

1 **Abundance-based reconstitution of microbial pan-genomes**  
2 **from whole-metagenome shotgun sequencing data**  
3

4 Florian Plaza Oñate<sup>1,2</sup>, Alessandra C. L. Cervino<sup>1</sup>, Frédéric Magoulès<sup>3</sup>, S. Dusko Ehrlich<sup>2</sup> &  
5 Matthieu Pichaud<sup>1§</sup>  
6

7 <sup>1</sup> Enterome, 94-96 Avenue Ledru Rollin, 75011 Paris, France

8 <sup>2</sup> MGP MetaGénoPolis, INRA, Université Paris-Saclay, 78350 Jouy en Josas, France

9 <sup>3</sup> CentraleSupélec, Université Paris-Saclay, 92290 Châtenay-Malabry, France  
10

11 <sup>§</sup> Corresponding author  
12

13 Email addresses:

14 FPO: `fplaza-onate [at] enterome [dot] com`

15 MP: `matthieu.pichaud [at] gmail [dot] com`

## 16 **Abstract**

17 Analysis toolkits for whole-metagenome shotgun sequencing data achieved strain-level  
18 characterization of complex microbial communities by capturing intra-species gene content  
19 variation. Yet, these tools are hampered by the extent of reference genomes that are far from  
20 covering all microbial variability, as many species are still not sequenced or have only few  
21 strains available.

22

23 Binning co-abundant genes obtained from de novo assembly is a powerful reference-free  
24 technique for discovering and reconstituting gene repertoire of microbial species. While  
25 current methods accurately identify species core genes, they miss many accessory genes or  
26 split them in small separated clusters.

27

28 We introduce MSPminer, a computationally efficient software tool that reconstitutes  
29 Metagenomic Species Pan-genomes (MSPs) by binning co-abundant genes across large-scale  
30 metagenomic datasets. MSPminer relies on a new robust measure for grouping not only  
31 species core genes but accessory genes also. In MSPs, an empirical classifier distinguishes core  
32 from accessory and shared genes.

33

34 We applied MSPminer to the largest publicly available gene abundance table which is  
35 composed of 9.9M genes quantified in 1 267 stool samples. We show that MSPminer  
36 successfully reconstitutes in a matter of several hours gene repertoire of > 1600 microbial  
37 species (some hitherto unknown) and detects many more accessory genes than existing tools.  
38 By compiling the information from thousands of samples, species gene content variability is  
39 better accounted for and their quantification is subsequently more precise

## 40 Introduction

41 Metagenomics has revolutionized microbiology by allowing culture-independent  
42 characterization of microbial communities. Its advent has allowed an unprecedented genetic  
43 characterization of the human gut microbiota and emphasized its fundamental role in health  
44 and disease [1–7]. Shotgun metagenomics where whole-community DNA is randomly  
45 sequenced bypasses the biases and limitations of 16S rRNA sequencing [8,9] by providing high  
46 resolution taxonomic profiling as well as insights into the diverse physiological roles and the  
47 metabolic potential of the community [10,11].

48  
49 The analysis of large cohorts revealed a substantial inter-individual microbial gene content  
50 variability [12] and nucleotide polymorphism [13] which reflects that individuals are not only  
51 carriers of various species, but also of different strains of the same species [14,15]. The  
52 characterization of the accessory genes found in individual strains is crucial in many contexts  
53 as they can provide functional advantages such as complex carbohydrates metabolism [16],  
54 antibiotic resistance or pathogenicity [17,18].

55  
56 Recent analysis toolkits for shotgun metagenomics data achieved strain-level resolution when  
57 coverage is sufficient. To this end, they either capture intra-species single-nucleotide  
58 polymorphisms (SNPs) in pre-identified marker genes [19,20], gene content variation [21] or  
59 both [22]. However, these tools are hampered by the extent of the reference genomes.

60  
61 Indeed, the microbial variability extends far beyond the content of reference genomes making  
62 metagenomic samples an untapped reservoir of information. First, it has been estimated that  
63 on average 50% of the species present in the human gut microbiota of western individuals  
64 lack reference genome and this proportion rises to 85% in individuals with traditional lifestyles  
65 [22]. Even if recent advancements of culture-based methods have proven that a substantial  
66 proportion of these species are actually culturable [23,24], the number of unknown species is  
67 probably still important. In addition, these techniques remain laborious and time consuming.  
68 Second, although species of public health interest (e.g. *Escherichia coli*, *Salmonella enterica* or  
69 *Clostridium difficile*) are represented by hundreds or even thousands of strains in genome  
70 databases [25], only few strains are available for the great majority of commensal species.

71 Consequently, accessory genes associated with microbial phenotypic traits may be missing in  
72 gene repertoires constructed from reference genomes.

73

74 Metagenomic assembly where overlapping reads are merged into longer sequences called  
75 contigs, is a powerful reference-free technique for overcoming the limitations of reference-  
76 based methods. Indeed, it allows the genetic survey of non-sequenced species and strains  
77 with previously unknown accessory genes. However, assembly remains a computationally  
78 challenging task [26] and despite the many dedicated tools proposed [27–30] the process only  
79 recovers incomplete genomes scattered in multiple contigs. To retrieve an exhaustive set of  
80 references, metagenomic assembly is performed on multiple samples. Then, non-redundant  
81 reference gene catalogs [31] are created and used as proxy for disease-related analyses [3,6]  
82 or descriptive purposes [12,32].

83

84 Metagenome-wide association studies have successfully identified associations between the  
85 gut microbiome and disease based on universal marker genes [33], clade specific marker  
86 genes [34] or whole gene sets [35]. The latter is the only method that allows an extensive  
87 investigation of microorganism core, accessory genes and mobile elements. However, testing  
88 millions of genes is biased and lacks statistical power. Testing all the variables is not adapted  
89 to situations where some of them are highly correlated [36], which is expected for genes  
90 coming from the same biological entity. The resulting list of significant genes will be biased  
91 towards organisms with the most genes in the pool as they have more chances of being picked  
92 up.

93

94 Considering that the genes originating from the same biological entity should have  
95 proportional abundances across samples, binning co-abundant genes has been proposed.  
96 However, clustering millions of genes is a very computationally intensive task and pairwise  
97 comparison of all gene abundance profiles is not feasible. To reduce the number of  
98 comparisons, some authors have performed binning on the subset of genes that were  
99 statistically significant [3,4,6] which restricts the analysis to the genes that are significant by  
100 themselves and does not improve the statistical power of the analysis. Others have proposed  
101 methods to perform the clustering of complete gene references based either on the Markov  
102 cluster algorithm [37] or a variant of the Canopy clustering algorithm [38].

103

104 Although direct proportionality is expected between co-abundant genes, these methods rely  
105 either on Pearson's or Spearman's correlation coefficients which respectively assess a linear  
106 association with a potentially non-null intercept or any monotonic association. Thus, these  
107 coefficients are too loose and spurious associations can be discovered. In addition, they are  
108 biased by sparse genes with many null counts [39], non-normal gene counts distributions [40]  
109 and presence of outliers [41].

110

111 Current clustering strategies group species core genes and highly prevalent accessory genes  
112 into the same cluster, but miss accessory genes of medium and low prevalence or assign them  
113 to small separate clusters [42]. Dependency between clusters of essential genes and accessory  
114 clusters can be evaluated downstream using the Fisher's exact test [38], which compares their  
115 presence/absence patterns across samples. Yet, this strategy does not account for the co-  
116 abundance of genes and is poorly discriminative when considering accessory clusters that are  
117 rare or associated with very prevalent species. In addition, it is not suitable for detecting  
118 clusters shared between several species.

119

120 To overcome these limitations, we have developed MSPminer, the first tool that discovers,  
121 delineates and structures Metagenomic Species Pan-genomes (MSPs) from large-scale  
122 shotgun metagenomics datasets without referring to genomes from isolated strains.  
123 MSPminer presents several significant improvements over existing methods. First, it relies on  
124 a new robust measure of proportionality for detection of co-abundant but not necessarily co-  
125 occurring genes as expected for non-core genes. Second, genes grouped in a MSP are  
126 empirically classified as core, accessory and shared genes.

127

128 To illustrate its usefulness, we applied MSPminer to the largest publicly available gene  
129 abundance table which is composed of 9.9M genes quantified in 1 267 samples [12]. We show  
130 that MSPminer successfully groups genes from the same species and identifies additional  
131 genes. Gene variability of microbial species is better captured and their quantification is  
132 subsequently more precise. MSPminer is a computationally efficient multithreaded program  
133 implemented in C++ that can process large datasets with millions of genes and thousands of  
134 samples in just a few hours on a single node server.

## 135 **Results**

### 136 **New measures of proportionality**

137 The gene repertoire of microbial species is composed of core genes present in all strains and  
138 accessory genes present in only some of them [43]. In a shotgun metagenomic sequencing  
139 context, we assumed that core genes of a microbial species should have a directly proportional  
140 number of mapped reads across samples (co-abundance) and should be consistently observed  
141 in samples if sequencing depth allows (co-occurrence). Remarkably, core genes and an  
142 accessory gene should have directly proportional counts only in the subset of samples where  
143 they are both detected (Figure 1). To group the core genes of a species and then identify its  
144 accessory genes, we developed a measure that evaluates proportionality between gene  
145 counts using samples where the number of mapped reads is high enough (see Methods).

146

147 To evaluate this new measure of proportionality, we generated an abundance table that  
148 simulates the counts of genes from a single virtual species across 300 samples (see Methods).  
149 We considered that each sample was carrying a unique strain of the species with specific gene  
150 content. Genes present in all the samples were labeled as core and those detected in a subset  
151 as accessory. We used this dataset to compare the performance of the Pearson correlation  
152 coefficient, the Spearman correlation coefficient and the proposed measure of proportionality  
153 for detecting a relation between the abundance profile of the species core genome and all its  
154 genes including accessories (Figure 2). Pearson and Spearman correlation coefficients  
155 decrease all the more as the prevalence of a tested gene decreases while the proposed  
156 measure remains high, as only samples where both the species core and its accessory gene  
157 are detected are used for calculation. Therefore, the association between core genes and  
158 many accessory genes will be missed using the correlation coefficients. However, accessory  
159 genes observed in similar subsets of samples may be grouped into small distinct clusters as  
160 their abundance profiles should be correlated.

161

162 Finally, we derived a robust version of the measure to identify associated genes despite the  
163 presence of samples with inconsistent counts named hereafter outliers. We evaluated this  
164 robust measure of proportionality against the non-robust version described above by adding  
165 an increasing percentage of outliers to the genes abundance profiles. For a given percentage

166 of outliers, each of these genes was compared to the outlier-free abundance profile of the  
167 species core genome. This simulation showed that the non-robust measure of proportionality  
168 decreases all the more as the percentage of outliers increases whereas the robust measure  
169 remains high; demonstrating that proportionality is still detected (Supplementary Figure 1).  
170 However, the robust measure decreases significantly when the percentage of outliers is high  
171 and the gene prevalence is low.

## 172 **Reconstitution of Metagenomic Species Pan-genomes of the human gut microbiota**

173 We developed MSPminer, a program that uses measures of proportionality to group co-  
174 abundant genes into Metagenomic Species Pan-genomes (MSPs, Figure 3). MSPminer  
175 empirically distinguishes core from accessory genes based on their presence absence patterns  
176 (see Methods) and tags genes observed in samples where the core is not detected as shared  
177 (Figure 4). Finally, non-core genes observed in the same subset of samples are grouped into  
178 modules of co-occurring genes.

179

180 We applied MSPminer to the largest publicly available gene abundance table provided with  
181 the Integrated Gene Catalog of the human gut microbiome [12]. In this table, 9 879 896 genes  
182 are quantified across 1 267 stool samples from individuals of various geographical origin  
183 (Europe, USA and China) and diverse health status (healthy, obese, diabetic, with  
184 inflammatory bowel disease etc.). 6 971 229 (70.6%) genes with counts greater than 6 in at  
185 least 3 samples were kept. Among these, 3 262 914 (46.8%) were organized into 1 677  
186 Metagenomic Species Pan-Genomes (MSPs) with at least 200 core genes (Supplementary  
187 Table 1).

## 188 *Taxonomy*

189 By considering the lowest taxonomic rank assigned to MSPs, 278 (16.6%) were annotated at  
190 species level, 85 (5%) at genus level, 119 (7.1%) at phylum level and the remaining 1 213 MSPs  
191 (71.5%) could not be annotated (Supplementary Figure 2) indicating that a clear majority of  
192 MSPs correspond to species not represented in reference genomes databases. Only 4 MSPs  
193 were annotated as Eukaryotes and corresponded to intestinal parasites of the *Blastocystis*  
194 genus. All others MSPs were assigned to prokaryotic species of which only 3 were Archaea and  
195 corresponded to the species *Methanobrevibacter smithii*, *Candidatus Methanomethylophilus*

196 *alvus* and *Methanosphaera stadtmanae*. Among the MSPs annotated as Bacteria, the phyla  
197 *Firmicutes* (268 MSPs), *Bacteroidetes* (118 MSPs), *Proteobacteria* (50 MSPs) and  
198 *Actinobacteria* (29 MSPs) were the most represented (Supplementary Figure 3) as expected  
199 from human gut metagenomes [44]. Some species were represented by multiple MSPs such  
200 as *Faecalibacterium prausnitzii* (5 MSPs), *Bacteroides fragilis* (2 MSPs), *Methanobrevibacter*  
201 *smithii* (2 MSPs) or *Hungatella hathewayi* (2 MSPs) suggesting a high nucleotide and gene  
202 content variability between strains attributed so far to the same species. Conversely, some  
203 MSPs had their core genes attributed to multiple species, usually from the same genus  
204 (Supplementary Table 3). Although some inconsistencies in the taxonomy assignment have  
205 previously been reported [22,45], it is possible that several MSPs regroup genes from highly  
206 related species with high average nucleotide identity and similar gene content. MSPs  
207 annotated at species level had a consistent taxonomical annotation as 97% of the core genes  
208 (median) were assigned to the same species (Supplementary Figure 4). Consistency of the  
209 taxonomic annotation was lower for accessory genes (62%) caused mainly by unannotated  
210 genes.

#### 211 *MSP content*

212 Most MSPs were small (median number of genes = 1 784) even if 53 had more than 5 000  
213 genes (Supplementary Figure 5). As expected, a strong positive correlation (Pearson's  $r = 0.8$ )  
214 between the total number of genes in a MSP and its number of accessory genes was observed  
215 (Supplementary Figure 6). Interestingly, four outliers corresponding to the intestinal parasites  
216 previously described had a high number of core genes and few accessory genes. This suggests  
217 that Eukaryotic genomes have a larger number of genes and a lower gene content variability  
218 than Prokaryotes. Among the MSPs with the more accessory genes (Supplementary Table 2),  
219 many corresponded to species reported as highly variable such as *Escherichia coli* [46],  
220 *Klebsiella pneumoniae* [47] or *Clostridium bolteae* [48]. As previously observed in population  
221 genomics studies comparing multiple strains of the same species [49,50], the prevalence of  
222 accessory genes in MSPs often follows a bimodal distribution (Supplementary Figure 7)  
223 showing either a high or low prevalence but rarely intermediate. Thus, the number of  
224 accessory genes in a MSP is correlated (Spearman's  $\rho = 0.86$ ) with its prevalence  
225 (Supplementary Figure 8). Indeed, the more a MSP is detected in many samples, the more  
226 exhaustively MSPminer will recover its accessory genes, especially the rare ones. Many MSPs



227 annotated at species level had accessory genes previously unobserved in available genomes  
228 (Supplementary Table 2)

### 229 *Prevalence*

230 As for the genes in the catalog, most MSPs were detected in very few samples (Supplementary  
231 Figure 9). Only 40 MSPs were detected in at least 70% of the samples showing that the  
232 common microbial core of the human gut microbiota is limited to a few dozen species  
233 (Supplementary Table 2). No clear relation between the prevalence of the MSPs and their  
234 mean abundance was found (Supplementary Figure 10). However, 2 MSPs corresponding to  
235 *Bacteroides vulgatus* and *Bacteroides uniformis* were both very prevalent (detected in 97.6%  
236 and 94.6% of the samples respectively) and very abundant (mean relative abundance of 7.8%  
237 and 4.4% respectively). Interestingly, many rare MSPs were abundant in the few samples  
238 which carried them. Many of these MSPs were annotated as bacteria of the *Lactobacillus*  
239 genus most likely consumed as probiotics. However, some others correspond to known  
240 invasive species associated with severe dysbiosis such as *Fusobacterium nucleatum* [51] or  
241 *Clostridium clostridioforme* [52].

### 242 *Census of universal single copy marker genes*

243 To check that MSPs correspond to real microbial species and evaluate the completeness of  
244 their set of core genes, we identified in each of them 40 universal single copy marker genes  
245 (SCM) [45]. 878 MSPs (54%) had at least 30 SCM and 403 (24%) had all of them (Supplementary  
246 Figure 11 A and Supplementary Table 2). As housekeeping genes, SCMs are essential to the  
247 microbe survival and should be found among core genes. Indeed, 92% of the SCMs were core  
248 genes in their respective MSP and the rest was mainly high prevalent accessory genes  
249 (Supplementary Figure 11 B). This shows the classification of genes as core or accessory  
250 performed by MSPminer is reliable.

### 251 *Comparison to sequenced genomes*

252 We compared the MSPs to 642 sequenced genomes for which at least 10% of their constituent  
253 genes were detected in the Integrated Genes Catalog of the human gut microbiome [12]  
254 (Supplementary Table 4). In total, these genomes covered 398 species representing 114  
255 different genera. 624 (97.1%) were unambiguously assigned to 281 different MSPs. By keeping  
256 only one representative per species, 47.1% (resp. 60.4%) of the genomes had at least 75% of

257 their genes grouped in their corresponding MSP considering either all their genes or only those  
258 that were in the catalog (Supplementary Figure 12). In compliance with the results of the  
259 taxonomic analysis, highly related species were assigned to the same MSP such as *Escherichia*  
260 *coli*, *Escherichia fergusonii* and all genomes of the *Shigella* genus. Conversely, some species  
261 grouping highly divergent strains were represented by several MSPs including  
262 *Faecalibacterium prausnitzii* or *Bacteroides fragilis*.

263

264 To give another perspective on the MSPs, we compared the complete genome of  
265 *Parabacteroides distasonis* ATCC 8503 [53] to its corresponding MSP (Figure 7). Among the 3  
266 850 genes predicted in the genome, 3 781 (98%) had at least a close homolog in the Integrated  
267 Gene Catalogue and 3 442 (89%) were found in the msp\_0011. As expected, almost all the  
268 core genes from the MSP were found in the genome (1 867 / 1 921, 97%), as well as accessory  
269 genes with a prevalence higher than 80% (522 / 599, 87%). Only a small fraction of less  
270 prevalent accessory genes was found in the genome (1371 / 5 090, 27%) (Supplementary  
271 Figure 13). Genes grouped in the same modules tended to be physically close, in coherence  
272 with genome organization of prokaryotes (Supplementary Figure 14). Remarkably, some  
273 singleton genes were surrounded by genes from the same module, which shows that the  
274 stringent grouping criteria used by MSPminer may split genes with slightly different  
275 presence/absence patterns into different modules while they could have been grouped  
276 (Supplementary Figure 15). Finally, few genomic regions contained genes that were not  
277 assigned to the MSP (Supplementary Table 5A). Interestingly, some of these regions were  
278 annotated as mobile elements (Supplementary Table 5B). Although some could be false  
279 negatives, many were appropriately excluded as they were observed in too few samples or  
280 their counts did not meet the proportionality-based grouping criterion (Supplementary Figure  
281 16).

## 282 **Comparison to the Canopy clustering algorithm**

283 The Canopy clustering algorithm [38] was compared to MSPminer by applying both tools to  
284 the metagenomic dataset described above. In total, Canopy grouped 2 691 408 genes into 3  
285 463 Co-Abundance gene Groups (CAGs) while MSPminer grouped 3 267 132 genes (+17.6%)  
286 into 1 677 MSPs (~ two-fold less objects).

287

288 178 MSPs encompassing 154 617 genes had no equivalent among the CAGs. Most MSPs were  
289 rare as 75% were detected in fewer than 5 samples but had a significant size as 50% were  
290 composed of at least 700 genes (Supplementary Figure 17). Remarkably, for 75% of them, the  
291 3 samples with the highest counts represented at least 90% of the sum of MSP abundance on  
292 all the samples. By default, Canopy discards such cases to avoid detection of spurious  
293 correlations but MSPminer limits this risk by applying a variance-stabilizing transformation  
294 and a stringent association criterion. In addition, Canopy grouped most core genes of a MSP  
295 into a single CAG while many accessory genes were missed or assigned to small separate CAGs  
296 (Figure 5). In agreement with the results of the simulation, most of the missed accessory genes  
297 had a medium or low prevalence. As they contained many unexpected zeros, the correlations  
298 with the core of their respective species were below the limit set in Canopy (Figure 6).

### 299 **MSPs quantification for biomarkers discovery**

300 To demonstrate that MSPminer was useful for biomarkers discovery, we first looked for  
301 differentially abundant MSPs according to the geographical origin of samples. We discovered  
302 94 MSPs differentially abundant between Westerners and Chinese ( $q$ -value  $< 10^{-3}$ ,  $\log_2$  fold  
303 change  $\geq 1$ ) including 72 more abundant in Westerners and 22 in Chinese (Supplementary  
304 Figure 18 and Supplementary Table 6A). Among the discriminant MSPs, all those assigned to  
305 the *Proteobacteria* phylum (*Klebsiella pneumoniae*, *Escherichia coli* and *Bifidobacterium*  
306 *wadsworthia*) were more abundant in Chinese which is consistent with previously published  
307 results [12]. Interestingly, two MSPs assigned to *Faecalibacterium prausnitzii* were significant  
308 but one was more abundant in Westerners and the other in Chinese. This shows that some  
309 strains of this species are associated with geographical origin of samples. In addition, we  
310 discovered 75 MSPs differentially abundant between Europeans and Americans ( $q$ -value  $<$   
311  $10^{-3}$ ,  $\log_2$  fold change  $\geq 1$ ) of which 70 were more abundant among Europeans.  
312 (Supplementary Figure 19 and Supplementary Table 6B). This result is consistent with previous  
313 studies showing lower gut microbiota diversity among Americans compared to Europeans  
314 [33].

315

316 Secondly, we checked if MSPs could be used to perform strain-level analysis. To do this, we  
317 tested if some accessory genes in the MSPs were more prevalent in samples of a given  
318 geographical origin. By way of example, we found 680 accessory genes associated with

319 geographical origin (chi-squared test,  $p$ -value  $< 10^{-10}$ ) in the msp\_0011 corresponding to  
320 *Parabacteroides distasonis* (Supplementary Figure 20 and Supplementary Table 7).  
321 Remarkably, genes involved in cell filamentation (V1.UC58-4\_GL0042624 and V1.UC18-  
322 0\_GL0014340) were more prevalent in Chinese than Westerners. More generally, many  
323 significant genes were more prevalent among Westerners or Chinese although European and  
324 American samples were separated in the analysis. This result suggests that strains of  
325 *Parabacteroides distasonis* carried by Chinese are distinct from those of Westerners.

## 326 **Discussion**

### 327 **Direct proportionality hypothesis and limits**

328 MSPminer relies on a new robust measure to detect genes with directly proportional counts.  
329 Even if this relation is more stringent than those assessed by Pearson's or Spearman's  
330 correlation coefficients, it was successfully used to reconstitute Metagenomic Species Pan-  
331 genomes of the human gut microbiota. In fact, most genes from sequenced genomes were  
332 grouped into a single MSP showing that direct proportionality is the most common relation  
333 between genes from the same biological entity. MSPminer misses some genes for which  
334 counts are not ruled by this relation. Indeed, proportionality is disrupted when gene copy  
335 number varies across samples [14] (Supplementary Figure 21), when a sample contains  
336 multiples strains [19,20] and when a gene is shared between several MSPs (Supplementary  
337 Figure 22) because of horizontal gene transfer [54] or grouping of highly similar orthologs.  
338 Nevertheless, the first two cases have most likely a limited impact as the majority of strains  
339 tend to have the same gene copy numbers [14] and samples often carry a dominant strain  
340 [20]. Regarding shared genes, their signals are a linear combination of the MSPs that carry  
341 them. Thus, they will be identified only if these MSPs are mostly detected in separate sets of  
342 samples.

### 343 **Computing performance**

344 MSPminer can process large datasets made up of thousands of samples and millions of genes  
345 in just a few hours on a regular single node server. The program is only limited by the amount  
346 of RAM available on the machine on which it is executed as the input count matrix must be  
347 fully loaded into memory. Consequently, RAM consumption grows linearly according to the  
348 number of samples and genes in the count matrix (Supplementary Figure 23).

349 MSPminer achieves good parallel efficiency (Supplementary Figure 24) through two  
350 parallelization strategies. First, a novel Map/Reduce programming model assigns genes to as  
351 many subsets as the number of available samples. Genes with greatest counts in the same  
352 sample are first compared, which not only decreases the number of comparisons to perform  
353 but increases the probability that related genes are placed in the same bin compared to  
354 random assignment (Supplementary Figure 25). Each subset of genes is processed in parallel  
355 and synchronization is only required before the reduction/merging step (see Methods). To  
356 avoid comparisons of all pairs of genes, others used an iterative algorithm where a random  
357 gene is compared to all the others until a significant proportion of genes is clustered [38]. This  
358 method allows fast identification of big clusters but struggles with the large number of small  
359 clusters and singleton genes. Furthermore, parallelism implementation is more complex as  
360 synchronization is required to detect duplicate clusters. In a second phase, MSPminer  
361 performs pairwise comparison of clusters to detect those corresponding to MSPs core gene  
362 sets. Then, full MSPs including accessory and shared genes are retrieved in parallel from the  
363 signal of their respective core genes. Here, no synchronization is required as core genes sets  
364 are supposedly independent.

### 365 **Quality of MSPs**

366 The quality of the MSPs is impacted by all the upstream steps required for generating the  
367 count matrix, as well as with the biological and ecological characteristics of the dataset. At the  
368 sequencing level, the number of reads (sequencing depth) generated for each sample impacts  
369 the detection and coverage of subdominant species, while reads length affects the quality of  
370 the assembly and the ability to assign a read to a gene without ambiguity. At the  
371 bioinformatics level, assembly, gene prediction, gene redundancy removal, mapping and  
372 counting require expertise to select the most appropriate strategies, tools and parameters.  
373 Indeed, assemblers returning chimeric contigs which combine sequences from highly related  
374 species, inaccurate predictors generating truncated or merged genes, redundancy removal  
375 with a common threshold for all genes (95% of nucleotide identity) lead to genes of variable  
376 quality in catalogues. Then, when quantifying genes, keeping only uniquely mapped reads  
377 underestimates the abundance of some genes whereas considering shared reads can generate  
378 false positives. Genes grouped in MSPs were significantly longer than those that were not  
379 (median length of 780 bp vs 498 bp, Wilcoxon rank-sum test p-value= 0) (Supplementary

380 Figure 26) as longer genes have a higher and less dispersed counts. Nevertheless, end-to-end  
381 mapping probably plays a role as it may fail to fully align a read whose size is approximately  
382 the same as the target gene. Finally, at the biology level, a high number of samples with varied  
383 phenotypes will improve the comprehensiveness and quality of MSPs. Indeed, as the number  
384 of samples grows, MSPminer will identify rare species and will extend the list of accessory  
385 genes of the MSPs corresponding to species with an open pan-genome. In addition, highly  
386 prevalent accessory genes will be reclassified from core to accessory as observed while  
387 sequencing an increasing number of strains of a species [50].

### 388 **Applications**

389 As illustrated in this paper, MSPminer supports the analysis of metagenomic data at species-  
390 level by identifying and quantifying the MSPs present in samples. Subsequently, MSPs  
391 associated with a given phenotype (e.g. the geographical origin) can be investigated both  
392 quantitatively and qualitatively. Here, information from unknown or non-sequenced species  
393 can be exploited. Compared to methods relying on marker genes [33,34], MSPminer improves  
394 the estimation of species abundance by automatically detecting among core genes those with  
395 the highest specificity, the highest counts and lowest dispersion.

396 Moreover, in each MSP, genes or modules of accessory genes associated with the tested  
397 phenotype can be explored opening the way to a strain-level analysis. Thus, biomarkers  
398 corresponding to functional traits specific to certain strains can be discovered.

399

400 MSPminer also provides microbial population genetics from large cohorts which can support  
401 culture-dependent methods by identifying species of particular interest, such as those with no  
402 reference genome available or with reference genomes distant from the strains actually  
403 present. Reciprocally, MSPminer will benefit from advances in culture-dependent methods  
404 which provide reference genomes of low abundance species detectable by shotgun  
405 sequencing but difficult to assemble [23,24].

### 406 **Further developments**

407 Several improvements of MSPminer are considered. Algorithms that could identify relevant  
408 associations currently missed by MSPminer will be evaluated. For instance, deconvolution

409 algorithms [14] may discover genes shared between several MSPs while kernel density  
410 estimators [55] may be useful for detecting genes with highly variable copy numbers.  
411 To increase MSPminer specificity, a statistical test determining how much an association  
412 between the core genome of a MSP and a gene is unexpected would be highly useful. Such a  
413 statistical test could assess both co-occurrence with a Fisher exact test and co-abundance with  
414 correlation test or equivalent. Currently, zero counts are either classified as structural or  
415 undetermined while comparing genes abundance profiles. However, a statistical model  
416 determining the probability that a zero is structural would allow overcoming threshold effects  
417 and classifying with more accuracy a gene as core, accessory or shared.  
418 Alternatives to the median for computing the representative of a MSP are envisaged. The sum  
419 which cumulates counts from multiple genes is a prime candidate, as clusters would be  
420 quantified with higher accuracy and increased quantification range particularly in samples  
421 where its abundance is low. However, one should carefully account for outliers.  
422 Finally, the robustness of the measure of proportionality could be improved. For instance, a  
423 robust linear regression [56] may replace the median for estimating the coefficient of  
424 proportionality between the abundance profiles of two genes while the median absolute  
425 deviation (MAD) could improve the detection of outliers [57].



## 426 **Methods**

### 427 **Measure of proportionality between two genes**

428 Let  $M$  be a  $n \times m$  matrix where  $n$  is the number of genes and  $m$  the number of samples.

429  $M$  is composed of counts  $c_{ij}$  representing the number of reads mapped on gene  $i$  in sample  $j$ .

430

431 Let  $g_i = (c_{i1}, c_{i2}, \dots, c_{im})$   $i \in [1, n]$  be the vector of the number of mapped reads on the gene  
432  $i$  across the  $m$  samples.

433

434 The distribution of count data  $g_i$  has the following properties:

- 435 1. Variance tends to be proportional to the counts.
- 436 2. It ranges over several orders of magnitude due to uneven sequencing depth and  
437 variable relative abundance of the gene between samples.
- 438 3. It usually contains many zeros as the majority of genes are observed in a few samples.  
439 For instance, the count matrix used in this study contains 92% of zeros.
- 440 4. It is prone to outliers. (see Discussion)

441

442 Let  $g_x$  and  $g_y$  ( $x, y \in [1, n]$  and  $x \neq y$ ) denote the vectors of the number of mapped reads  
443 on two distinct genes. A robust measure is proposed to assess direct proportionality between  
444  $g_y$  and  $g_x$  (formally written  $g_x \propto g_y$ ) which accounts for the points mentioned above.

445

446 Let  $\alpha$  be the coefficient of proportionality between  $g_x$  and  $g_y$ .  $\alpha$  is a strictly positive constant  
447 expected to be roughly equal to the ratio of  $g_x$  and  $g_y$  length. However, it can be impacted  
448 by other factors such as uneven coverage or gene duplication (Supplementary Figure 27).  
449 Therefore, instead of relying on genes length,  $\alpha$  was robustly estimated as follow:

$$450 \quad \sqrt{\alpha} = \text{median} \left( \frac{\sqrt{c_{ys}}}{\sqrt{c_{xs}}} \right) \forall s \in [1, m] \text{ such as } c_{ys} \geq t \text{ and } c_{xs} \geq t \text{ with } t = 6 \text{ by default}$$

451

452 A square root transformation was applied to stabilize variance as suggested by several authors  
453 for count data [58,59]. For a comparison of some data transformations, refer to  
454 Supplementary Figure 28. The median was used to tolerate some outliers.

455



456 To estimate the coefficient of proportionality, only samples where both genes counts were  
457 above a threshold  $t$  were kept. This has the following advantages:

- 458 1. It discards samples where both genes are absent as they do not provide any quantitative  
459 information for the estimation.
- 460 2. It discards samples with overdispersed low counts which do not allow a precise  
461 estimation of the coefficient of proportionality.
- 462 3. It discards samples where only one gene has a null count. In such sample, the zero  
463 count can be either a sampling zero that corresponds to an undetected gene because  
464 of sampling or technical effects or a structural zero that corresponds to unobserved  
465 gene actually absent in the sample. Distinguish structural from sampling zeros is crucial  
466 to classify a gene as core or accessory. Here, zeros below the threshold  $t$  were of an  
467 undetermined type (yellow points in Figure 4) whereas those above were classified as  
468 structural (red points in Figure 4).

469

470 When  $\alpha > 1$ ,  $g_y$  yields more counts than  $g_x$ . As a result, a null count from  $g_x$  can be  
471 misclassified as a structural zero.

472

473 When  $\alpha \neq 1$ , different quantification thresholds for  $g_x$  and  $g_y$  respectively named  $t_x$  and  $t_y$   
474 were used to reflect the different yields for  $g_x$  and  $g_y$ :

475 
$$\alpha \geq 1 \rightarrow (t_x = t \text{ and } t_y = \alpha \cdot t)$$

476 
$$\alpha < 1 \rightarrow (t_x = \frac{t}{\alpha} \text{ and } t_y = t)$$

477 *Non-robust version*

478 The relationship of proportionality between the two genes was evaluated by a modified  
479 version of the Lin's concordance correlation coefficient [60] by considering only samples where  
480 both genes had non-null counts.

481

482 The Lin's concordance correlation coefficient originally designed to assess relationships of the  
483 type  $y = x$  was defined as:

484 
$$\frac{2 \cdot cov(x, y)}{\sigma_x^2 + \sigma_y^2 + (\bar{x} - \bar{y})^2}$$

485 where  $\bar{x}$  and  $\bar{y}$  are the means,  $\sigma_x^2$  and  $\sigma_y^2$  are the variances and  $cov(x, y)$  is the covariance of  
486 the variables  $x$  and  $y$ .

487

488 To assess relationships of the type  $y = kx$  where  $k$  is a constant, the variable  $x$  was  
489 substituted by  $kx$  in the formula hereabove:

490 
$$\frac{2k \cdot cov(x, y)}{k^2 \cdot \sigma_x^2 + \sigma_y^2 + (k \cdot \bar{x} - \bar{y})^2}$$

491

492 After substituting the variables  $x$  and  $y$  by  $\sqrt{g_x}$  and  $\sqrt{g_y}$  respectively and  $k$  by  $\sqrt{\alpha}$ , the final  
493 formula was:

494 
$$\frac{2\sqrt{\alpha} \cdot cov(\sqrt{g_x}, \sqrt{g_y})}{\alpha \cdot \sigma_{\sqrt{g_x}}^2 + \sigma_{\sqrt{g_y}}^2 + (\sqrt{\alpha} \cdot \sqrt{g_x} - \sqrt{g_y})^2}$$

495 *Robust version*

496 This measure was used to assess the proportionality between two genes in presence of few  
497 outliers. If not treated specifically, even a small number of outliers may decrease significantly  
498 the concordance coefficient calculated.

499

500 Residuals were computed in samples where both genes have non-null counts with the  
501 following formula:

502 
$$abs(\sqrt{g_y} - \alpha \cdot \sqrt{g_x})$$

503

504 Let  $Q_1$  and  $Q_3$  be the first and third quartiles of the residuals. Let  $IQR$  be the interquartile  
505 range defined by  $IQR = Q_3 - Q_1$ . Samples with residuals above  $Q_3 + 1.5 \cdot IQR$  were  
506 considered as outliers (purple points in Figure 4).

507

508 Let  $m$  be the number of samples for which residuals were computed. If there were more than  
509  $(m - 5) \cdot 0.3$  outliers (percentage of outliers asymptotically equal to 30%), the robust  
510 measure of proportionality was not computed. Otherwise, the concordance coefficient was  
511 calculated on non-outlier samples.

512 *Simulation*

513 A simulated gene abundance table which quantifies genes of a single virtual species across  
514 300 virtual samples was generated.

515 The species pan-genome consisted of 2 000 core genes and 2 900 accessory genes of variable  
516 length. Genes length was generated from a negative binomial (mean = 1 000; overdispersion  
517 = 0.3) and values below 100 and above 5 000 were rejected.

518 Each sample was carrier of a strain with specific gene content. For each prevalence ranging  
519 from 10 to 299 samples, 10 accessory genes were considered. The subset of genes in which  
520 an accessory gene was detected was drawn randomly.

521 The number of mapped reads per sample was generated from a log normal law (mean=log(1  
522 000 000), sd=1) and values below 100 000 and above 10 000 000 were rejected.

523 The theoretical number of counts attributed to each gene was calculated according to its  
524 length and its presence or not in the strain. The observed gene counts were generated by  
525 using a negative binomial distribution (mean=theoretical gene count; overdispersion=0.05) to  
526 approach real metagenomic data.

527 Outliers were added in each gene by selecting 5%, 10% and 20% of their non-null samples and  
528 multiplying each observed count by either  $\frac{1}{4}$ ,  $\frac{1}{3}$ , 2, 3 or 4.

529 **Metagenomic Species Pan-genome generation**

530 The Supplementary Figure 29 gives an overview of the MSPminer workflow.

531 *Input data*

532 MSPminer processes raw gene counts tables. In this study, we used the table provided with  
533 the Integrated Gene Catalog of the human gut microbiome (1267sample.gene.pairNum.table)  
534 available on GigaDB [61].

535 *Data filtering*

536 Rare genes which do not support enough quantitative information for further processing were  
537 discarded. By default, genes with counts greater than 6 in at least 3 samples were kept.

538 *Data transformation*

539 A square root transformation was applied to gene counts. This transformation stabilizes gene  
540 counts variance and limits the skewness of gene counts distribution.

541 *Gene binning*

542 Genes with the highest counts in the same sample were binned. To limit bias due to variable  
543 sequencing depth, raw read counts were normalized by the number of mapped reads prior to  
544 bin assignment (Supplementary Figure 30). Note that normalized counts were used in this step  
545 only.

546 *Seeds creation*

547 This step identifies sets of co-abundant and co-occurring genes called *seeds* hereafter.  
548 Seeds were created in parallel in each bin by a greedy approach. First, genes were compared  
549 pairwise. All pairs of genes with a non-robust measure of proportionality of at least 0.8 and  
550 no structural zeros were saved in a list. Then, the list was sorted by decreasing measure of  
551 proportionality. The pair of genes with the highest measure of proportionality was selected as  
552 a centroid. Genes related to one of the centroid genes were grouped together in a new seed.

553 *Seed representative*

554 For each seed, a pseudo gene referred as *representative* was computed as follow. First, the  
555 seed representative was defined as the median vector of the counts of all its genes. Then, each  
556 gene of the seed was compared to the seed representative using the measure proportionality.  
557 The final seed representative corresponded to the median vector of the counts of the 30 genes  
558 with the highest measure of proportionality.

559 *Seeds merging*

560 Some related genes may have been assigned to different bins, for instance, in a situation  
561 where samples with the highest counts had close values. Therefore, a merging step was  
562 performed. First, seeds from all the bins were pooled and sorted by decreasing size. Then, the  
563 representative of the largest seed was compared to the representatives of the other seeds.  
564 Seeds with a non-robust measure of proportionality of at least 0.8 and no structural zeros  
565 counts were merged with the largest seed to form the final seed. Merged seeds were removed  
566 from the list and the procedure was iterated until or there were no more seeds to process.  
567 After merging, seeds with less than 200 genes were discarded.

568 *Core seeds identification*

569 In this step, core seeds were identified among final seeds, based on the assumption that in a  
570 set of related seeds, the largest corresponds to a species core genome and the others are  
571 modules of either accessory or shared genes.

572 First, seeds were sorted by decreasing number of genes. The largest seed was defined as a  
573 new core seed. Then, the representative of the core seed was compared to the representative  
574 of all remaining seeds. The seeds with a robust measure of proportionality of at least 0.8 were  
575 considered as related to the core seed and discarded from the list of potential cores. The  
576 procedure was iterated until there was no more seed to process.

577 *Metagenomic Species Pan-genome generation*

578 The representatives of each core seed were compared to all the genes. Because core seeds  
579 were identified all at once in the previous step, the MSPs generation was run in parallel. Genes  
580 with a robust measure of proportionality of at least 0.8 were considered as related to the core  
581 seed.

582

583 Let  $g_x$  be the median vector of the number of mapped reads on the core seed and  $g_y$  the  
584 vector of the number of mapped reads on a gene related to the core seed. The related gene  
585 was assigned to one of the 4 following categories:

586 1. core genes were detected in the same samples as the core seed (Figure 4A).

587  $(g_y \infty g_x)$  and  $(\forall s \in [1, m] c_{xs} \geq t_x \rightarrow c_{ys} \neq 0 \text{ and } c_{ys} \geq t_y \rightarrow c_{xs} \neq 0)$

588 2. accessory genes were detected in a subset of samples where the core seed was  
589 detected (Figure 4B).

590  $(g_y \infty g_x)$  and  $(\exists s \in [1, m] c_{xs} \geq t_x \text{ and } c_{ys} = 0)$  and  $(\forall s \in [1, m] c_{ys} \geq t_y \rightarrow c_{xs} \neq 0)$

591 3. shared core genes were detected in all the samples where the core seed was detected  
592 plus some samples where it was not (Figure 4C).

593  $(g_y \infty g_x)$  and  $(\forall s \in [1, m] c_{xs} \geq t_x \rightarrow c_{ys} \neq 0)$  and  $(\exists s \in [1, m] c_{ys} \geq t_y \text{ and } c_{xs} = 0)$

594 4. shared accessory genes were detected in a subset of samples where the core seed was  
595 detected plus some samples where it was not (Figure 4D).

596  $(g_y \infty g_x)$  and  $(\exists s \in [1, m] c_{xs} \geq t_x \text{ and } c_{ys} = 0)$  and  $(\exists s \in [1, m] c_{ys} \geq t_y \text{ and } c_{xs} = 0)$

597

598 In each category, a clustering procedure similar to the one used to create seeds was run. It  
599 identified modules of co-occurring genes that may be interpreted as functional units, i.e.  
600 operons. Unclustered genes were saved as singleton modules.

### 601 **Comparison to the Canopy clustering algorithm**

602 The implementation of the Canopy clustering algorithm was downloaded at  
603 [https://www.cbs.dtu.dk/projects/CAG/Supplementary\\_Software\\_canopy\\_clustering.zip](https://www.cbs.dtu.dk/projects/CAG/Supplementary_Software_canopy_clustering.zip). The  
604 gene count table normalized following the procedure described in [3] was taken as an input.  
605 Default parameters were used.  
606 MSPs were projected on Co-Abundance gene Groups (CAGs) and reciprocally with an in-house  
607 script.

### 608 **Biomarkers discovery**

#### 609 *Identification of MSPs associated with geographical origin*

610 A two-tailed Wilcoxon rank-sum test was used on relative median abundance of the 30 best  
611 representative core genes of each MSP (1 696 tested variables). The obtained p-values were  
612 adjusted by the Benjamini-Hochberg procedure. In addition, a  $\log_2$  ratio was computed  
613 between the median abundances of the MSP in the two populations tested. MSPs with an  
614 adjusted p-value inferior to  $10^{-2}$  and a  $\log_2$  ratio superior to 1 were considered significant.

#### 615 *Identification of accessory genes associated with geographical origin*

616 For each accessory gene of a MSP, a 2x2 contingency table counting in both populations the  
617 number of samples where the gene was present or absent was built. Only samples where the  
618 MSP core genome was detected were kept. A gene was considered as present in a sample if  
619 at least two reads were mapped on it. Then, a chi-squared test was performed on each  
620 contingency table. Accessory genes with a p-value inferior to  $10^{-10}$  were considered  
621 significant.  $\log_2$  presence ratios equal to +infinity or -infinity were replaced by +10 or -10  
622 respectively.

### 623 **Taxonomic annotation of the gene catalog**

624 Genes were aligned at the nucleotide level using BLASTn [62] (version 2.6.0) against KEGG  
625 GENOME [63] (Release 82.0, April 2017) and RefSeq [64] (Release 81, March 2017). Hits that  
626 covered less than 80% of the query gene or with a e-value superior to 0.01 were discarded.

627 Thresholds of 95%, 80% and 65% of nucleotide identity were respectively used for taxonomic  
628 annotation at species, genus and phylum level.

629 At a given taxonomic level, “no consensus” was reported if the selected hits did not share the  
630 same annotation. Finally, results from KEGG were preferred to those from RefSeq.

631

632 MSPs were assigned to the lowest level taxon representing more than 50% of the annotations  
633 of their core genes. Taxonomic annotations of MSPs with at least 50% of their core genes  
634 annotated at species level (including “no consensus”) and less than 80% assigned to the most  
635 represented species were considered ambiguous (c.f. Supplementary Table 3)

### 636 **Functional annotation of the gene catalog**

637 Translated genes were annotated with eggNOG-mapper [65] (version 0.12.7) based on eggNOG  
638 orthology assignments [66]. Sequence similarity searches were performed using HMMER [67].

639

640 The 40 universal single copy marker genes were discovered using fetchMG v1.0 [68]

### 641 **Comparison of the MSPs to sequenced genomes**

642 Genomes used to build the integrated catalog of the human gut microbiome [12], HMP  
643 reference genomes [69] and genomes from species detected while performing the taxonomic  
644 annotation of the MSPs were downloaded from GenBank [70]. When not provided, CDS were  
645 predicted with Prodigal [71]. Genes from the reference catalog were aligned against the  
646 genomes with BLASTn [62] (version 2.6.0; arguments: -perc\_identity 95 -ungapped).  
647 Alignments of less than 100 nucleotides were discarded. Hits found for a single gene at  
648 neighboring positions on the target genome were merged. 642 genomes with less than 10%  
649 of their constituent genes detected in the reference catalog were kept. Genes were annotated  
650 with the related MSP information when available. Hits from most abundant MSP were kept  
651 and overlapping hits from less abundant MSPs were discarded. The local GC-content of  
652 genome was computed using a sliding window of approximately 100 nucleotides. Finally these  
653 data were plotted using Circos [72].

## 654 **Declarations**

### 655 **Funding**

656 This work was funded by Enterome, the ANRT (Association Nationale de la Recherche et de la  
657 Technologie) via the grant CIFRE 2014/0057 and INRA MetaGenoPolis via the grant  
658 “Investissements d'avenir” ANR-11-DPBS-0001.

### 659 **Authors' contributions**

660 FPO and MP designed the software, performed the analyses and wrote the manuscript. FPO  
661 implemented the software. FM, AC and SDE supervised the project and revised the  
662 manuscript.

### 663 **Competing interests**

664 The authors declare that they have no competing interests.

### 665 **Additional files**

666 Supplementary Table 1: Tab-separated file listing the genes and modules in the MSPs  
667 Supplementary Table 2: XLS file describing the MSPs (taxonomic annotation, number of genes,  
668 number of universal marker genes, prevalence and abundance)  
669 Supplementary Table 3: XLS file listing the MSPs with ambiguous annotation at species level.  
670 Supplementary Table 4: XLS file summarizing the comparison of the MSPs to 642 sequenced  
671 genomes.  
672 Supplementary Table 5: XLS file listing the genomic regions of *Parabacteroides distasonis* ATCC  
673 8503 containing genes not assigned to the msp\_0011.  
674 Supplementary Table 6: XLS file listing the MSPs associated with the geographic origin of  
675 samples.  
676 Supplementary Table 7: XLS file listing the accessory genes of the msp\_0011 (*Parabacteroides*  
677 *distasonis*) associated with the geographic origin of samples.

678



## 679 **References**

- 680 1. O'Hara AM, Shanahan F. The gut flora as a forgotten organ. *EMBO Rep.* [Internet].  
681 2006;7:688–93. Available from:  
682 <http://embor.embopress.org/cgi/doi/10.1038/sj.embor.7400731>
- 683 2. Morgan XC, Tickle TL, Sokol H, Gevers D, Devaney KL, Ward D V, et al. Dysfunction of the  
684 intestinal microbiome in inflammatory bowel disease and treatment. *Genome Biol.* [Internet].  
685 2012;13:R79. Available from: [http://genomebiology.biomedcentral.com/articles/10.1186/gb-](http://genomebiology.biomedcentral.com/articles/10.1186/gb-2012-13-9-r79)  
686 [2012-13-9-r79](http://genomebiology.biomedcentral.com/articles/10.1186/gb-2012-13-9-r79)
- 687 3. Qin J, Li Y, Cai Z, Li S, Zhu J, Zhang F, et al. A metagenome-wide association study of gut  
688 microbiota in type 2 diabetes. *Nature* [Internet]. 2012;490:55–60. Available from:  
689 <http://www.nature.com/doifinder/10.1038/nature11450>
- 690 4. Le Chatelier E, Nielsen T, Qin J, Prifti E, Hildebrand F, Falony G, et al. Richness of human gut  
691 microbiome correlates with metabolic markers. *Nature* [Internet]. 2013;500:541–6. Available  
692 from: <http://www.nature.com/doifinder/10.1038/nature12506>
- 693 5. Kostic AD, Xavier RJ, Gevers D. The Microbiome in Inflammatory Bowel Disease: Current  
694 Status and the Future Ahead. *Gastroenterology* [Internet]. 2014;146:1489–99. Available from:  
695 <http://linkinghub.elsevier.com/retrieve/pii/S0016508514002200>
- 696 6. Qin N, Yang F, Li A, Prifti E, Chen Y, Shao L, et al. Alterations of the human gut microbiome  
697 in liver cirrhosis. *Nature* [Internet]. 2014;513:59–64. Available from:  
698 <http://www.nature.com/doifinder/10.1038/nature13568>
- 699 7. Feng Q, Liang S, Jia H, Stadlmayr A, Tang L, Lan Z, et al. Gut microbiome development along  
700 the colorectal adenoma–carcinoma sequence. *Nat. Commun.* [Internet]. 2015;6:6528.  
701 Available from: <http://www.nature.com/doifinder/10.1038/ncomms7528>
- 702 8. Větrovský T, Baldrian P. The Variability of the 16S rRNA Gene in Bacterial Genomes and Its  
703 Consequences for Bacterial Community Analyses. *PLoS One.* 2013;8.
- 704 9. Brooks JP, Edwards DJ, Harwich MD, Rivera MC, Fettweis JM, Serrano MG, et al. The truth

- 705 about metagenomics: quantifying and counteracting bias in 16S rRNA studies. *BMC Microbiol.*  
706 [Internet]. 2015;15:66. Available from: <http://www.biomedcentral.com/1471-2180/15/66>
- 707 10. Ranjan R, Rani A, Metwally A, McGee HS, Perkins DL. Analysis of the microbiome:  
708 Advantages of whole genome shotgun versus 16S amplicon sequencing. *Biochem. Biophys.*  
709 *Res. Commun.* [Internet]. Elsevier Ltd; 2016;469:967–77. Available from:  
710 <http://dx.doi.org/10.1016/j.bbrc.2015.12.083>
- 711 11. Jovel J, Patterson J, Wang W, Hotte N, O’Keefe S, Mitchel T, et al. Characterization of the  
712 Gut Microbiome Using 16S or Shotgun Metagenomics. *Front. Microbiol.* [Internet]. 2016;7.  
713 Available from: <http://journal.frontiersin.org/Article/10.3389/fmicb.2016.00459/abstract>
- 714 12. Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, et al. An integrated catalog of reference  
715 genes in the human gut microbiome. *Nat. Biotechnol.* [Internet]. 2014;32:834–41. Available  
716 from: <http://www.nature.com/doifinder/10.1038/nbt.2942>
- 717 13. Schloissnig S, Arumugam M, Sunagawa S, Mitreva M, Tap J, Zhu A, et al. Genomic variation  
718 landscape of the human gut microbiome. *Nature* [Internet]. Nature Publishing Group;  
719 2012;493:45–50. Available from: <http://www.nature.com/doifinder/10.1038/nature11711>
- 720 14. Greenblum S, Carr R, Borenstein E. Extensive Strain-Level Copy-Number Variation across  
721 Human Gut Microbiome Species. *Cell* [Internet]. 2015;160:583–94. Available from:  
722 <http://linkinghub.elsevier.com/retrieve/pii/S0092867415000136>
- 723 15. Zhu A, Sunagawa S, Mende DR, Bork P. Inter-individual differences in the gene content of  
724 human gut bacterial species. *Genome Biol.* [Internet]. 2015;16:82. Available from:  
725 <http://genomebiology.com/2015/16/1/82>
- 726 16. Larsbrink J, Rogers TE, Hemsworth GR, McKee LS, Tauzin AS, Spadiut O, et al. A discrete  
727 genetic locus confers xyloglucan metabolism in select human gut Bacteroidetes. *Nature*  
728 [Internet]. 2014;506:498–502. Available from:  
729 <http://www.nature.com/doifinder/10.1038/nature12907>
- 730 17. Loman NJ, Constantinidou C, Christner M, Rohde H, Chan JZ-M, Quick J, et al. A Culture-  
731 Independent Sequence-Based Metagenomics Approach to the Investigation of an Outbreak of

- 732 Shiga-Toxigenic *Escherichia coli* O104:H4. *JAMA* [Internet]. 2013;309:1502. Available from:  
733 <http://jama.jamanetwork.com/article.aspx?doi=10.1001/jama.2013.3231>
- 734 18. Scaria J, Ponnala L, Janvilisri T, Yan W, Mueller LA, Chang Y-F. Analysis of Ultra Low Genome  
735 Conservation in *Clostridium difficile*. Horsburgh MJ, editor. *PLoS One* [Internet].  
736 2010;5:e15147. Available from: <http://dx.plos.org/10.1371/journal.pone.0015147>
- 737 19. Luo C, Knight R, Siljander H, Knip M, Xavier RJ, Gevers D. ConStrains identifies microbial  
738 strains in metagenomic datasets. *Nat. Biotechnol.* [Internet]. 2015;33:1045–52. Available  
739 from: <http://www.nature.com/doi/10.1038/nbt.3319>
- 740 20. Truong DT, Tett A, Pasolli E, Huttenhower C, Segata N. Microbial strain-level population  
741 structure & genetic diversity from metagenomes. *Genome Res.* 2017;27:626–38.
- 742 21. Scholz M, Ward D V, Pasolli E, Tolio T, Zolfo M, Asnicar F, et al. Strain-level microbial  
743 epidemiology and population genomics from shotgun metagenomics. *Nat. Methods*  
744 [Internet]. 2016;13:435–8. Available from:  
745 <http://www.nature.com/doi/10.1038/nmeth.3802>
- 746 22. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. An integrated metagenomics pipeline  
747 for strain profiling reveals novel patterns of bacterial transmission and biogeography. *Genome*  
748 *Res.* [Internet]. 2016;26:1612–25. Available from:  
749 <http://genome.cshlp.org/lookup/doi/10.1101/gr.201863.115>
- 750 23. Browne HP, Forster SC, Anonye BO, Kumar N, Neville BA, Stares MD, et al. Culturing  
751 of ?unculturable? human microbiota reveals novel taxa and extensive sporulation. *Nature*  
752 [Internet]. 2016;533:543–6. Available from:  
753 <http://www.nature.com/doi/10.1038/nature17645>
- 754 24. Lagier J-C, Khelaifia S, Alou MT, Ndongo S, Dione N, Hugon P, et al. Culture of previously  
755 uncultured members of the human gut microbiota by culturomics. *Nat. Microbiol.* [Internet].  
756 2016;1:16203. Available from: <http://www.nature.com/articles/nmicrobiol2016203>
- 757 25. Wattam AR, Abraham D, Dalay O, Disz TL, Driscoll T, Gabbard JL, et al. PATRIC, the bacterial  
758 bioinformatics database and analysis resource. *Nucleic Acids Res.* 2014;42.

- 759 26. Pop M. Genome assembly reborn: Recent computational challenges. *Brief. Bioinform.*  
760 2009;10:354–66.
- 761 27. Boisvert S, Raymond F, Godzaridis É, Laviolette F, Corbeil J. Ray Meta: scalable de novo  
762 metagenome assembly and profiling. *Genome Biol.* [Internet]. 2012;13:R122. Available from:  
763 <http://genomebiology.biomedcentral.com/articles/10.1186/gb-2012-13-12-r122>
- 764 28. Peng Y, Leung HCM, Yiu SM, Chin FYL. IDBA-UD: A de novo assembler for single-cell and  
765 metagenomic sequencing data with highly uneven depth. *Bioinformatics.* 2012;28:1420–8.
- 766 29. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. MEGAHIT: an ultra-fast single-node solution for  
767 large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*  
768 [Internet]. 2015;31:1674–6. Available from:  
769 [https://academic.oup.com/bioinformatics/article-](https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btv033)  
770 [lookup/doi/10.1093/bioinformatics/btv033](https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btv033)
- 771 30. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. metaSPAdes: a new versatile  
772 metagenomic assembler. *Genome Res.* [Internet]. 2017; Available from:  
773 <http://genome.cshlp.org/lookup/doi/10.1101/gr.213959.116>
- 774 31. Sharpton TJ. An introduction to the analysis of shotgun metagenomic data. *Front. Plant*  
775 *Sci.* [Internet]. 2014 [cited 2017 Jan 9];5. Available from:  
776 <http://journal.frontiersin.org/article/10.3389/fpls.2014.00209/abstract>
- 777 32. Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, et al. A human gut microbial  
778 gene catalogue established by metagenomic sequencing. *Nature* [Internet]. 2010;464:59–65.  
779 Available from: <http://www.nature.com/doi/10.1038/nature08821>
- 780 33. Sunagawa S, Mende DR, Zeller G, Izquierdo-Carrasco F, Berger SA, Kultima JR, et al.  
781 Metagenomic species profiling using universal phylogenetic marker genes. *Nat. Methods*  
782 [Internet]. 2013;10:1196–9. Available from:  
783 <http://www.nature.com/doi/10.1038/nmeth.2693>
- 784 34. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, et al. MetaPhlan2 for  
785 enhanced metagenomic taxonomic profiling. *Nat. Methods* [Internet]. 2015;12:902–3.

- 786 Available from: <http://www.nature.com/doifinder/10.1038/nmeth.3589>
- 787 35. Wang J, Jia H. Metagenome-wide association studies: fine-mining the microbiome. *Nat.*  
788 *Rev. Microbiol.* [Internet]. 2016;14:508–22. Available from:  
789 <http://www.nature.com/doifinder/10.1038/nrmicro.2016.83>
- 790 36. Schwartzman A, Lin X. The effect of correlation in false discovery rate estimation.  
791 *Biometrika* [Internet]. 2011;98:199–214. Available from:  
792 <https://academic.oup.com/biomet/article-lookup/doi/10.1093/biomet/asq075>
- 793 37. Karlsson FH, Tremaroli V, Nookaew I, Bergström G, Behre CJ, Fagerberg B, et al. Gut  
794 metagenome in European women with normal, impaired and diabetic glucose control. *Nature*  
795 [Internet]. 2013;498:99–103. Available from:  
796 <http://www.nature.com/doifinder/10.1038/nature12198>
- 797 38. Nielsen HB, Almeida M, Juncker AS, Rasmussen S, Li J, Sunagawa S, et al. Identification and  
798 assembly of genomes and genetic elements in complex metagenomic samples without using  
799 reference genomes. *Nat. Biotechnol.* [Internet]. 2014 [cited 2014 Jul 9];32:822–8. Available  
800 from: <http://www.ncbi.nlm.nih.gov/pubmed/24997787>
- 801 39. Huson LW. Performance of Some Correlation Coefficients When Applied to Zero-Clustered  
802 Data. *J. Mod. Appl. Stat. Methods* [Internet]. 2007;6:530–6. Available from:  
803 <http://digitalcommons.wayne.edu/jmasm%5Cnhttp://digitalcommons.wayne.edu/jmasm/vol6/iss2/17>  
804 16/iss2/17
- 805 40. Kowalski CJ. On the Effects of Non-Normality on the Distribution of the Sample Product-  
806 Moment Correlation Coefficient. *Appl. Stat.* [Internet]. 1972;21:1. Available from:  
807 <http://www.jstor.org/stable/10.2307/2346598?origin=crossref>
- 808 41. Osborne JW, Overbay A. The power of outliers (and why researchers should always check  
809 for them). *Pract. Assessment, Res. Eval.* 2004;9:1–8.
- 810 42. Almeida M, Pop M, Le Chatelier E, Prifti E, Pons N, Ghoulane A, et al. Capturing the most  
811 wanted taxa through cross-sample correlations. *ISME J.* [Internet]. 2016;10:2459–67.  
812 Available from: <http://www.nature.com/doifinder/10.1038/ismej.2016.35>

- 813 43. Medini D, Donati C, Tettelin H, Massignani V, Rappuoli R. The microbial pan-genome. Curr.  
814 Opin. Genet. Dev. 2005. p. 589–94.
- 815 44. Arumugam M, Raes J, Pelletier E, Le Paslier D, Yamada T, Mende DR, et al. Enterotypes of  
816 the human gut microbiome. Nature [Internet]. 2011;474:666–666. Available from:  
817 <http://www.nature.com/doifinder/10.1038/nature10187>
- 818 45. Mende DR, Sunagawa S, Zeller G, Bork P. Accurate and universal delineation of prokaryotic  
819 species. Nat. Methods [Internet]. 2013;10:881–4. Available from:  
820 <http://www.nature.com/doifinder/10.1038/nmeth.2575>
- 821 46. Tenaillon O, Skurnik D, Picard B, Denamur E. The population genetics of commensal  
822 *Escherichia coli*. Nat. Rev. Microbiol. [Internet]. 2010;8:207–17. Available from:  
823 <http://www.nature.com/doifinder/10.1038/nrmicro2298>
- 824 47. Holt KE, Wertheim H, Zadoks RN, Baker S, Whitehouse CA, Dance D, et al. Genomic analysis  
825 of diversity, population structure, virulence, and antimicrobial resistance in *Klebsiella*  
826 *pneumoniae*, an urgent threat to public health. Proc. Natl. Acad. Sci. [Internet].  
827 2015;112:E3574–81. Available from:  
828 <http://www.pnas.org/lookup/doi/10.1073/pnas.1501049112>
- 829 48. Dehoux P, Marvaud JC, Abouelleil A, Earl AM, Lambert T, Dauga C. Comparative genomics  
830 of *Clostridium bolteae* and *Clostridium clostridioforme* reveals species-specific genomic  
831 properties and numerous putative antibiotic resistance determinants. BMC Genomics  
832 [Internet]. 2016;17:819. Available from:  
833 <http://bmcbgenomics.biomedcentral.com/articles/10.1186/s12864-016-3152-x>
- 834 49. Donati C, Hiller NL, Tettelin H, Muzzi A, Croucher NJ, Angiuoli S V, et al. Structure and  
835 dynamics of the pan-genome of *Streptococcus pneumoniae* and closely related species.  
836 Genome Biol. [Internet]. 2010;11:R107. Available from:  
837 <http://genomebiology.biomedcentral.com/articles/10.1186/gb-2010-11-10-r107>
- 838 50. Touchon M, Hoede C, Tenaillon O, Barbe V, Baeriswyl S, Bidet P, et al. Organised Genome  
839 Dynamics in the *Escherichia coli* Species Results in Highly Diverse Adaptive Paths. Casades S J,  
840 editor. PLoS Genet. [Internet]. 2009;5:e1000344. Available from:

- 841 <http://dx.plos.org/10.1371/journal.pgen.1000344>
- 842 51. Castellarin M, Warren RL, Freeman JD, Dreolini L, Krzywinski M, Strauss J, et al.  
843 *Fusobacterium nucleatum* infection is prevalent in human colorectal carcinoma. *Genome Res.*  
844 2012;22:299–306.
- 845 52. Finegold SM, Song Y, Liu C, Hecht DW, Summanen P, Könönen E, et al. *Clostridium*  
846 *clostridioforme*: A mixture of three clinically important species. *Eur. J. Clin. Microbiol. Infect.*  
847 *Dis.* 2005;24:319–24.
- 848 53. Xu J, Mahowald MA, Ley RE, Lozupone CA, Hamady M, Martens EC, et al. Evolution of  
849 Symbiotic Bacteria in the Distal Human Intestine. *PLoS Biol.* 2007;5:1574–86.
- 850 54. Dagan T, Artzy-Randrup Y, Martin W. Modular networks and cumulative impact of lateral  
851 transfer in prokaryote genome evolution. *Proc. Natl. Acad. Sci.* [Internet]. 2008;105:10039–  
852 44. Available from: <http://www.pnas.org/cgi/doi/10.1073/pnas.0800679105>
- 853 55. Silverman BW. Using Kernel Density Estimates to Investigate Multimodality. *J. R. Stat. Soc.*  
854 *Ser. B.* 1981;43:97–9.
- 855 56. Maronna RA, Martin RD, Yohai VJ. Robust Statistics: Theory and Methods. *Ann. Stat.*  
856 2006;30:17–23.
- 857 57. Leys C, Ley C, Klein O, Bernard P, Licata L. Detecting outliers: Do not use standard deviation  
858 around the mean, use absolute deviation around the median. *J. Exp. Soc. Psychol.*  
859 2013;49:764–6.
- 860 58. Bland JM, Altman DG. Statistics Notes: Transforming data. *Bmj* [Internet]. 1996;312:770–  
861 770. Available from: <http://www.bmj.com/cgi/doi/10.1136/bmj.312.7033.770>
- 862 59. Osborne J. Notes on the use of data transformations. *Pract. Assess.* [Internet]. 2002;8:4–  
863 13. Available from:  
864 <c:%5CUsers%5CTeresa%5CDocuments%5CTeseFinal%5C2010%5Cserinus2010%5CPDF'sPara>  
865 <Organizar%5CPDF's2010%5Ctransformation.pdf>
- 866 60. Lin LI-K. A Concordance Correlation Coefficient to Evaluate Reproducibility. *Biometrics*



- 867 [Internet]. 1989;45:255. Available from:  
868 <http://www.jstor.org/stable/2532051?origin=crossref>
- 869 61. Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, et al. Supporting data for the paper: “An  
870 integrated catalog of reference genes in the human gut microbiome” [Internet]. GigaScience  
871 Database; 2014. Available from: <http://gigadb.org/dataset/100064>
- 872 62. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. J.  
873 Mol. Biol. [Internet]. 1990;215:403–10. Available from:  
874 <http://www.sciencedirect.com/science/article/pii/S0022283605803602>
- 875 63. Ogata H, Goto S, Sato K, Fujibuchi W, Bono H, Kanehisa M. KEGG: Kyoto encyclopedia of  
876 genes and genomes. *Nucleic Acids Res.* 1999. p. 29–34.
- 877 64. Pruitt KD, Tatusova T, Maglott DR. NCBI reference sequences (RefSeq): A curated non-  
878 redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.*  
879 2007;35.
- 880 65. Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, von Mering C, et al. Fast  
881 Genome-Wide Functional Annotation through Orthology Assignment by eggNOG-Mapper.  
882 *Mol. Biol. Evol.* [Internet]. 2017;278:631–7. Available from: [https://oup.silverchair-  
883 cdn.com/oup/backfile/Content\\_public/Journal/mbe/PAP/10.1093\\_molbev\\_msx148/3/msx1  
884 48.pdf?Expires=1499795422&Signature=W7WZu1nfufNHcrf1QFCH~13BjwFo82LccucPBoCvq  
885 GvLaClAQFlzzY6iT3cfr4A1bFPxYnhBaz6D-wBoorppxD5SfKpgd63KRkl2HHDcb2BEBArfT4f](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/mbe/PAP/10.1093_molbev_msx148/3/msx148.pdf?Expires=1499795422&Signature=W7WZu1nfufNHcrf1QFCH~13BjwFo82LccucPBoCvqGvLaClAQFlzzY6iT3cfr4A1bFPxYnhBaz6D-wBoorppxD5SfKpgd63KRkl2HHDcb2BEBArfT4f)
- 886 66. Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC, et al. EGGNOG 4.5:  
887 A hierarchical orthology framework with improved functional annotations for eukaryotic,  
888 prokaryotic and viral sequences. *Nucleic Acids Res.* 2016;44:D286–93.
- 889 67. Eddy SR. Accelerated profile HMM searches. *PLoS Comput. Biol.* 2011;7.
- 890 68. Kultima JR, Sunagawa S, Li J, Chen W, Chen H, Mende DR, et al. MOCAT: A Metagenomics  
891 Assembly and Gene Prediction Toolkit. *PLoS One.* 2012;7.
- 892 69. Nelson KE, Weinstock GM, Highlander SK, Worley KC, Creasy HH, Wortman JR, et al. A  
893 Catalog of Reference Genomes from the Human Microbiome. *Science* (80-. ). [Internet].



894 2010;328:994–9. Available from:  
895 <http://www.sciencemag.org/cgi/doi/10.1126/science.1183605>

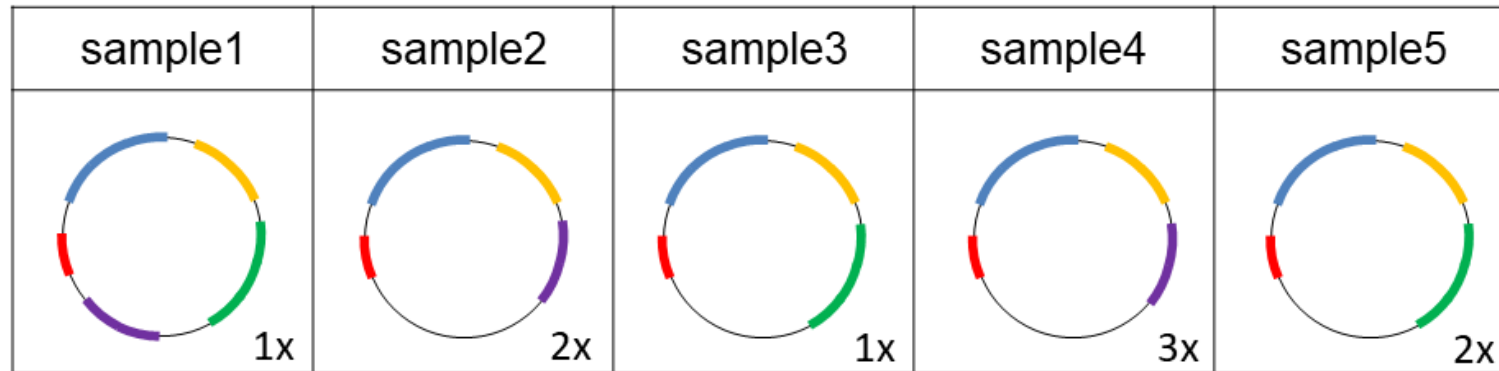
896 70. Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, et al. GenBank.  
897 Nucleic Acids Res. 2013;41.

898 71. Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic gene  
899 recognition and translation initiation site identification. BMC Bioinformatics. 2010;11:119.






900 72. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, et al. Circos: An  
901 information aesthetic for comparative genomics. Genome Res. [Internet]. 2009;19:1639–45.  
902 Available from: <http://genome.cshlp.org/cgi/doi/10.1101/gr.092759.109>

903

904 **Figures**



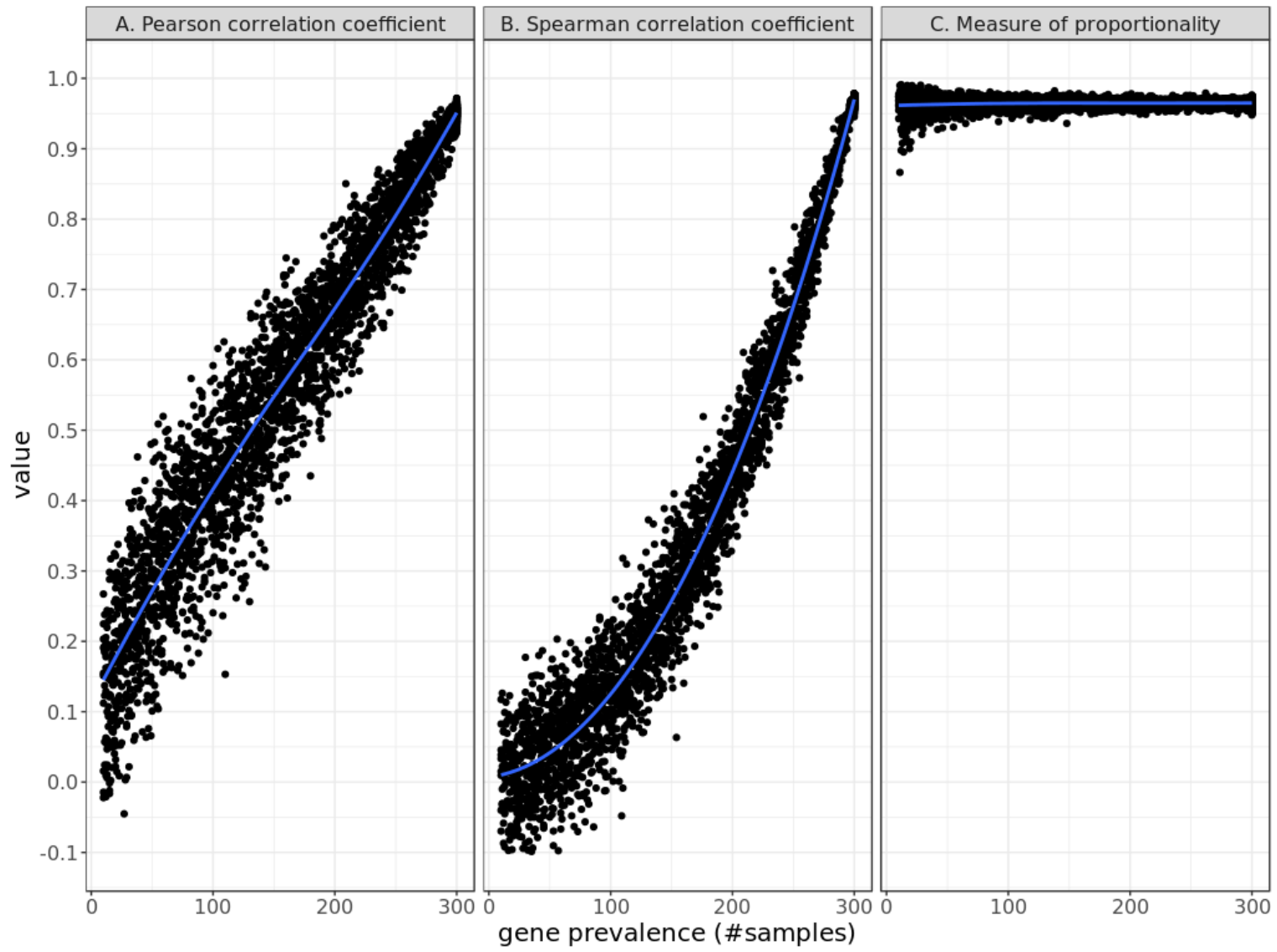
Shotgun sequencing

		sample1	sample2	sample3	sample4	sample5	
core gene 1		1	2	1	3	2	Co-abundant genes
core gene 2		3	6	3	9	6	
core gene 3		2	4	2	6	4	
accessory gene 1		3	0	3	0	6	<b>Partially</b> co-abundant with core genes
accessory gene2		0	4	2	6	0	

905

906

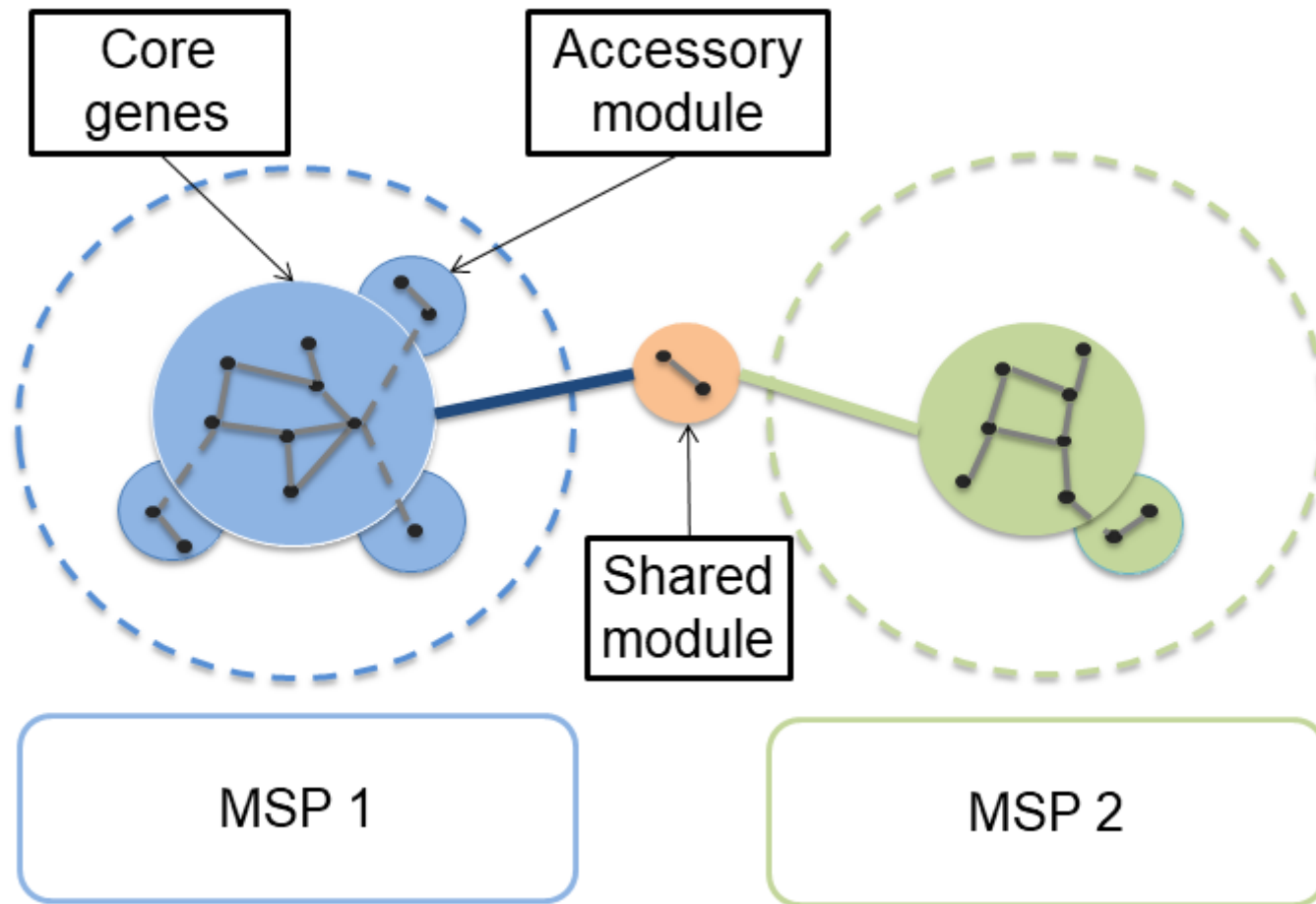
Figure 1



907

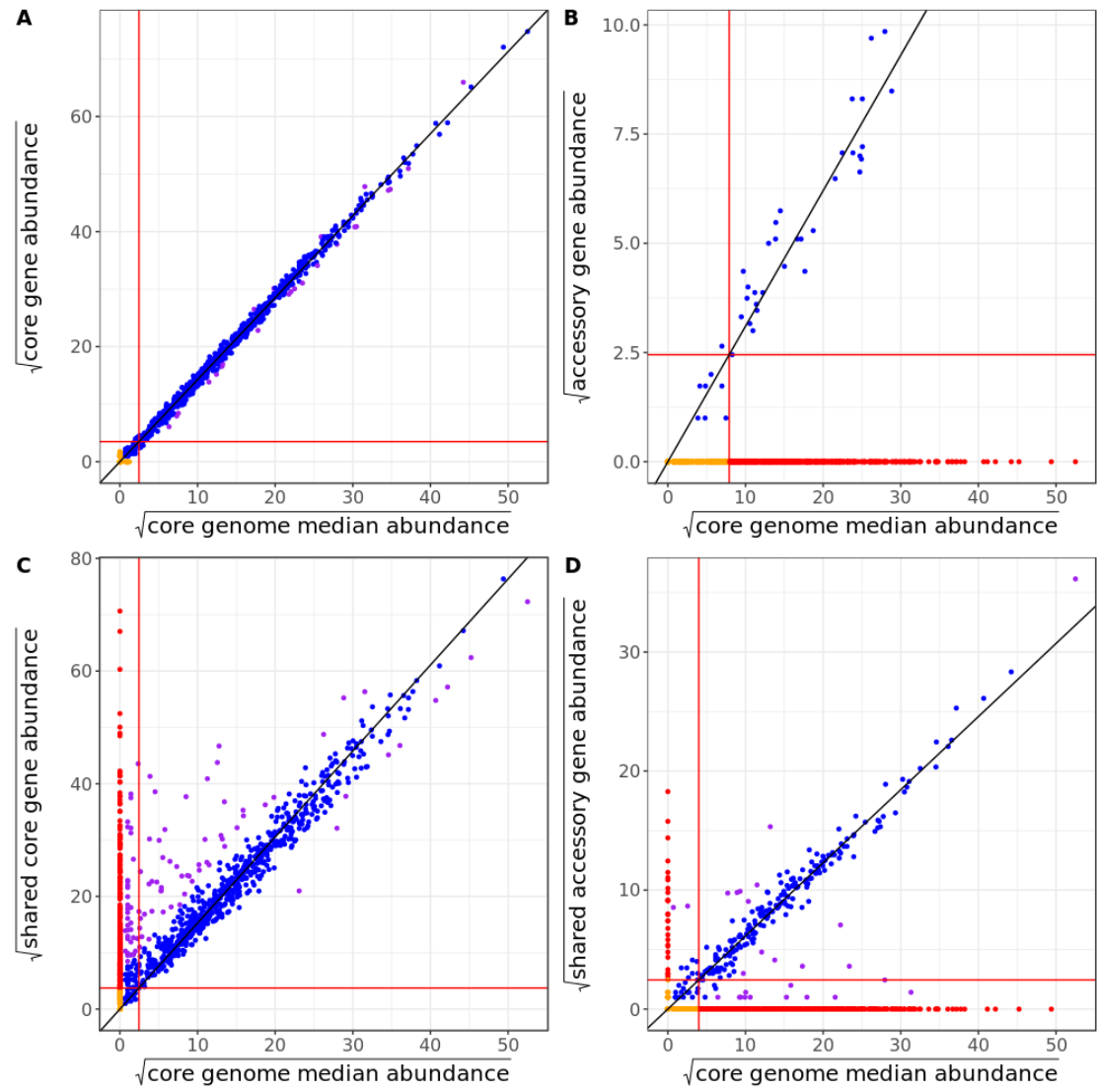
908

Figure 2



909  
910

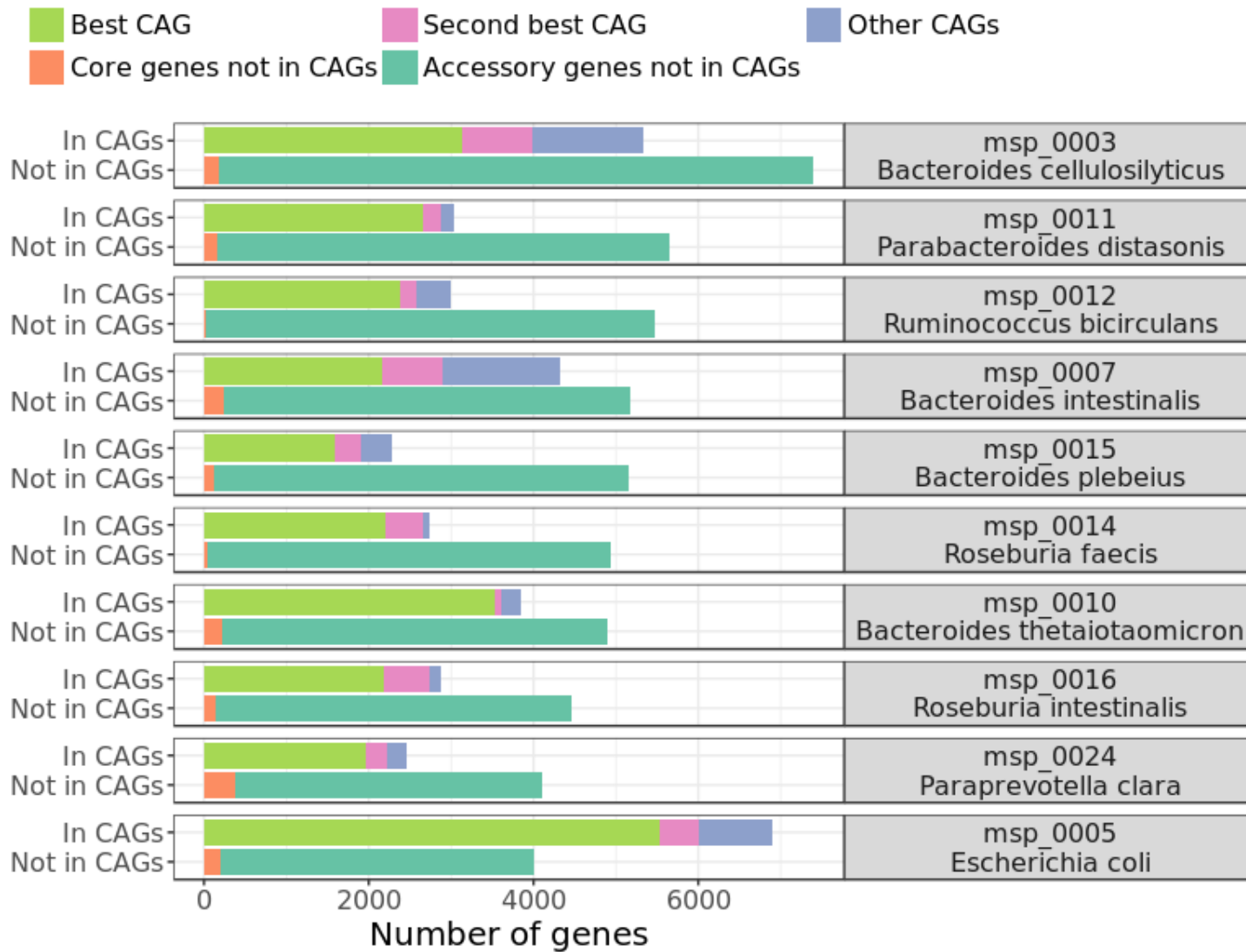
Figure 3



911

912

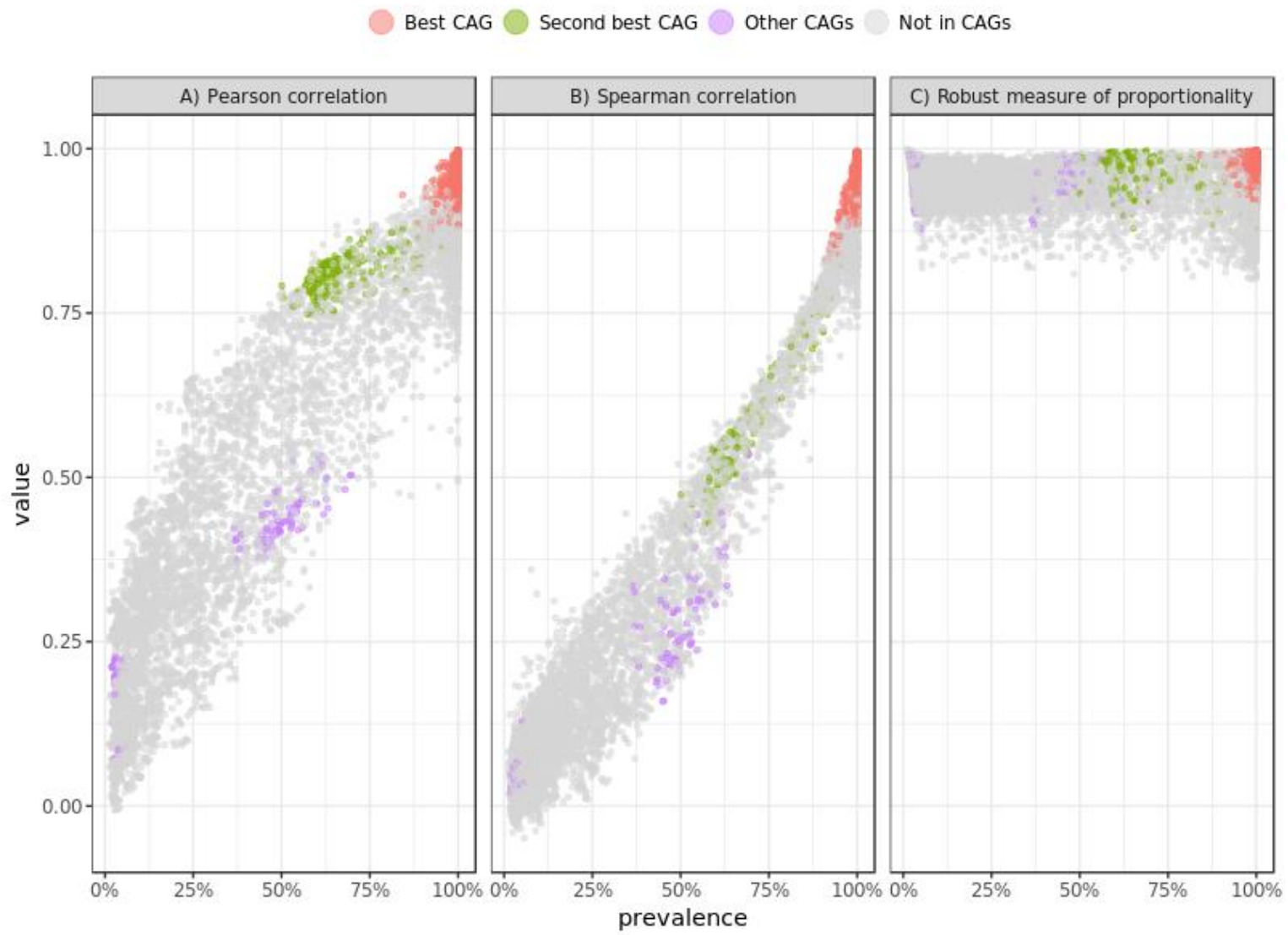
Figure 4



913

914

Figure 5



915

916

Figure 6

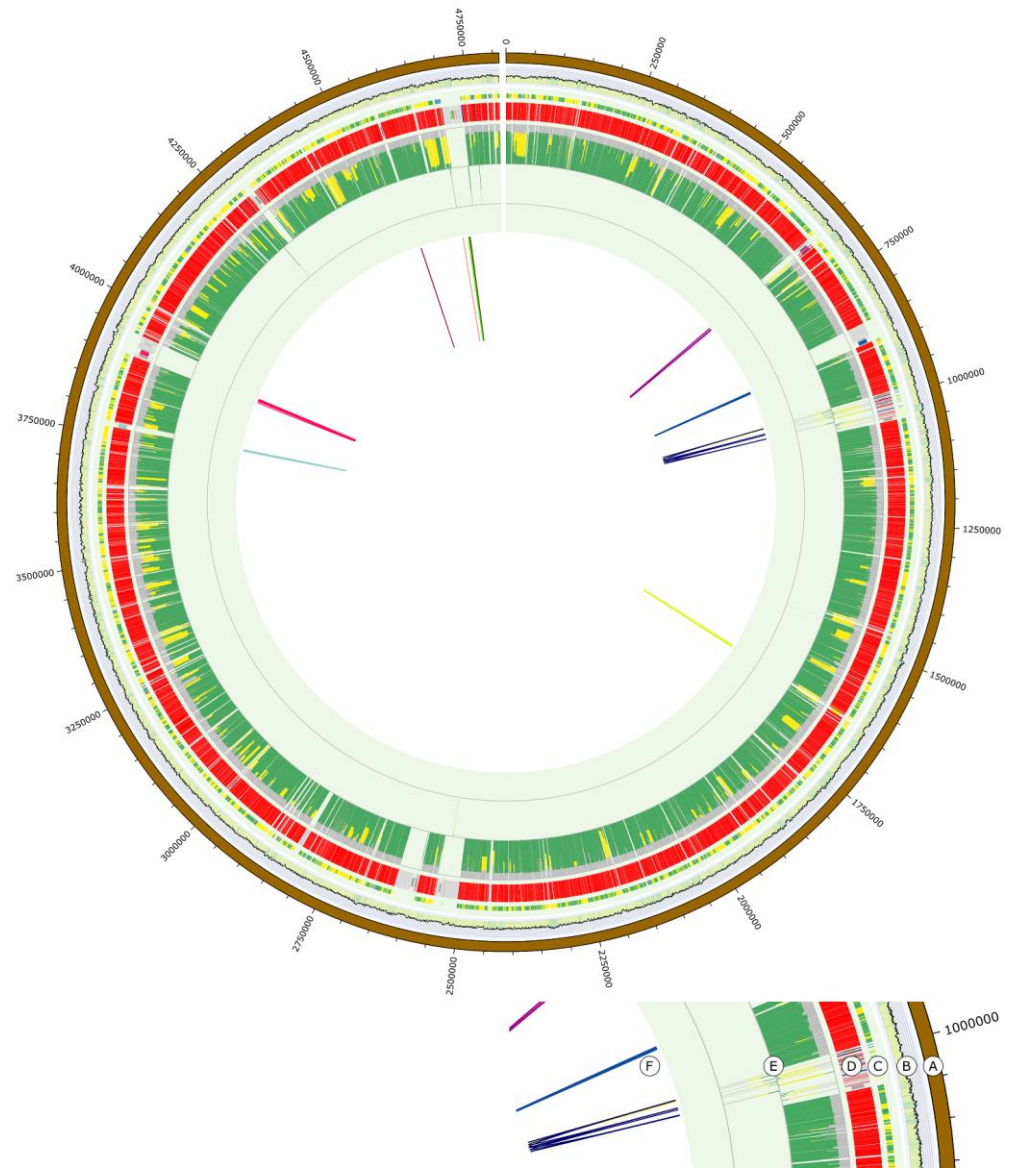


Figure 7

917

918



## 919 **Figures legends**

### 920 **Figure 1: Simple model illustrating the rationale behind the method**

921 5 samples carry different strains from the same species. Three core genes (red, blue, orange)  
922 are present in all the strains. Two accessory genes (green, purple) are present only in some  
923 strains. The abundance of the species in each sample ranges from 1 to 3 copies.

924 After shotgun sequencing, a raw gene abundance matrix is built.

925 A strict proportionality relationship is expected between two core genes, the proportionality  
926 coefficient being equal to the ratio of their length. In contrast, such relationship between a  
927 core and an accessory gene should be observed only in the subset of samples where the  
928 accessory gene is present.

### 929 **Figure 2: Comparison of the genes abundance profiles of the virtual species to the** 930 **median signal of its core genome**

931 using:

- 932 A. the Pearson correlation coefficient
- 933 B. the Spearman correlation coefficient
- 934 C. the measure of proportionality

### 935 **Figure 3: Structure of Metagenomic Species Pan-genomes (MSPs)**

### 936 **Figure 4: Illustration of the four types of genes in a MSP**

937 The core genome median abundance of the msp\_0043 (*Ruminococcus bromii*) is compared to:

- 938 A. the gene MH0003\_GL0010264 classified as core. The gene is detected in all the  
939 samples where MSP is detected.
- 940 B. the gene MH0025\_GL0082295 classified as accessory. The gene is missing in 516  
941 samples where the MSP is detected.
- 942 C. the gene 657321.RBR\_R\_22270 classified as its shared core. The gene is present in all  
943 the samples where MSP is detected but also in 286 samples where the MSP is not.
- 944 D. the gene MH0205\_GL0102923 classified as shared accessory. The gene is missing in  
945 454 samples where the MSP is detected but present in 28 samples where the MSP is  
946 not.

947 **Figure 5: Comparison of the gene content of some MSPs and their corresponding**  
948 **CAGs**

949 **Figure 6: Comparison of the genes abundance profiles of the msp\_0011**  
950 **(*Parabacteroides distasonis*) to the median signal of its core genome.**

951 Three measures are compared:

- 952 A. the Pearson correlation coefficient
- 953 B. the Spearman correlation coefficient
- 954 C. the measure of proportionality.

955 Grey points correspond to genes unclassified by Canopy whereas those colored were grouped  
956 in CAGs.

957 **Figure 7: Circos representation of the mapping of the msp\_0011 on the genome of**  
958 ***P. distasonis* strain ATCC 8503**

959 Description of layers from outside to inside:

- 960 A. Position on chromosome
- 961 B. GC-content (format: histogram)
- 962 C. gene or module type (format: highlight):
  - 963 • green: core
  - 964 • yellow: accessory
  - 965 • blue: shared core
  - 966 • purple: shared accessory
- 967 D. MSP (format: highlight):
  - 968 • Bandwidth:
    - 969 · wide: gene grouped in a MSP
    - 970 · narrow: gene grouped in a seed
  - 971 • color code:
    - 972 · red: gene grouped in the most represented MSP
    - 973 · other color + grey: gene grouped in another MSP or a seed
- 974 E. Sample assignment (format: histogram):
  - 975 • facing outwards if the gene is related to the most represented MSP, facing  
976 inwards otherwise.

- 977
- color code:
    - 978 · grey: samples where the MSP module is not detected
    - 979 · green: samples where the MSP core and the gene are detected
    - 980 · yellow: samples where the MSP core is detected but not the gene
    - 981 · purple: samples where the gene is detected but not the MSP core
- 982 F. relation between genes associated to a MSP different from the most represented
- 983 (format: edges). Genes from the same alien MSP are linked by edges.
- 984