Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina**  
875 **sequence data.** *Bioinformatics* 2014, **30**:2114–2120.
- 876 41. Mikheenko A, Saveliev V, Gurevich A: **MetaQUAST: evaluation of metagenome**  
877 **assemblies.** *Bioinformatics* 2015:btv697.
- 878 42. Seemann T: **Prokka: rapid prokaryotic genome annotation.** *Bioinformatics* 2014,  
879 **30**:2068–2069.
- 880 43. Laczny CC, Sternal T, Plugaru V, Gawron P, Atashpendar A, Margossian HH, Coronado  
881 S, der Maaten L van, Vlassis N, Wilmes P: **VizBin - an application for reference-**  
882 **independent visualization and human-augmented binning of metagenomic data.**  
883 *Microbiome* 2015, **3**:1.
- 884 44. Mende DR, Waller AS, Sunagawa S, Järvelin AI, Chan MM, Arumugam M, Raes J, Bork  
885 P: **Assessment of metagenomic assembly using simulated next generation sequencing**  
886 **data.** *PLoS One* 2012, **7**:e31386.
- 887 45. Schürch AC, Schipper D, Bijl MA, Dau J, Beckmen KB, Schapendonk CME, Raj VS,  
888 Osterhaus ADME, Haagmans BL, Tryland M, Smits SL: **Metagenomic survey for viruses**  
889 **in Western Arctic caribou, Alaska, through iterative assembly of taxonomic units.** *PLoS*  
890 *One* 2014, **9**:e105227.

- 891 46. Reyes A, Blanton L V., Cao S, Zhao G, Manary M, Trehan I, Smith MI, Wang D, Virgin  
892 HW, Rohwer F, Gordon JI: **Gut DNA viromes of Malawian twins discordant for severe**  
893 **acute malnutrition.** *Proc Natl Acad Sci U S A* 2015, **112**:11941–11946.
- 894 47. Li H, Durbin R: **Fast and accurate short read alignment with Burrows-Wheeler**  
895 **transform.** *Bioinformatics* 2009, **25**:589–595.
- 896 48. Quinlan AR, Hall IM: **BEDTools: a flexible suite of utilities for comparing genomic**  
897 **features.** *Bioinformatics* 2010, **26**:841–842.
- 898 49. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,  
899 Durbin R: **The Sequence Alignment/Map format and SAMtools.** *Bioinformatics* 2009,  
900 **25**:2078–2079.
- 901 50. Garrison E, Marth G: **Haplotype-based variant detection from short-read sequencing.**  
902 *bioRxiv* 2012:9.
- 903 51. Rimmer A, Phan H, Mathieson I, Iqbal Z, Twigg SRF, Wilkie AOM, McVean G, Lunter  
904 G: **Integrating mapping-, assembly- and haplotype-based approaches for calling**  
905 **variants in clinical sequencing applications.** *Nat Genet* 2014, **46**:912–918.
- 906 52. Patel RK, Jain M: **NGS QC Toolkit: a toolkit for quality control of next generation**  
907 **sequencing data.** *PLoS One* 2012, **7**:e30619.
- 908 53. Kanehisa M: **KEGG: Kyoto Encyclopedia of Genes and Genomes.** *Nucleic Acids Res*  
909 2000, **28**:27–30.
- 910 54. Ondov BD, Bergman NH, Phillippy AM: **Interactive metagenomic visualization in a**  
911 **Web browser.** *BMC Bioinformatics* 2011, **12**:385.
- 912 55. Jia B, Xuan L, Cai K, Hu Z, Ma L, Wei C: **NeSSM: a Next-generation Sequencing**  
913 **Simulator for Metagenomics.** *PLoS One* 2013, **8**:e75448.

914 56. Huang X, Madan A: **CAP3: A DNA sequence assembly program**. *Genome Res* 1999,  
915 **9**:868–877.

916 57. Fu L, Niu B, Zhu Z, Wu S, Li W: **CD-HIT: accelerated for clustering the next-**  
917 **generation sequencing data**. *Bioinformatics* 2012, **28**:3150–3152.

918 58. Zhu W, Lomsadze A, Borodovsky M: **Ab initio gene identification in metagenomic**  
919 **sequences**. *Nucleic Acids Res* 2010, **38**:e132.

920 59. Varrette S, Bouvry P, Cartiaux H, Georgatos F: **Management of an Academic HPC**  
921 **Cluster : The UL Experience**. *Proc 2014 Int Conf High Perform Comput Simul* 2014:959–  
922 967.

923

## 924 **Tables**

925 **Table 1.** Co-assemblies versus separate single-omic assemblies.

Sample	Assembly	No. of contigs (all)	Total length	N50	No. of predicted genes (unique)	Genome fraction (%)	MG mapped reads (%)	MT mapped reads (%)
<b>SM</b>	IMP	84052	198407791	10154	201480	65.3	97.6	90.2
	IMP-MEGAHIT	108011	208588327	8243	212669	70	98	91.39
	MG-only	74498	183252937	10902	187382	60.8	96.39	80.36
	MT-only	14723	9497250	961	7357	3.5	8.58	29.03
<b>HF</b>	IMP	131728	174193511	4227	182924	NA	87.5	96.8
	IMP-MEGAHIT	149163	181987482	3696	191503	NA	88	95.7
	MG-only	121018	146897961	3564	155068	NA	35.1	34.35
	MT-only	43466	46274315	2248	48527	NA	4.27	4.73
<b>WW</b>	IMP	167295	125102684	1501	126748	NA	30.9	61.6
	IMP-MEGAHIT	208345	143114883	1285	143354	NA	32.1	62.8
	MG-only	88818	77058077	2100	81084	NA	0.73	12.2
	MT-only	47237	23251573	808	17851	NA	0.05	0.53

926

927 Characteristics of the metagenomic (MG) and metatranscriptomic (MT) co-assemblies (IMP  
928 and IMP-MEGAHIT) against MG-only and MT-only assemblies based on simulated mock  
929 (SM) community dataset, human fecal (HF) community dataset and wastewater (WW) sludge  
930 community dataset. N50 statistics are reported based on a 500bp cut-off while read mappings  
931 were performed on all contigs.

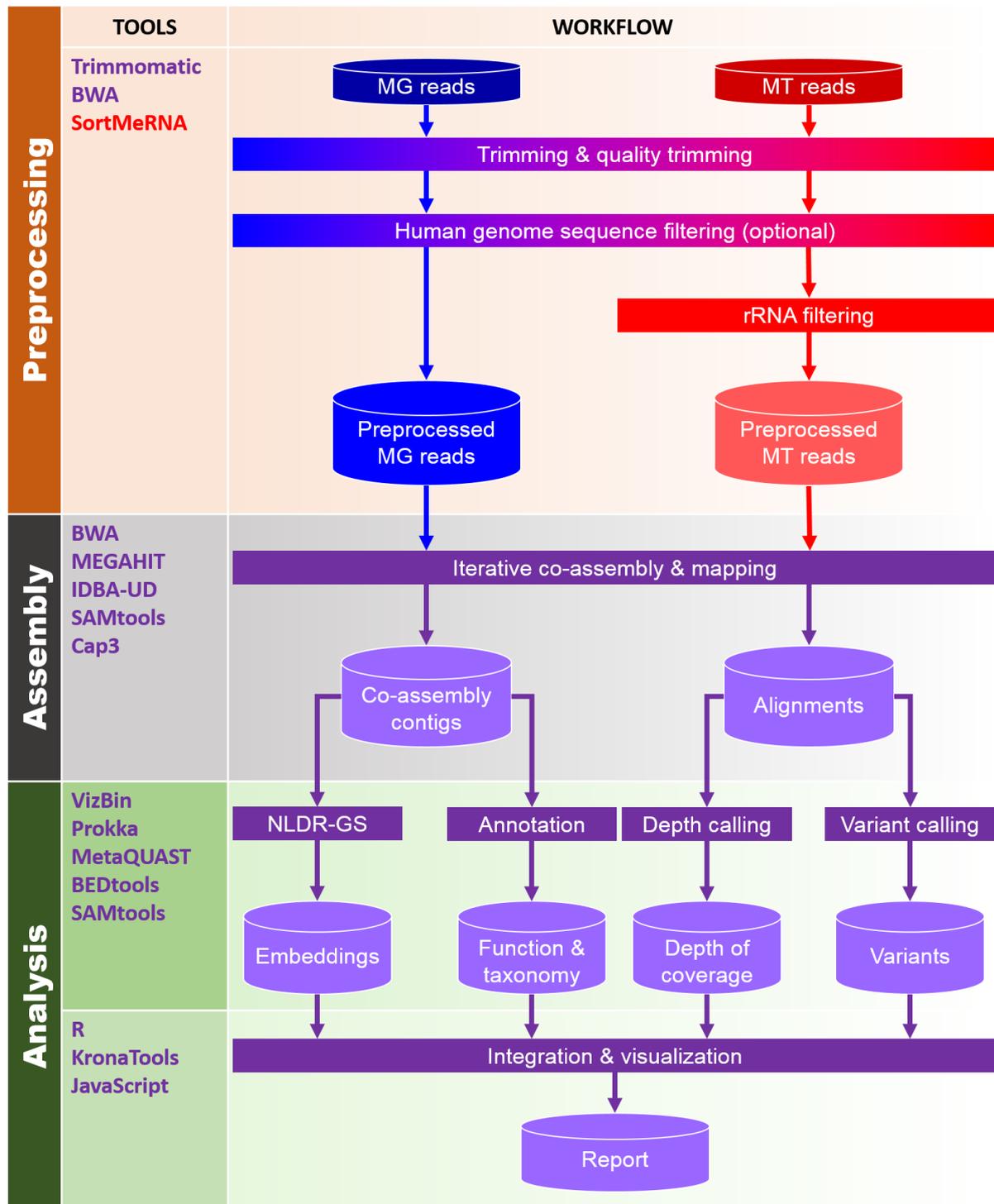
932 **Table 2.** Contigs reconstructed from the metatranscriptomic data with a likely  
 933 viral/bacteriophage origin/function.

Sample	Contig ID	Contig length	Avg. contig depth of coverage	Gene product	Avg. gene depth of coverage
HF	contig_34	6468	20926.50	Virus coat protein (TMV like)	30668.30
				Viral movement protein (MP)	26043.40
				RNA dependent RNA polymerase	22578.40
				Viral methyltransferase	18816.60
	contig_13948	2074	46.00	RNA dependent RNA polymerase	40.98
				Viral movement protein (MP)	55.76
WW	contig_6405	4062	46.24	Tombusvirus p33	43.42
				Viral RNA dependent RNA polymerase	41.76
				Viral coat protein (S domain)	36.48
	contig_7409	3217	20.62	Viral RNA dependent RNA polymerase	18.24
				Viral coat protein (S domain)	20.88
	contig_7872	2955	77.01	hypothetical protein	112.11
Phage maturation protein				103.02	

934

935 Contigs of  $\geq 1$ kb and average depth of coverage  $\geq 20$  were selected.

936 **Figures**



937

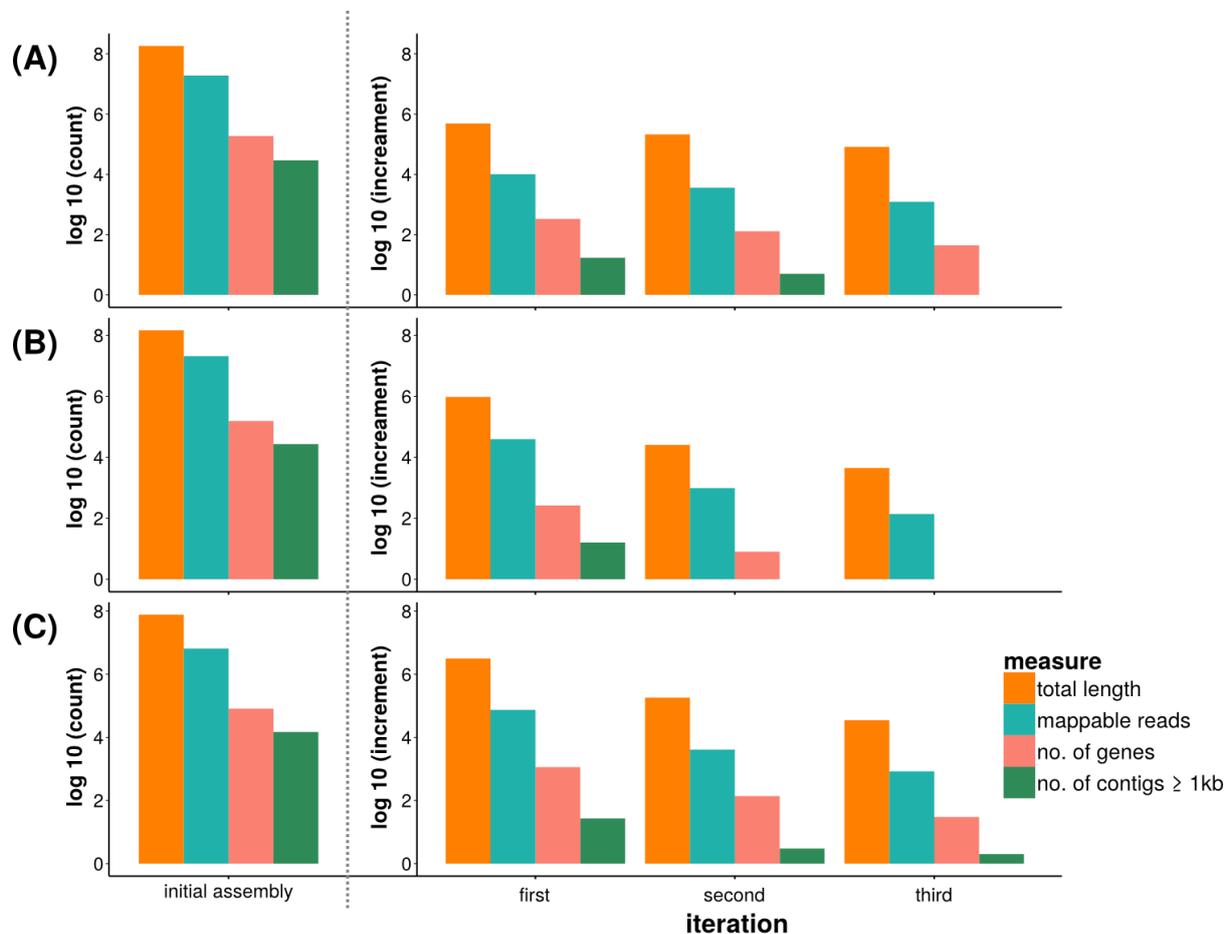


938 **Fig. 1.** Schematic overview of the IMP pipeline. Cylinders represent input and output while

939 rectangles represent processes. MG: Metagenomic data, MT: Metatranscriptomic data, rRNA:

940 ribosomal RNA and NLDR-GS: non-linear dimensionality reduction of genomic signatures.

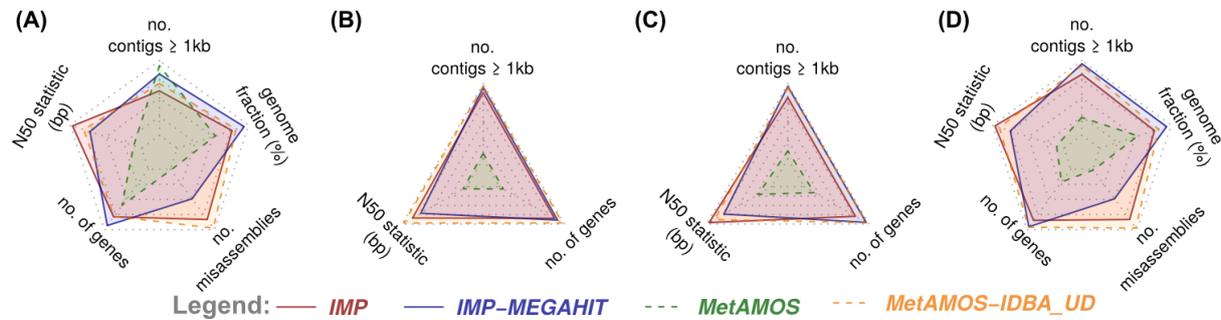
941 IMP handles raw paired-end (unprocessed) MG and MT data. Processes, input and output  
942 specific to MG and MT data are labeled blue and red, respectively. Processes and output that  
943 involve integration of MG and MT data are represented in purple. Details about the “iterative  
944 co-assembly” are available in [Additional file 2: Figure S1](#). IMP is launched as Docker container  
945 with Ubuntu as the operating system and uses Snakemake for workflow automation.



946

947 **Fig. 2.** Information contained within the iterative metagenomic assemblies. Quantitative  
948 assessment of the initial metagenomic (MG) assembly as well as incremental information and  
949 data usage from additional MG assembly iterations, employing unmappable reads from: (A)  
950 simulated mock (SM) community, (B) human fecal (HF) community and (C) wastewater  
951 (WW) sludge community.

952



953

954 **Fig. 3.** Assessment of the IMP-based iterative co-assemblies in comparison to established

955 methods. Radar charts summarizing the characteristics of the co-assemblies generated using

956 IMP and MetAMOS pipelines on: (A) simulated mock (SM) community (B) human fecal (HF)

957 community and (C) wastewater (WW) sludge community. (D) Summary radar chart reflecting

958 the cumulative measures obtained using the different datasets. The solid lines represent IMP

959 assemblies while the dashed lines represent MetAMOS assemblies, both pipelines executed

960 using two different assemblers. The assemblies are assessed based on number of contigs  $\geq$  1kb,

961 N50 statistics (contiguity), number of predicted genes (unique). N50 statistics are reported

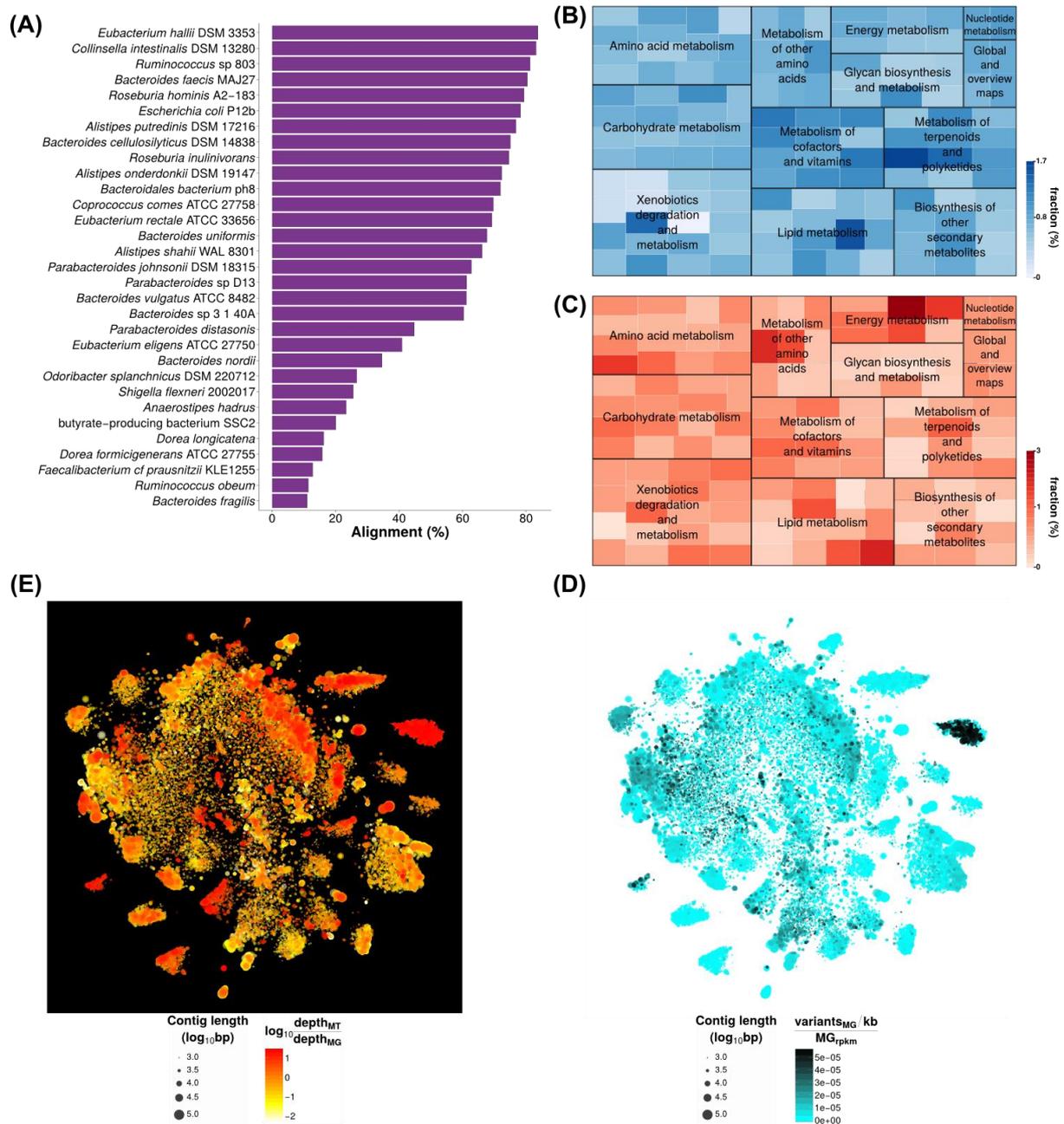
962 using a 500bp cut-off. Additional reference-based assessments for SM assemblies include

963 recovered genome fraction (%) and proportion of misassemblies. Higher values within the

964 radar charts (furthest from center) represent best performance, except for misassemblies, where

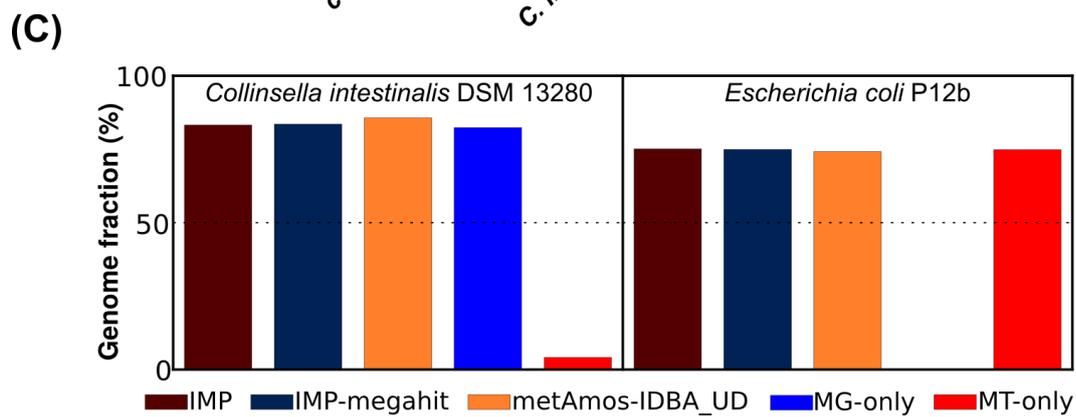
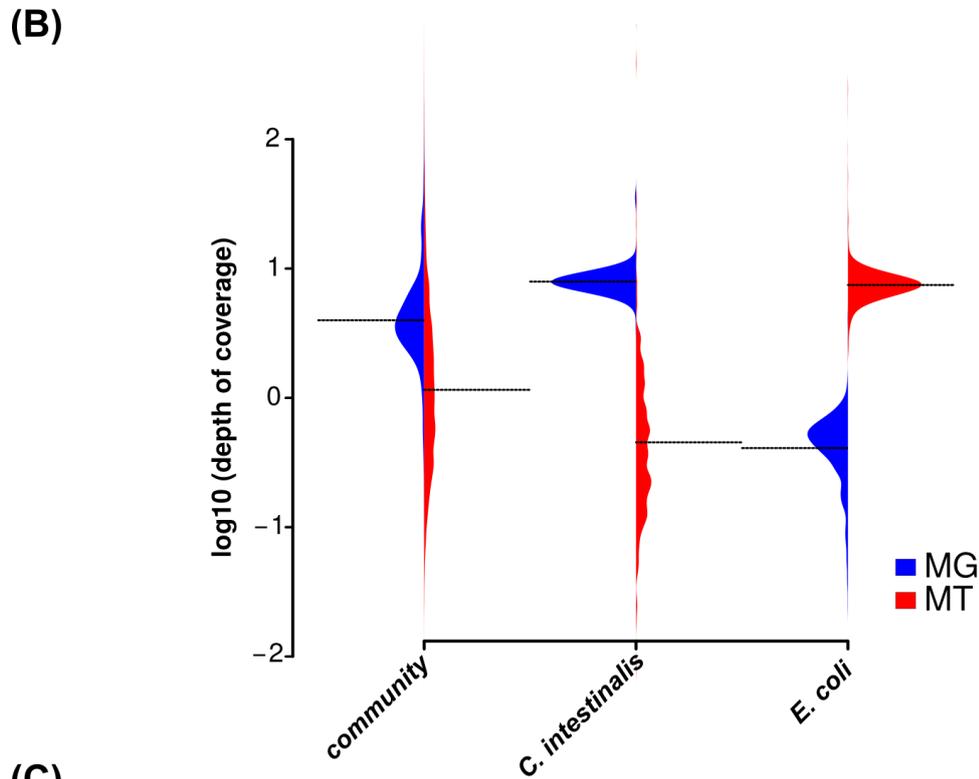
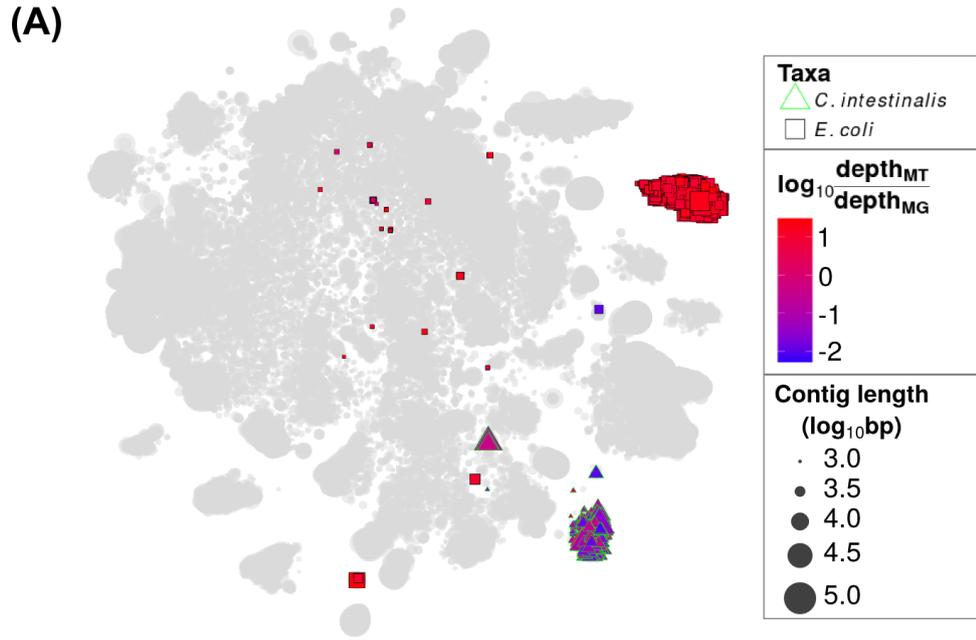
965 lower (closer to center) values indicate best performance.

966



967

968 **Fig. 4.** Examples of information retrievable from the IMP output of human fecal metagenomic  
 969 and metatranscriptomic data. **(A)** Taxonomic composition reflecting the percentages of  
 970 genomes covered. Representation of the inferred abundances of certain metabolic pathways  
 971 based on **(B)** metagenomic data (functional potential) and **(C)** metatranscriptomic data (gene  
 972 expression). Augmented VizBin maps of contigs  $\geq 1$  kb, representing **(D)** contig-level  
 973 metagenomic variant densities and **(E)** contig-level ratios of MT to MG average depth of  
 974 coverage.



976 **Fig. 5.** Metagenomic and metatranscriptomic data integration. **(A)** Augmented VizBin map  
977 highlighting contig subsets with sequences that are most similar to *Escherichia coli* P12b and  
978 *Colinsella intestinalis* DSM 13280 genomes. **(B)** Beanplot of metagenomic (MG) and  
979 metatranscriptomic (MT) average contig-level depth of coverage for the entire microbial  
980 community and two subsets (population-level genomes) of interest. **(C)** Recovered portion of  
981 subsets associated to the aforementioned taxa compared to a MG-only and MT-only  
982 assemblies.  
983

984

## 985 **Additional files**

986 The following additional data are available with the online version of this article.

987

### 988 **Additional file 1:**

989 File format: HTML

990 Title of data: Supplementary IMP HTML reports.

991 Description of data: HTML S1 and S2 are reports produced by IMP for the analysis of the  
992 human fecal (HF) microbial community and wastewater (WW) sludge microbial community  
993 datasets.

994

### 995 **Additional file 2:**

996 File format: PDF

997 Title of data: Supplementary figures.

998 Description of data: Figures S1 to S4. Detailed figure legends available within file.

999

### 1000 **Additional file 3:**

1001 File format: MS Excel (XLSX)

1002 Title of data: Supplementary tables.

1003 Description of data: Tables S1 to S9. Detailed table legends available within file.