

## Benchmarking taxonomic assignments based on 16S rRNA gene profiling of the microbiota from commonly sampled environments

--Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-18-00043	
<b>Full Title:</b>	Benchmarking taxonomic assignments based on 16S rRNA gene profiling of the microbiota from commonly sampled environments	
<b>Article Type:</b>	Technical Note	
<b>Funding Information:</b>	European Bioinformatics Institute	Dr Alexandre Almeida
	H2020 Research Infrastructures (676559)	Not applicable
<b>Abstract:</b>	<p><b>Background:</b> Taxonomic profiling of ribosomal RNA (rRNA) sequences has been the accepted norm for inferring the composition of complex microbial ecosystems. QIIME and mothur have been the most widely used taxonomic analysis tools for this purpose, with MAPseq and QIIME 2 being two recently released alternatives. However, no independent and direct comparison between these four main tools has been performed. Here, we compared MAPseq, mothur, QIIME, and QIIME 2 using synthetic simulated datasets comprised of some of the most abundant genera found in the human gut, ocean and soil environments. We evaluate their accuracy when paired with both different reference databases and variable sub-regions of the 16S rRNA gene.</p> <p><b>Findings:</b> We show that QIIME 2 with the SILVA 128 database provided the best recall at the genus level, and the lowest distance estimates between the observed and simulated samples. However, MAPseq showed the highest precision, with miscall rates consistently below 2%. Notably, QIIME 2 was the most computationally expensive tool, with CPU time and memory usage almost two and 30 times higher than MAPseq, respectively. Using the SILVA database generally yielded a higher recall than using Greengenes, while assignment results of different 16S rRNA variable sub-regions varied by up to 35% between samples analysed with the same pipeline.</p> <p><b>Conclusions:</b> Our results support the use of either QIIME 2 or MAPseq for optimal 16S rRNA gene profiling, and we suggest that the choice between the two should be based on the level of recall, precision and/or computational performance required.</p>	
<b>Corresponding Author:</b>	Alexandre Almeida	
	UNITED KINGDOM	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>		
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Alexandre Almeida	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Alexandre Almeida	
	Alex L Mitchell	
	Aleksandra Tarkowska	
	Robert D Finn	
<b>Order of Authors Secondary Information:</b>		
<b>Opposed Reviewers:</b>		
<b>Additional Information:</b>		

Question	Response
<p>Are you submitting this manuscript to a special series or article collection?</p>	<p>No</p>
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	<p>Yes</p>
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

1 **Benchmarking taxonomic assignments based on 16S rRNA gene profiling**  
2 **of the microbiota from commonly sampled environments**

3

4 Alexandre Almeida<sup>1,2,\*</sup>, Alex L. Mitchell<sup>1</sup>, Aleksandra Tarkowska<sup>1</sup> and Robert D. Finn<sup>1</sup>

5

6 <sup>1</sup>EMBL-EBI European Bioinformatics Institute, Wellcome Genome Campus, Hinxton,  
7 Cambridge CB10 1SD, UK; <sup>2</sup>Wellcome Trust Sanger Institute, Wellcome Genome Campus,  
8 Hinxton CB10 1SA, UK

9

10 \*Corresponding author:

11 Alexandre Almeida,

12 EMBL-EBI European Bioinformatics Institute,

13 Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK

14 Tel + 44 (0) 1223 494 468 E-mail: [aalmeida@ebi.ac.uk](mailto:aalmeida@ebi.ac.uk)

15

16 **Keywords:** 16S rRNA gene, human gastrointestinal tract, ocean, microbiome, soil, taxonomy

17

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 18 **Abstract**

19 **Background:** Taxonomic profiling of ribosomal RNA (rRNA) sequences has been the  
20 accepted norm for inferring the composition of complex microbial ecosystems. QIIME and  
21 mothur have been the most widely used taxonomic analysis tools for this purpose, with  
22 MAPseq and QIIME 2 being two recently released alternatives. However, no independent  
23 and direct comparison between these four main tools has been performed. Here, we compared  
24 MAPseq, mothur, QIIME, and QIIME 2 using synthetic simulated datasets comprised of  
25 some of the most abundant genera found in the human gut, ocean and soil environments. We  
26 evaluate their accuracy when paired with both different reference databases and variable sub-  
27 regions of the 16S rRNA gene.

29 **Findings:** We show that QIIME 2 with the SILVA 128 database provided the best recall at  
30 the genus level, and the lowest distance estimates between the observed and simulated  
31 samples. However, MAPseq showed the highest precision, with miscall rates consistently  
32 below 2%. Notably, QIIME 2 was the most computationally expensive tool, with CPU time  
33 and memory usage almost two and 30 times higher than MAPseq, respectively. Using the  
34 SILVA database generally yielded a higher recall than using Greengenes, while assignment  
35 results of different 16S rRNA variable sub-regions varied by up to 35% between samples  
36 analysed with the same pipeline.

38 **Conclusions:** Our results support the use of either QIIME 2 or MAPseq for optimal 16S  
39 rRNA gene profiling, and we suggest that the choice between the two should be based on the  
40 level of recall, precision and/or computational performance required.

## 43 Findings

### 44 Background

45 Genome sequencing has provided an unprecedented view of the microbial diversity of  
46 ecosystems from wide-ranging environments. For example, the commensal flora of the  
47 human gut has been extensively explored for potential associations with the onset of many  
48 human diseases [1–3]. Similarly, the rich microbial diversity of environments such as soil and  
49 oceans have been studied in depth, yielding important ecological inferences [4–6]. There are  
50 now a substantial number of such microbial community datasets deposited in sequence  
51 archives (for example, the European Nucleotide Archive currently holds over 600 000  
52 environmental samples [7]) and the rate of deposition is increasing. Drawing relevant  
53 biological correlations from this vast amount of data requires accurate and reliable tools and  
54 methods.

55  
56 One of the crucial steps in almost all microbiome-based analyses is inference of community  
57 composition through taxonomic classification. For a few decades now [8], the common  
58 approach for taxonomic assignment of microbial species has been the classification of  
59 ribosomal RNA (rRNA) sequences. Currently, the most widely used tools for this purpose are  
60 the mothur [9] and “Quantitative Insights Into Microbial Ecology” (QIIME) software  
61 packages [10]. Both tools take individual genetic markers (e.g. the 16S rRNA gene,  
62 conserved across the prokaryotic domains) and compare them to a reference database,  
63 assigning a taxonomic lineage to each of the queried sequences. Greengenes [11], NCBI [12],  
64 RDP [13] and SILVA [14] are some of the most widely used rRNA sequence databases.  
65 Ultimately, the success of these analyses is not only dependent on the breadth and diversity  
66 of annotated sequences available in public repositories, but also on the accuracy of the  
67 classification algorithms used by each of the tools. By default, QIIME makes use of the

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

68 UCLUST clustering method [15] to assign biological sequences to a reference database,  
69 while mothur wraps the naïve Bayesian RDP classifier, developed by Wang, *et al.* [16], for  
70 sequence classification. Two other tools — MAPseq [17] and QIIME 2 (<https://qiime2.org/>)  
71 — have recently been released, providing additional assignment methods. QIIME 2 also  
72 makes use of a naïve Bayes classifier [18], and MAPseq is a *k-mer* search approach that  
73 outputs confidence estimates at different taxonomic ranks.

74  
75 A community-driven initiative known as the “Critical Assessment of Metagenome  
76 Interpretation” (CAMI) benchmarked a range of software tools for the analysis of shotgun  
77 metagenomic datasets [19]. In regard to amplicon-based approaches, previous studies have  
78 mainly evaluated the classification methods of QIIME and mothur, highlighting some of their  
79 advantages and pitfalls [20–22]. However, until now, no independent study has compared the  
80 accuracy of these methods to MAPseq and QIIME 2 whilst also taking into account potential  
81 differences arising from the use of distinct reference databases. Furthermore, for genotyping  
82 the 16S rRNA gene there is also much debate within the scientific community on the most  
83 informative variable sub-region to target [23]. Strong arguments have been made towards  
84 sequencing specific or combined sub-regions, such as the V4 [24] and V3-V4 [25], while  
85 difficulty in amplifying bacterial species, such as those from the *Actinobacteria* group, has  
86 prompted the development of more specialized primers [26,27]. The impact of variable  
87 region choice on the taxonomic classification performance of different tools or databases is  
88 therefore also important to assess.

89  
90 The use of mock communities in microbiome studies has revealed that different experimental  
91 conditions and methods dramatically affect the quality of the results [28–31]. In contrast, *in*

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

92 *silico* benchmarking approaches provide an agnostic view on the efficiency of the  
93 computational pipelines, independently of experimental variation and technical biases.

94  
95 Therefore, in this work we have leveraged a set of simulated 16S rRNA gene sequences  
96 representative of genera commonly found in the human gut, ocean and soil environments, to  
97 evaluate the accuracy of MAPseq, mothur, QIIME and QIIME 2 with different reference  
98 databases, and according to some of the most commonly targeted sub-regions of the 16S  
99 rRNA gene. We show that, regardless of the database used, QIIME 2 outperformed all other  
100 tools in terms of overall recall at both genus and family levels, as well as in distance  
101 estimations between the observed and predicted samples. Considerable performance  
102 differences were observed between using distinct 16S rRNA gene sub-regions, while limited  
103 software-dependent variation was seen between different reference databases. We believe this  
104 work will help inform microbial ecologists about important decisions to take when designing  
105 new 16S rRNA-based community studies.

### 106 107 **Composition of the simulated datasets**

108 The microbiota colonizing the human gut, ocean and soil environments are some of the most  
109 frequently studied microbial communities. Hence, to provide data with direct practical  
110 applications, we focused on simulating datasets containing a diverse set of genera commonly  
111 found in these three ecosystems (Additional file: Fig. S1). Representative genera were  
112 selected after identifying the 80 most abundant genera across publicly available metagenomes  
113 from human gut, ocean and soil [7]. Then, for each biome, four different communities were  
114 generated with two levels of diversity: samples A100 and B100 with a random set of 100  
115 species belonging to these genera; and A500 and B500 with 500 species. Final datasets  
116 comprised a total of 66, 66 and 76 different genera from the human gut, oceanic and soil

117 environments, respectively. For the purpose of this benchmarking, we simulated the datasets  
118 with a similar relative abundance per genus to avoid introducing any taxon-specific biases  
119 (Additional file: Fig. S1).

120  
121 To simulate a realistic scenario, where variation can arbitrarily occur and sequences may not  
122 have an exact representative in public databases, we randomly mutated 2% of the positions of  
123 each 16S rRNA sequence retrieved. Using an *in silico* PCR, we then extracted each sub-  
124 region using commonly used primer sequences (Additional file: Table S1). Notably, the  
125 percentage of bacterial sequences from the Greengenes and SILVA databases matching the  
126 primers selected for V1-V2 was dramatically lower (37.6%) than that of V3-V4 (99.2%), V4  
127 (99.1%) and V4-V5 (99.4%) (Additional file: Fig. S2). The 16S rRNA V1 sub-region had  
128 been previously found to be truncated in a substantial number of reference sequences [23].  
129 Our results confirm this observation and again raise caution at the use of the 16S V1-V2  
130 rRNA primer sequences for complex and diverse samples, due to the reduced number of  
131 reference sequences available.

132

### 133 **Taxonomic assignment**

134 Microbiome studies frequently strive to associate microbial diversity signatures with a  
135 phenotype of interest. However, focusing solely on high-level taxonomic ranks can severely  
136 underestimate the degree of variation observed between sample groups. To circumvent this,  
137 highly discriminative approaches are needed to be able to pinpoint the most significant taxa  
138 warranting further validation. For assessing the performance of MAPseq, mothur, QIIME and  
139 QIIME 2 with different reference databases (Additional file: Fig. S3), we limited our analyses  
140 to classification at the lineage level instead of operational taxonomic units (OTUs), as it  
141 allows a more consistent and easier interpretation of the results. Species assignment of every



142 queried sequence would be the desired outcome, but the limited resolution of the 16S rRNA  
143 locus precludes an accurate classification at this level. In fact, databases such as RDP do not  
144 report taxon names below genus. In this work, we calculated the degree of recall and  
145 precision at the genus and family ranks, as in our opinion they provide the best compromise  
146 between classification accuracy and resolution.

147  
148 By comparing the level of recall across all software tools, we found that QIIME 2 recovered  
149 the largest proportion of sequences from the expected genera (Table 1 and Fig. 1). Combined  
150 with the SILVA database, this resulted in the highest recall (sensitivity) for human gut  
151 (67.1%) and soil samples (67.7%), while the Greengenes database outperformed in the case  
152 of the oceanic microbiome (79.8%). In fact, all tools except QIIME saw a decrease in recall  
153 when using SILVA specifically for the classification of the oceanic dataset. Globally,  
154 however, SILVA most frequently provided a better genus recall than Greengenes (five out of  
155 nine comparisons across MAPseq, QIIME and QIIME 2, Fig. 1). In terms of correctly  
156 identified taxa, MAPseq in conjunction with SILVA detected the greatest number of expected  
157 genera in all three biomes (Fig. 1). At the family level, all tools presented a substantially  
158 higher recall (Table 1), with QIIME 2 reaching 93.9% in the human gut sample, 95.8% with  
159 the ocean set and 91.3% with the soil sample.

160  
161 Although the level of recall is a crucial metric in choosing the most appropriate taxonomic  
162 classification pipeline, it is equally important to ensure a low frequency of false-positive  
163 assignments. We evaluated the degree of precision (specificity) by the percentage of  
164 sequences assigned to the wrong taxon out of all the detected taxa. Accuracy was high for all  
165 the tools, with precision estimates of at least 85% across all analysis pipelines (Fig. 2A). In  
166 terms of total number of sequences, this translated to less than 9% of the reads misassigned at

167 the genus level (Additional files: Tables S2, S3, S4 and Fig. S4). MAPseq with the SILVA  
168 database consistently outperformed all other tools, with a precision above 96% for the three  
169 tested biomes (Fig. 2A), equating to less than 2% of miscalled sequences (Additional file:  
170 Fig. S4).

171

172 To combine both recall and precision into a single metric, we calculated the F-score for all  
173 taxonomic assignments (Fig. 2A and Additional file: Fig. S5). At both genus and family  
174 levels, we found that QIIME 2 had the highest score across the samples representative of the  
175 three different biomes, with the SILVA database coming out on top for the human gut  
176 (genus: 78.9%, family: 96.6%) and soil (genus: 78%, family: 94.1%) environments in  
177 particular, but the Greengenes database performing better with the oceanic dataset (genus:  
178 87.6%, family: 97.2%). After fractioning the data according to different sub-regions of the  
179 16S rRNA gene, we then repeated the same analysis (Fig. 2B). This revealed that the  
180 performance of each tool varied up to 35% depending on the 16S rRNA sub-region targeted.  
181 Notably, the V1-V2 or V3-V4 sub-regions performed the best across most of the pipelines  
182 (Fig. 2B). In our study, each synthetic species had a genetically close full-length 16S rRNA  
183 sequence represented in the databases, so our tests were probably not significantly affected  
184 by the reduced number of V1-V2 reference sequences available.

185

186 The ongoing surge in genome sequencing is producing thousands of novel sequences each  
187 year. Therefore, efficient tools that can scale up to provide analysis of tens of thousands of  
188 samples is increasingly important. With this in mind, we compared the computational  
189 performance of MAPseq, mothur, QIIME and QIIME 2 throughout the whole classification  
190 pipeline of our simulated datasets. We analysed average memory usage and CPU time across  
191 the three biomes for the processing and assignment of 3 million quality-filtered sequences

192 against the SILVA 128 database (Fig. 3). MAPseq was the most memory-efficient tool, with  
193 mothur, QIIME and QIIME 2 requiring over 72, 6 and 27 times more memory resources,  
194 respectively (Fig. 3A). CPU time of QIIME 2 was the highest, close to twice that of MAPseq,  
195 and over 200 times longer than QIIME, which was the fastest (Fig. 3B). Of note is that each  
196 pipeline has its own processing procedure; both the mothur and QIIME 2 pipelines included a  
197 de-replication step of the query sequences prior to taxonomic assignment, which substantially  
198 reduces the number of sequences used for classification.

### 200 **Relative quantification and beta diversity**

201 One of the main aspects of any microbiome-based analysis is the assessment of the  
202 differential abundance and beta diversity across a set of sample groups. In this respect,  
203 accurate estimation of the relative abundance of each taxon is essential to find statistically  
204 significant patterns. To assess how accurately each tool was able to predict taxa relative  
205 abundances in each sample, we calculated dissimilarity scores (DS) for each genus present in  
206 the simulated dataset (Fig. 4). Interestingly, QIIME 2 showed the most accurate prediction in  
207 relation to the true genera composition, with an average DS of 0.32 when used in conjunction  
208 with the SILVA database (Table 1). In terms of the reference database used, analyses carried  
209 out with SILVA consistently yielded more accurate predictions than with the Greengenes  
210 database. Substantial differences in accuracy were observed across different genera, with  
211 sequences from the *Paraprevotella* genus — frequently present in human gut samples —  
212 more accurately predicted, in contrast to those from *Hyphomicrobium*, *Thalassobacter* and  
213 *Verrucomicrobium* — commonly found in oceanic biomes — which had the worst results  
214 (Fig. 4). These genera might either be underrepresented in the reference databases, or have a  
215 high degree of conservation with other closely related taxa, making accurate taxonomic  
216 assignments more challenging.

217 For a global assessment of the beta diversity across samples, we performed a principal  
218 coordinates analysis (PCoA) and calculated both Bray-Curtis and Jaccard distances between  
219 the observed and expected results. Both distance methods represent complementary  
220 approaches, as the Bray-Curtis metric corresponds to a quantitative evaluation of the  
221 dissimilarity across samples, whereas the Jaccard index is a qualitative measure of  
222 community similarity. We found that samples analysed with QIIME 2 were the closest (i.e.  
223 had the lowest distance estimate) to the true simulated datasets, with minor differences  
224 between the use of SILVA or Greengenes with both the Bray-Curtis and Jaccard methods  
225 (Table 1; Fig. 5).

## 227 **Discussion**

228 With the number of tools, databases and options available for taxonomic classification of  
229 marker sequences, it can be a daunting task to decide the optimal approach for analysis of a  
230 specific dataset. In this work, we have strived to help guide this decision-making process by  
231 independently assessing the performance of the most commonly used taxonomic assignment  
232 strategies with simulated samples comprised of genera found in frequently sampled  
233 environments.

234  
235 Overall, we show that all tools we tested performed moderately well, with high precision and  
236 modest-to-high recall rates at the genus level. QIIME 2 presents significant improvements  
237 over the other tools, particularly over the preceding version of QIIME, in regard to detection  
238 sensitivity at both family and genus level. The superiority of QIIME 2 also held true for the  
239 prediction of sample composition, as beta diversity estimates between the analysed and  
240 simulated communities were the closest using this method. Therefore, these data support the  
241 use of QIIME 2 to obtain the largest proportion of classified sequences at the most accurate

242 relative abundances. Nevertheless, the results did show MAPseq to be a more conservative  
243 but specific approach, meaning that fewer genera were misassigned. This tool also showed  
244 considerably better computational performance than QIIME 2, requiring approximately 30  
245 times less memory and almost half the CPU time to process the same dataset (even though  
246 QIIME 2 classifies substantially fewer query sequences due to a prior de-replication step).  
247 These results show that MAPseq provides a credible option if precision and computational  
248 performance or scale are a priority. It should be noted that selecting a single best software  
249 package is not a straightforward affair, and we expect that differences in performance will be  
250 observed with different real-world datasets. Furthermore, aside from the software packages  
251 we tested, other web-based tools such as BioMaS [22] are also available. But, they are  
252 usually restricted to the use of specific reference databases, making individual customizations  
253 and accurate comparisons more challenging.

254  
255 In addition to choosing the right tool, combining that with the appropriate reference database  
256 is equally important to ensure the best classification performance. Greengenes and SILVA  
257 have been the most widely-used and readily supported databases. Generally, the SILVA 128  
258 database performed better than Greengenes 13\_8 in terms of recall at both family and genus  
259 levels, as well as in predicting the true taxa composition of the simulated communities.  
260 Conversely, there was an almost universal decrease in its performance in the detection of  
261 ocean-specific taxa, so special care should be taken in the analysis of datasets sampled from  
262 this particular environment. Nonetheless, there are additional advantages to the use of  
263 SILVA: it is more frequently updated (Greengenes was last updated in May 2013); it includes  
264 rRNA sequences of eukaryotic organisms in addition to archaea and bacterial species; and  
265 has been shown to be more easily comparable and mapped to other taxonomies such as the  
266 NCBI [32]. In the case of MAPseq and mothur, the NCBI and RDP databases also performed

267 well, with higher recall but slightly lower precision scores compared to SILVA. Therefore,  
268 the SILVA, RDP or NCBI databases are all appropriate choices for a comprehensive and  
269 accurate taxonomic analyses.

270

271 The choice of primer sequences for taxonomic profiling of the 16S rRNA gene has been a  
272 matter of frequent debate. In common with previously reported observations [27], we show  
273 that targeting different sub-regions can considerably influence the taxonomic assignment  
274 performance (by up to 35% in our analyses). Overall, the V1-V2 and V3-V4 sub-regions  
275 performed the best across most of the tools. However, the V1-V2 primers did not match more  
276 than 60% of the reference sequences across SILVA and Greengenes, so we discourage its use  
277 for classification of complex community samples. As our simulated datasets were generated  
278 from close representatives containing full-length 16S rRNA genes, it is reasonable to assume  
279 that our analysis of the V1-V2 sub-region was not significantly hampered by this reduced  
280 number of reference sequences. Kozich *et al.* [24] have argued in favour of standardizing the  
281 use of the V4 sub-region for Illumina MiSeq sequencing, as it allows complete overlap of  
282 paired-end sequences, mitigating sequence errors introduced during PCR amplification or  
283 sequencing. Phylogenetic studies have also showed that the V4 sub-region is the closest  
284 representative of the phylogenetic signal of the whole 16S rRNA locus [23]. Here, we  
285 analysed the performance of some of the most commonly used sub-regions under a purely  
286 computational perspective, and conclude that amplification of the V3-V4 sub-region is most  
287 frequently the best option for a reliable taxonomic inference.

288

289 In summary, we have identified the major benefits and drawbacks of the most recent and  
290 popular taxonomic classification methods. Importantly, we show that the choice of software,  
291 database and sub-region significantly affects the quality of the classification results. Given

1 292 the impact of each of these variables, it is imperative to strive for consistency in the analysis  
2 293 of samples not only within individual studies, but across different projects as well. Services  
3  
4 294 like the EBI metagenomics [7] and MG-RAST [33] help provide a basis for standardization,  
5  
6  
7 295 but additional factors relating to the experimental design are up to individual users to decide.  
8  
9 296 Hence, in this fast-evolving field, we believe the work presented here will help the  
10  
11 297 microbiome research community make more informed decisions about the most appropriate  
12  
13 298 methodological approach to take in their own analysis pipeline.  
14  
15  
16  
17 299

## 20 300 **Methods**

### 23 301 **Generating simulated datasets**

25 302 Twelve sets of synthetic communities were generated for evaluating the accuracy of the  
26  
27 303 taxonomic assignment pipelines: four each for human gut, ocean and soil environments. First,  
28  
29 304 the 80 most abundant genera across publicly deposited samples from these biomes were  
30  
31 305 retrieved using the EBI metagenomics API (<https://www.ebi.ac.uk/metagenomics/api/>) [7].  
32  
33  
34 306 This list was then used to randomly select either 100 (datasets A100 and B100) or 500  
35  
36  
37 307 species (datasets A500 and B500) belonging to these genera, allowing a maximum of 20 and  
38  
39 308 50 species per genus, respectively. 16S rRNA gene sequences were extracted from the  
40  
41 309 European Nucleotide Archive (ENA) and 2% of the positions were randomly mutated to  
42  
43 310 create nucleotide diversity, using a custom python script ([https://github.com/Finn-Lab/Tax-](https://github.com/Finn-Lab/Tax-Benchmarking)  
44  
45 [Benchmarking](https://github.com/Finn-Lab/Tax-Benchmarking)). From these mutated sequences, an *in silico* PCR was carried out with an  
46  
47 311 additional python script ([https://github.com/simonrharris/in\\_silico\\_pcr](https://github.com/simonrharris/in_silico_pcr)), targeting commonly  
48  
49 312 used regions for 16S rRNA profiling (Additional file: Table S1): V1-V2, V3-V4, V4 and V4-  
50  
51 313 V5. Sequencing reads were simulated from these amplicon sequences in duplicate with ART  
52  
53 314 [34], generating ~ 10 000 and ~ 200 000 paired-end reads of 250 bp per region to have  
54  
55 315 samples representing both low and high levels of sequencing depth.  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 317 **Sequence classification**

1  
2 318 Initial pre-processing and quality control was performed following the mothur standard  
3  
4  
5 319 operating procedure (SOP) [24], accessed on November 2017. Briefly, the *make.contigs*  
6  
7 320 command was used to align, filter and merge the paired-end reads into contigs. Subsequently,  
8  
9  
10 321 we used the *screen.seqs* command to filter out any sequences with ambiguous base calls. This  
11  
12 322 final set of quality controlled sequences was then assigned into taxonomic lineages with  
13  
14 323 MAPseq v1.2.2 [17], mothur v1.39.5 [9], QIIME 1.9.1 [10], and QIIME 2 v2017.11  
15  
16 324 (<https://qiime2.org/>). For each software, we evaluated the settings and databases most  
17  
18  
19 325 frequently used and recommended for optimal taxonomic classification (Additional file: Fig.  
20  
21  
22 326 S3). With MAPseq, we tested the default NCBI database (mapref 2.2), as well as Greengenes  
23  
24 327 13\_8 and the SILVA 128 database re-mapped to an eight-level taxonomy (available in  
25  
26 328 [ftp://ftp.ebi.ac.uk/pub/databases/metagenomics/mapseq\\_silva128](ftp://ftp.ebi.ac.uk/pub/databases/metagenomics/mapseq_silva128)). Each set of reference  
27  
28  
29 329 sequences was analysed following the internal clustering by MAPseq. Options *-tophits 80*  
30  
31  
32 330 and *-topotus 40* were used in combination with the *-outfmt simple* option. For QIIME 1.9.1,  
33  
34 331 the *pick\_closed\_reference\_otus.py* script was used with both the default Greengenes database  
35  
36 332 (13\_8) and with SILVA 128 clustered at 97% identity. Taxonomic assignment with mothur  
37  
38  
39 333 was carried out according to the MiSeq SOP [24], excluding the chimera detection and  
40  
41 334 removal steps, using the available pre-formatted SILVA 128 database for alignment and  
42  
43  
44 335 either the RDP version 16 or SILVA 128 for sequence classification. Lastly, for QIIME 2 we  
45  
46 336 first dereplicated the query sequences using the *vsearch dereplicate-sequences* function and  
47  
48  
49 337 then assigned them to the Greengenes (13\_8) or SILVA 128 (99% identity clusters) databases  
50  
51 338 using the *feature-classifier classify-sklearn* function [18].  
52

53 339

54 340

55 341

56  
57  
58  
59  
60  
61  
62  
63  
64  
65



342 **Analysis and visualization**

1  
2 343 TSV and BIOM files were generated from the MAPseq and QIIME 2 outputs and combined  
3  
4  
5 344 with the output BIOM files created by QIIME and mothur (make.biom command).  
6  
7 345 Taxonomy names obtained from each individual reference database were normalized so that  
8  
9  
10 346 each genus and family would be assigned to the same lineage. Results were visualized and  
11  
12 347 analysed with the phyloseq [35] and vegan R packages. The recall rate (sensitivity) for each  
13  
14 348 tool and database was estimated as the percentage of sequences assigned to the expected taxa  
15  
16  
17 349 for each biome, while precision (specificity) was calculated as the fraction of sequences from  
18  
19 350 these predicted taxa out of all those from the taxa observed. Finally, the F-score was  
20  
21  
22 351 calculated as follows:

$$23$$
$$24 \quad F\text{-score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$
$$25 \quad 352$$
$$26$$
$$27$$
$$28 \quad 353$$

29  
30 354 Distance estimates were calculated with either the Bray-Curtis or Jaccard dissimilarity  
31  
32 355 indices after grouping the taxonomic lineages at the genus level. Principal coordinate analysis  
33  
34  
35 356 (PCoA) were performed with the Bray-Curtis distance method. Dissimilarity scores (DS) on  
36  
37 357 the relative abundance (rel.ab) of each expected genus were calculated as:

$$38$$
$$39$$
$$40 \quad 358$$
$$41$$
$$42 \quad DS = \frac{|\text{rel. ab. (Observed)} - \text{rel. ab. (Expected)}|}{\text{rel. ab. (Expected)}}$$
$$43 \quad 359$$
$$44$$
$$45$$
$$46 \quad 360$$
$$47$$

48 361 Memory usage and CPU time was estimated as the total amount required for the processing  
49  
50 362 and assignment of all combined sequences against the SILVA 128 database, following the  
51  
52  
53 363 protocols described above.

54  
55 364

56  
57  
58 365

59  
60  
61  
62  
63  
64  
65

366 **Availability of supporting source code and requirements**

367 Project name: Taxonomy benchmarking

368 Project home page: <https://github.com/Finn-Lab/Tax-Benchmarking>

369 Operating system: Platform independent

370 Programming languages: Python 2.7, R 3.4.1

371 Other requirements: BioPython module, R libraries (ggplot2, phyloseq, vegan, scales, grid,

372 ape, RColorBrewer, data.table)

373 License: MIT

374

375 **Availability of supporting data**

376 The datasets supporting the conclusions of this article are available in the GigaDB repository

377 [36].

378

379 **Declarations**

380 **List of abbreviations**

381 DS: Dissimilarity score

382 GG: Greengenes

383 PCoA: Principal coordinates analysis

384 rRNA : ribosomal rRNA

385

386 **Ethics approval and consent to participate**

387 Not applicable

388

389

1  
2  
3  
4  
5  
6  
7 **390 Consent for publication**

8  
9  
10 **391** Not applicable

11  
12 **392**

13  
14  
15 **393 Competing interests**

16  
17 **394** The authors declare that they have no competing interests

18  
19 **395**

20  
21  
22 **396 Funding**

23  
24 **397** European Molecular Biology Laboratory (EMBL); European Commission within the

25  
26 **398** Research Infrastructures Programme of Horizon 2020 [676559] (ELIXIR-EXCELERATE)

27  
28 **399**

29  
30  
31 **400 Authors' contributions**

32  
33 **401** AA, ALM, AT and RDF performed the analyses. AA, ALM and RDF conceived the study

34  
35 **402** and wrote the manuscript. All authors have read and approved the final manuscript.

36  
37 **403**

38  
39  
40 **404 Acknowledgements**

41  
42 **405** We thank João Matias Rodrigues for providing useful comments on the utility of MAPseq in

43  
44 **406** relation to the benchmarking.

45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

407 **Figure legends**

408 **Figure 1.** Level of recall at the genus level, represented as taxa relative abundances, obtained  
409 with each analysis pipeline for the three different biomes (human gut, ocean and soil). The  
410 number of genera correctly identified by each pipeline is indicated above the graph.

411  
412 **Figure 2.** (A) Recall, precision and F-score estimates at the genus level for each tool and  
413 database tested. (B) F-scores calculated for some of the most commonly tested sub-regions of  
414 the 16S rRNA gene: V1-V2, V3-V4, V4 and V4-V5.

415  
416 **Figure 3.** Computational cost of each taxonomy assignment tool, estimated as the total  
417 memory usage (A) and CPU time (B) required for the processing and classification of ~ 3  
418 million sequences against the SILVA 128 database.

419  
420 **Figure 4.** Dissimilarity scores (DS) calculated for each genus included in the simulated  
421 datasets. Lower (brighter) values indicate a closer prediction to the true composition of the  
422 original sample. The black outline indicates the overall best scoring analysis pipeline for each  
423 environment.

424  
425 **Figure 5.** Principal coordinates analysis (PCoA) between all samples analysed in relation to  
426 the true, expected dataset, using the Bray-Curtis distance method.

427  
428 **Figure S1.** Composition of the synthetic communities per selected environment. Samples  
429 A100 and B100 are randomly generated sets of 100 species, while A500 and B500 were  
430 simulated from 500 different species.

1  
2 431 **Figure S2.** Percentage of bacterial sequences retrieved from the Greengenes and SILVA  
3 432 databases with an *in silico* PCR targeting different 16S rRNA gene sub-regions.  
4

5 433  
6  
7 434 **Figure S3.** Tools and databases benchmarked in our study. We tested at least two databases  
8  
9 435 per software tool. The reference databases used were either readily supported by the specific  
10  
11 436 tool and/or recommended by their developers. SILVA was compared across all tools;  
12  
13 437 MAPseq was specifically assessed with the NCBI database, its default reference; mothur was  
14  
15 438 not paired with Greengenes due to its poor-quality alignment and was analysed with RDP  
16  
17 439 instead.  
18  
19

20  
21 440  
22  
23 441 **Figure S4.** Number of genera misassigned in each analysis pipeline and their overall relative  
24  
25 442 abundance. Names and abundance values of each misclassified taxon are included as  
26  
27 443 additional files (Additional files: Tables S2, S3 and S4).  
28  
29

30 444  
31  
32 445 **Figure S5.** Recall, precision and F-score estimates at the family level for each tool and  
33  
34 446 database tested.  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

447 **References**

- 1  
2  
3 448 1. Forbes JD, Van Domselaar G, Bernstein CN. The gut microbiota in immune-mediated  
4  
5 449 inflammatory diseases. *Front. Microbiol.* 2016.
- 6  
7 450 2. Duvallet C, Gibbons SM, Gurry T, Irizarry RA, Alm EJ. Meta-analysis of gut microbiome  
8  
9 451 studies identifies disease-specific and shared responses. *Nat. Commun.* Nature Publishing  
10  
11 Group; 2017;8:1784.
- 12 452  
13  
14 453 3. Turnbaugh PJ, Ley RE, Mahowald MA, Magrini V, Mardis ER, Gordon JI. An obesity-  
15  
16 associated gut microbiome with increased capacity for energy harvest. *Nature.* Nature  
17 454  
18 Publishing Group; 2006;444:1027–131.
- 19 455  
20  
21 456 4. Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, et al. A communal  
22  
23 catalogue reveals Earth’s multiscale microbial diversity. *Nature.* 2017;551:457–63.
- 24 457  
25  
26 458 5. Yilmaz P, Yarza P, Rapp JZ, Glöckner FO. Expanding the World of Marine Bacterial and  
27  
28 Archaeal Clades. *Front. Microbiol.* 2016;6:1524.
- 29 459  
30  
31 460 6. Fierer N. Embracing the unknown: disentangling the complexities of the soil microbiome.  
32  
33 *Nat. Rev. Microbiol.* Nature Publishing Group; 2017;15:579–90.
- 34 461  
35  
36 462 7. Mitchell AL, Scheremetjew M, Denise H, Potter S, Tarkowska A, Qureshi M, et al. EBI  
37  
38 Metagenomics in 2017: enriching the analysis of microbial communities, from sequence  
39 463  
40 reads to assemblies. *Nucleic Acids Res.* 2017;
- 41 464  
42  
43 465 8. Pace NR, Stahl DA, Lane DJ, Olsen GJ. The Analysis of Natural Microbial Populations by  
44  
45 Ribosomal RNA Sequences. *Adv. Microb. Ecol.* Springer, Boston, MA; 1986. p. 1–55.
- 46 466  
47  
48 467 9. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al. Introducing  
49  
50 mothur: open-source, platform-independent, community-supported software for describing  
51 468  
52 and comparing microbial communities. *Appl. Environ. Microbiol.* American Society for  
53 469  
54 Microbiology; 2009;75:7537–41.
- 55 470  
56  
57 471 10. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al.

60  
61  
62  
63  
64  
65

472 QIIME allows analysis of high-throughput community sequencing data. *Nat. Methods.*  
1  
2 473 *Nature Publishing Group; 2010;7:335–6.*  
3  
4  
5 474 11. McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, et al. An  
6  
7 475 improved Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses  
8  
9 476 of bacteria and archaea. *ISME J. Nature Publishing Group; 2012;6:610–8.*  
10  
11  
12 477 12. Federhen S. The NCBI Taxonomy database. *Nucleic Acids Res. Oxford University Press;*  
13  
14 478 *2012;40:D136-43.*  
15  
16  
17 479 13. Cole JR, Wang Q, Fish JA, Chai B, McGarrell DM, Sun Y, et al. Ribosomal Database  
18  
19 480 Project: data and tools for high throughput rRNA analysis. *Nucleic Acids Res. Oxford*  
20  
21 481 *University Press; 2014;42:D633–42.*  
22  
23  
24 482 14. Yilmaz P, Parfrey LW, Yarza P, Gerken J, Pruesse E, Quast C, et al. The SILVA and  
25  
26 483 “All-species Living Tree Project (LTP)” taxonomic frameworks. *Nucleic Acids Res. Oxford*  
27  
28 484 *University Press; 2014;42:D643-8.*  
29  
30  
31 485 15. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics.*  
32  
33 486 *Oxford University Press; 2010;26:2460–1.*  
34  
35  
36 487 16. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian classifier for rapid assignment  
37  
38 488 of rRNA sequences into the new bacterial taxonomy. *Appl. Environ. Microbiol. American*  
39  
40 489 *Society for Microbiology (ASM); 2007;73:5261–7.*  
41  
42  
43 490 17. Matias Rodrigues JF, Schmidt TSB, Tackmann J, von Mering C. MAPseq: highly  
44  
45 491 efficient k-mer search with confidence estimates, for rRNA sequence analysis.  
46  
47 492 *Bioinformatics. 2017;*  
48  
49  
50  
51 493 18. Bokulich NA, Kaehler BD, Rideout JR, Dillon M, Bolyen E, Knight R, et al. Optimizing  
52  
53 494 taxonomic classification of marker gene amplicon sequences. *PeerJ (preprint). 2018;*  
54  
55  
56 495 19. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, et al. Critical  
57  
58 496 *Assessment of Metagenome Interpretation—a benchmark of metagenomics software. Nat.*  
59  
60  
61  
62  
63  
64  
65

497 Methods. Nature Publishing Group; 2017;14:1063–71.

1  
2  
3 498 20. Golob JL, Margolis E, Hoffman NG, Fredricks DN. Evaluating the accuracy of amplicon-  
4  
5 499 based microbiome computational pipelines on simulated human gut microbial communities.  
6  
7 500 BMC Bioinformatics. BioMed Central; 2017;18:283.

8  
9  
10 501 21. Lindgreen S, Adair KL, Gardner PP. An evaluation of the accuracy and speed of  
11  
12 502 metagenome analysis tools. Sci. Rep. 2015;

13  
14 503 22. Fosso B, Santamaria M, Marzano M, Alonso-Aleman D, Valiente G, Donvito G, et al.  
15  
16 504 BioMaS: a modular pipeline for Bioinformatic analysis of Metagenomic AmpliconS. BMC  
17  
18 505 Bioinformatics. BioMed Central; 2015;16:203.

19  
20  
21 506 23. Yang B, Wang Y, Qian P-Y. Sensitivity and correlation of hypervariable regions in 16S  
22  
23 507 rRNA genes in phylogenetic analysis. BMC Bioinformatics. 2016;

24  
25  
26 508 24. Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD. Development of a Dual-  
27  
28 509 Index Sequencing Strategy and Curation Pipeline for Analyzing Amplicon Sequence Data on  
29  
30 510 the MiSeq Illumina Sequencing Platform. Appl. Environ. Microbiol. 2013;

31  
32  
33 511 25. Klindworth A, Pruesse E, Schweer T, Rg Peplies J, Quast C, Horn M, et al. Evaluation of  
34  
35 512 general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-  
36  
37 513 based diversity studies. Nucleic Acids Res. 2012;41:e1.

38  
39  
40  
41 514 26. Walker AW, Martin JC, Scott P, Parkhill J, Flint HJ, Scott KP. 16S rRNA gene-based  
42  
43 515 profiling of the human infant gut microbiota is strongly influenced by sample processing and  
44  
45 516 PCR primer choice. Microbiome. BioMed Central; 2015;3:26.

46  
47  
48 517 27. Comeau AM, Douglas GM, Langille MGI. Microbiome Helper: a Custom and  
49  
50 518 Streamlined Workflow for Microbiome Research. Eisen J, editor. mSystems. American  
51  
52 519 Society for Microbiology Journals; 2017;2:e00127-16.

53  
54  
55 520 28. Brooks JP, Edwards DJ, Harwich MD, Rivera MC, Fettweis JM, Serrano MG, et al. The  
56  
57 521 truth about metagenomics: quantifying and counteracting bias in 16S rRNA studies. BMC  
58  
59  
60  
61  
62  
63  
64  
65



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

522 Microbiol. 2015;15:66.

523 29. Schloss PD, Gevers D, Westcott SL. Reducing the effects of PCR amplification and  
524 sequencing Artifacts on 16s rRNA-based studies. PLoS One. 2011;6.

525 30. Quince C, Walker AW, Simpson JT, Loman NJ, Segata N. Shotgun metagenomics, from  
526 sampling to analysis. Nat. Biotechnol. 2017;35:833–44.

527 31. Shakya M, Quince C, Campbell JH, Yang ZK, Schadt CW, Podar M. Comparative  
528 metagenomic and rRNA microbial diversity characterization using archaeal and bacterial  
529 synthetic communities. Environ. Microbiol. NIH Public Access; 2013;15:1882–99.

530 32. Balvočiute M, Huson DH. SILVA, RDP, Greengenes, NCBI and OTT — how do these  
531 taxonomies compare? BMC Genomics. 2017;18.

532 33. Keegan KP, Glass EM, Meyer F. MG-RAST, a Metagenomics Service for Analysis of  
533 Microbial Community Structure and Function. Methods Mol. Biol. 2016. p. 207–33.

534 34. Huang W, Li L, Myers JR, Marth GT. ART: a next-generation sequencing read simulator.  
535 Bioinformatics. 2012;28:593–59410.

536 35. McMurdie PJ, Holmes S, Watson M. phyloseq: An R Package for Reproducible  
537 Interactive Analysis and Graphics of Microbiome Census Data. Watson M, editor. PLoS One.  
538 Public Library of Science; 2013;8:e61217.

539 36. Sneddon TP, Li P, Edmunds SC. GigaDB: announcing the GigaScience database.  
540 Gigascience. Oxford University Press; 2012;1:11.

16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Table 1. Global metrics averaged across the analyses of simulated samples from human gut, ocean and soil.**

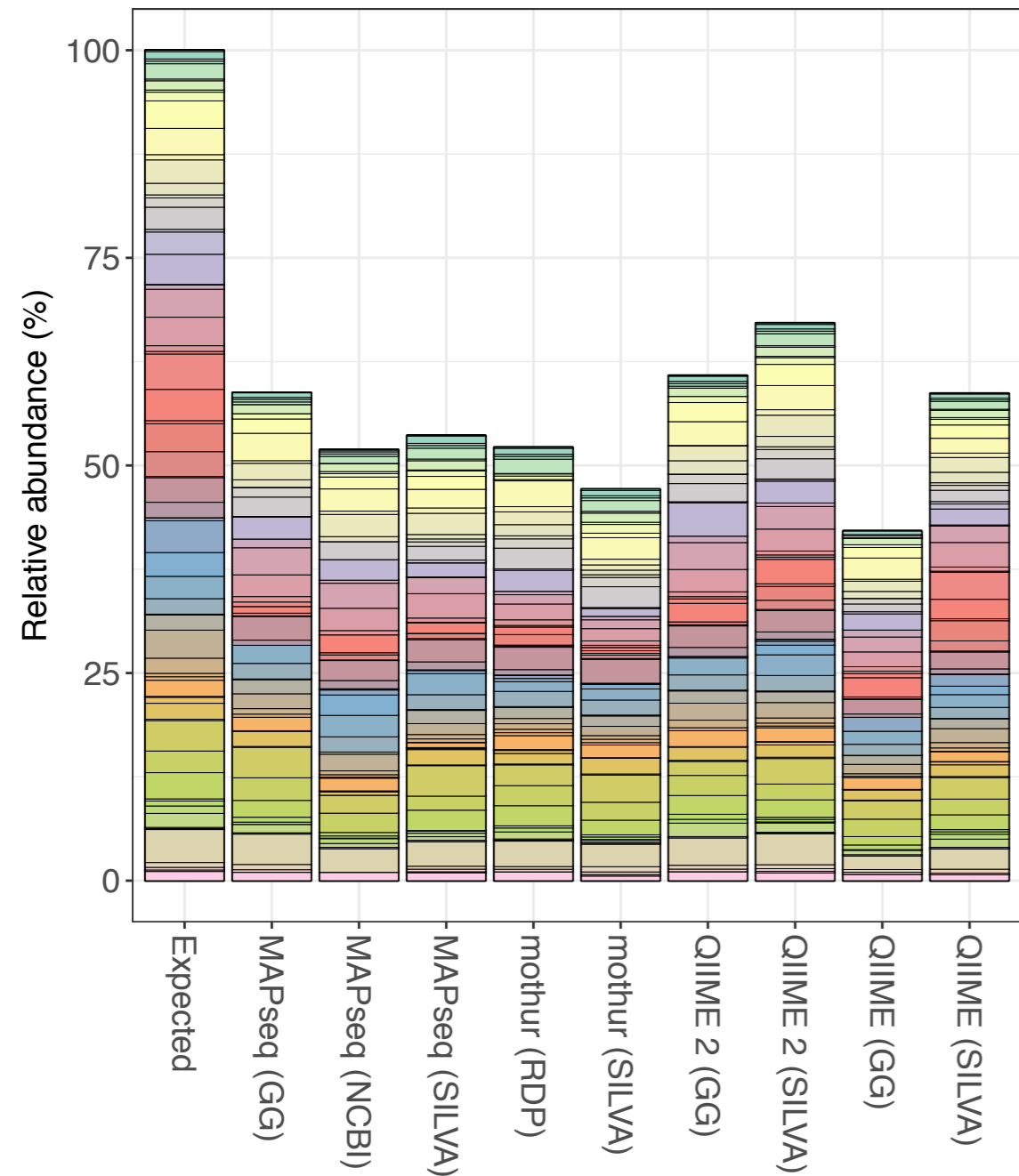
Software	Database	Family		Genus					
		Recall	Miscalled	Recall	Miscalled	Sub-region <sup>1</sup>	Mean DS	Bray-Curtis	Jaccard
MAPseq	Greengenes	87.8	2.4	58.6	2.5	V3-V4	0.435	0.284	0.441
MAPseq	NCBI	81.4	1.3	51.2	1.9	V3-V4	0.523	0.333	0.499
MAPseq	SILVA	66.9	<b>0.7</b>	46.2	<b>0.9</b>	V3-V4	0.484	0.375	0.543
mothur	RDP	84.8	3.2	49.2	4.8	V4-V5	0.430	0.364	0.532
mothur	SILVA	82.2	2.2	40.1	4.5	V4-V5	0.493	0.449	0.617
QIIME 2	Greengenes	92.5	1.8	69.1	3.4	V3-V4	0.372	0.211	<b>0.344</b>
QIIME 2	SILVA	<b>93.3</b>	1.9	<b>69.2</b>	4.4	V3-V4	<b>0.323</b>	<b>0.209</b>	0.345
QIIME	Greengenes	55.6	1.3	42.3	1.9	V4	0.619	0.420	0.592
QIIME	SILVA	61.8	2.1	53.3	6.4	V3-V4	0.473	0.343	0.508

Values in bold denote the best score.

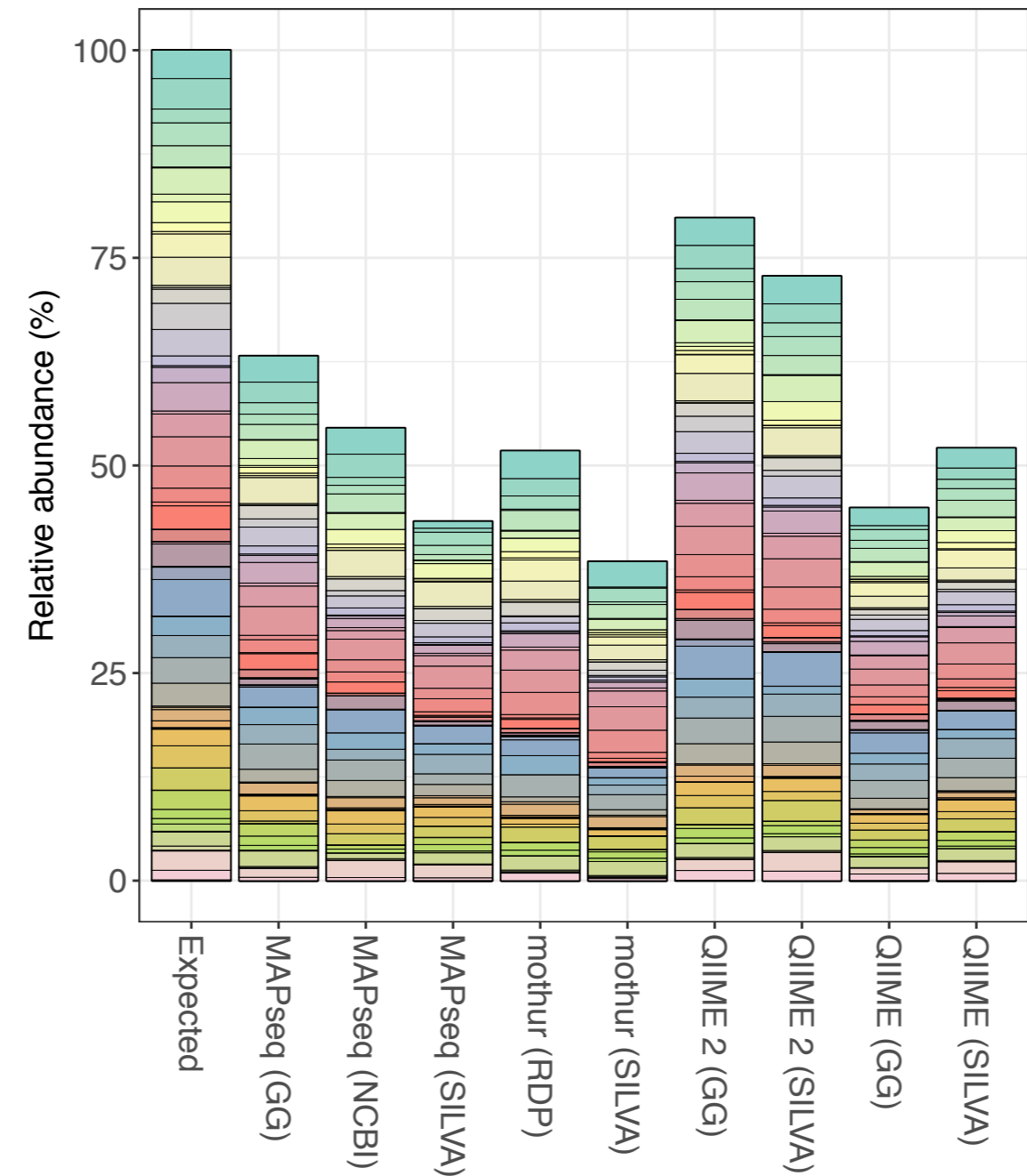
<sup>1</sup>Sub-region with the highest F-score, excluding V1-V2.

542

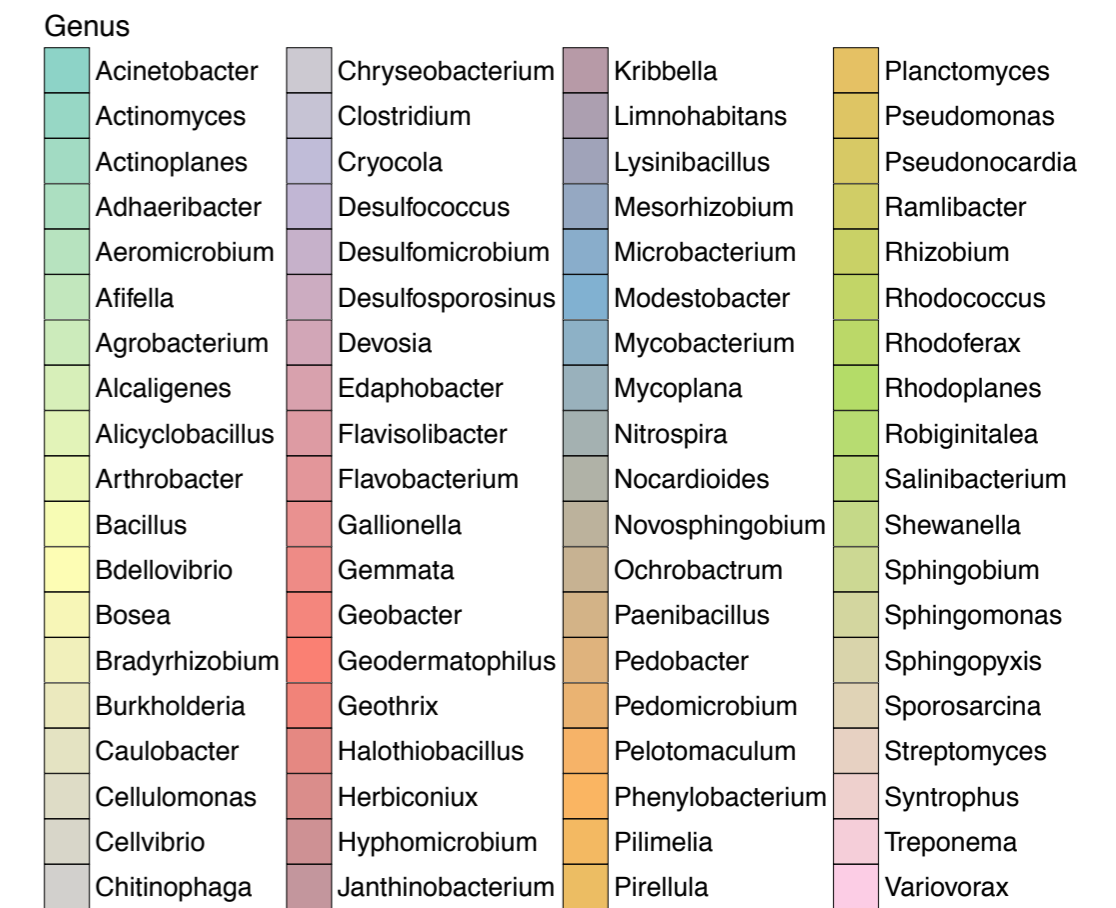
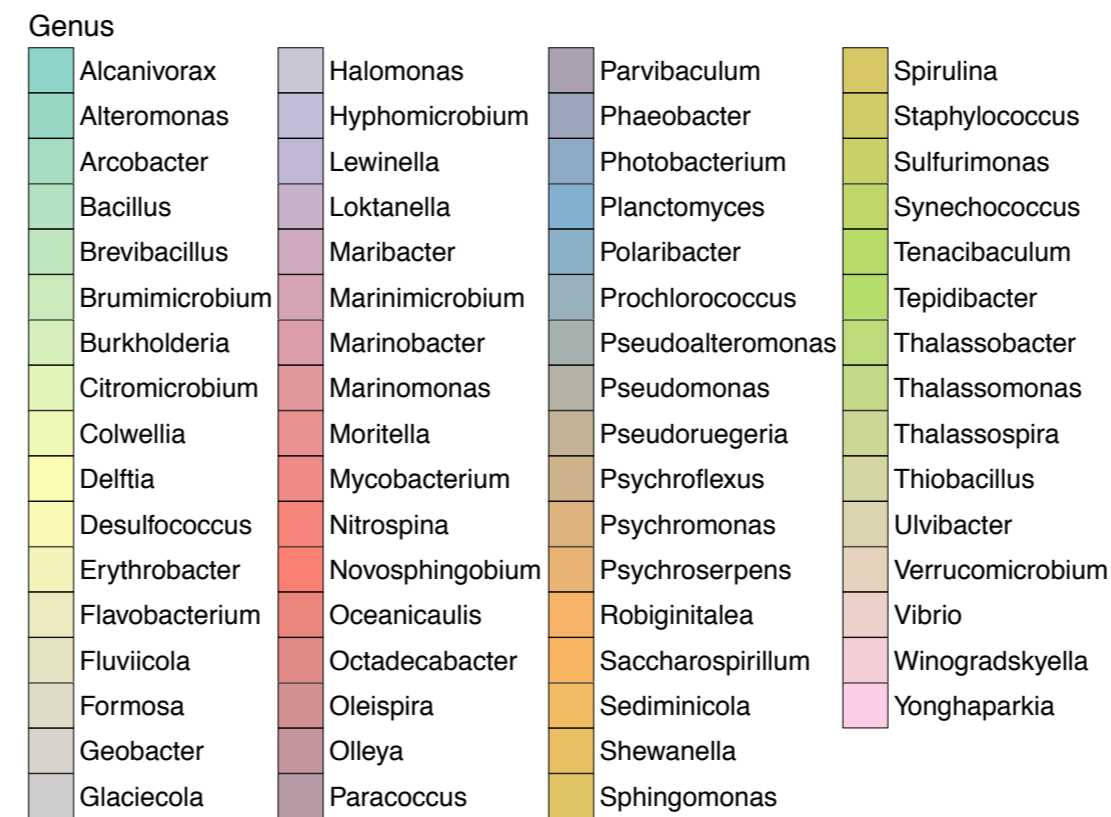
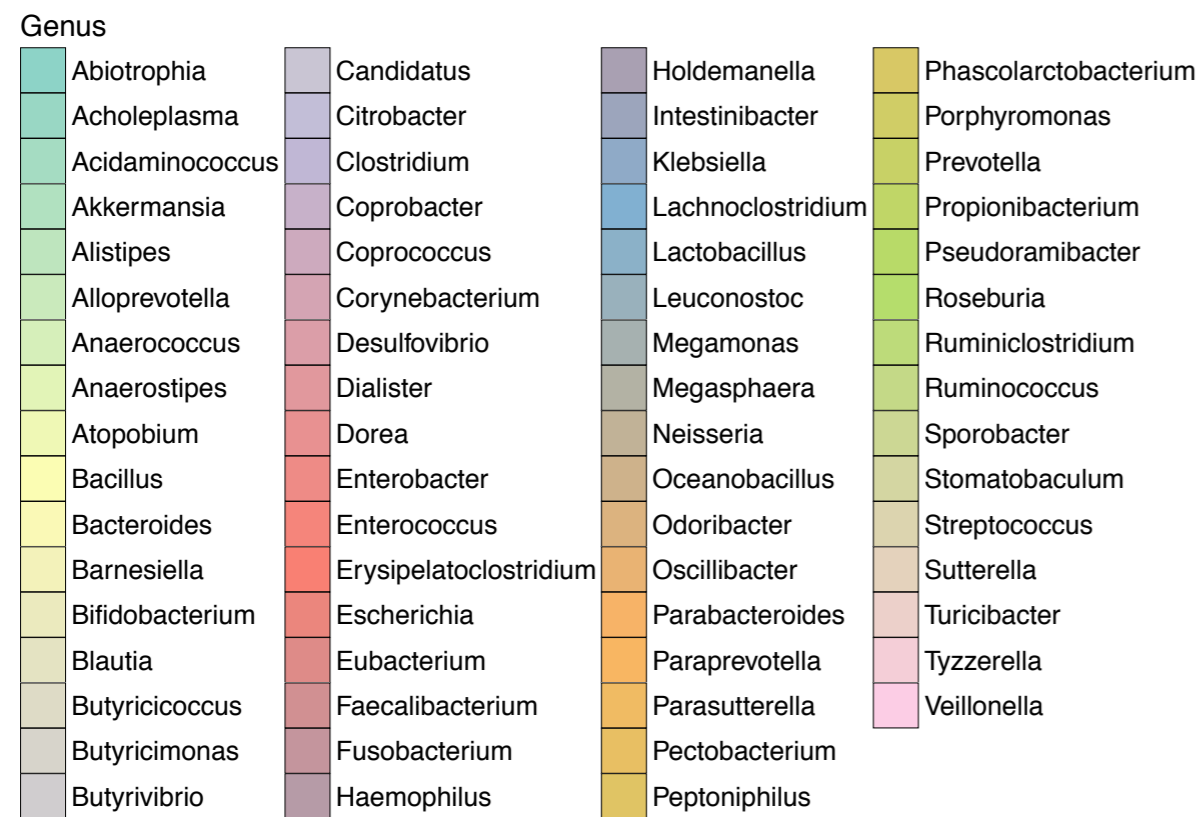
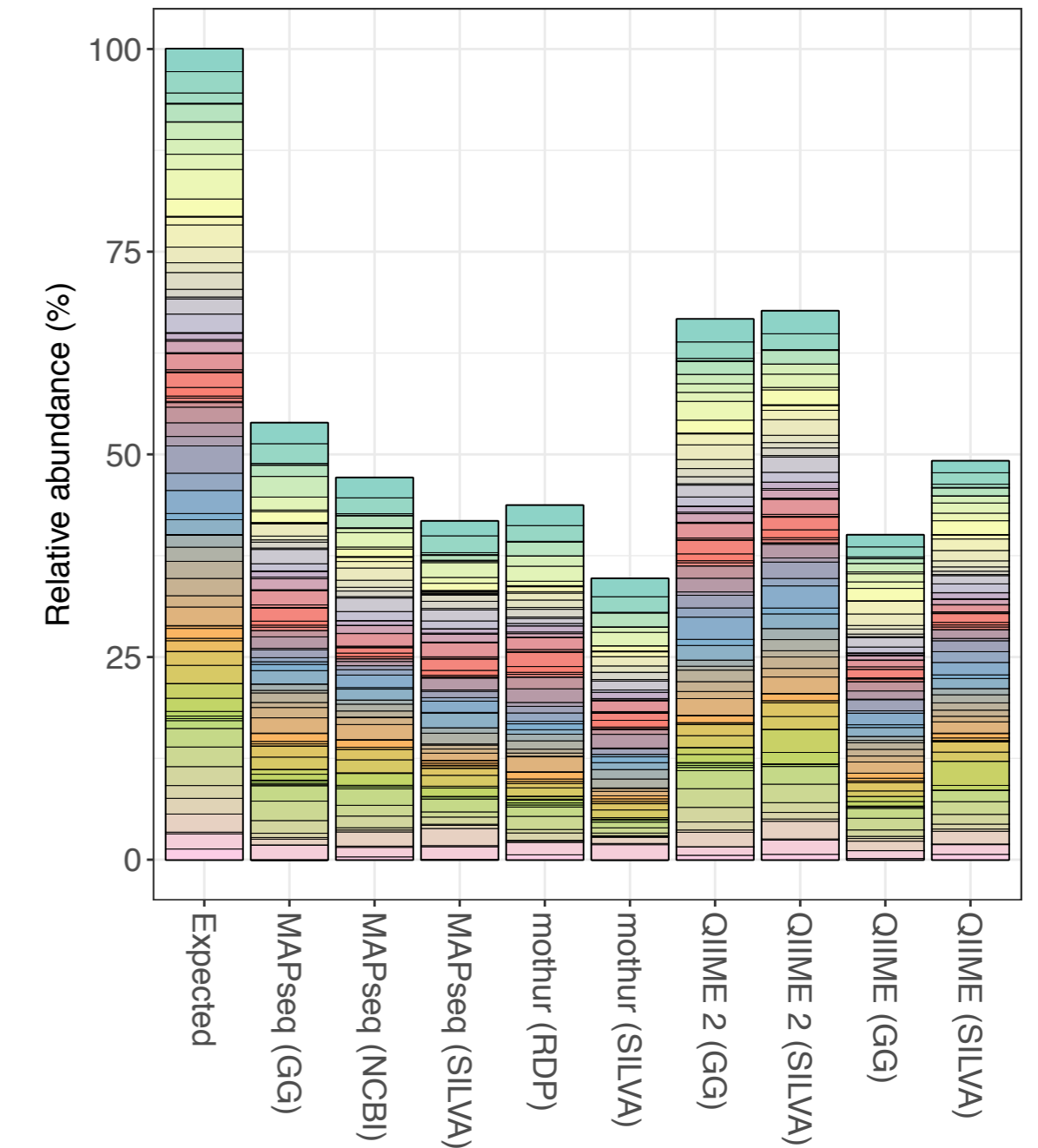
Human gut  
Genera 66 49 53 66 58 59 52 64 50 63

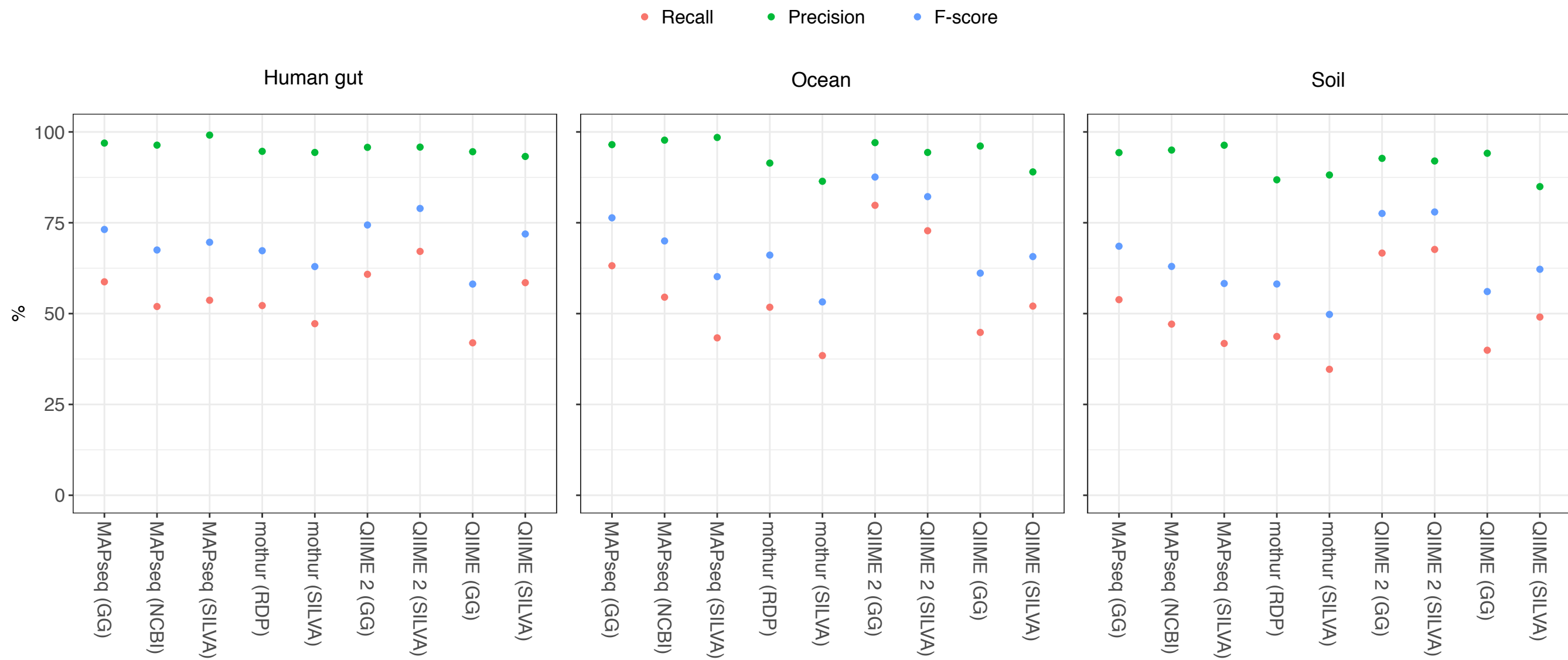
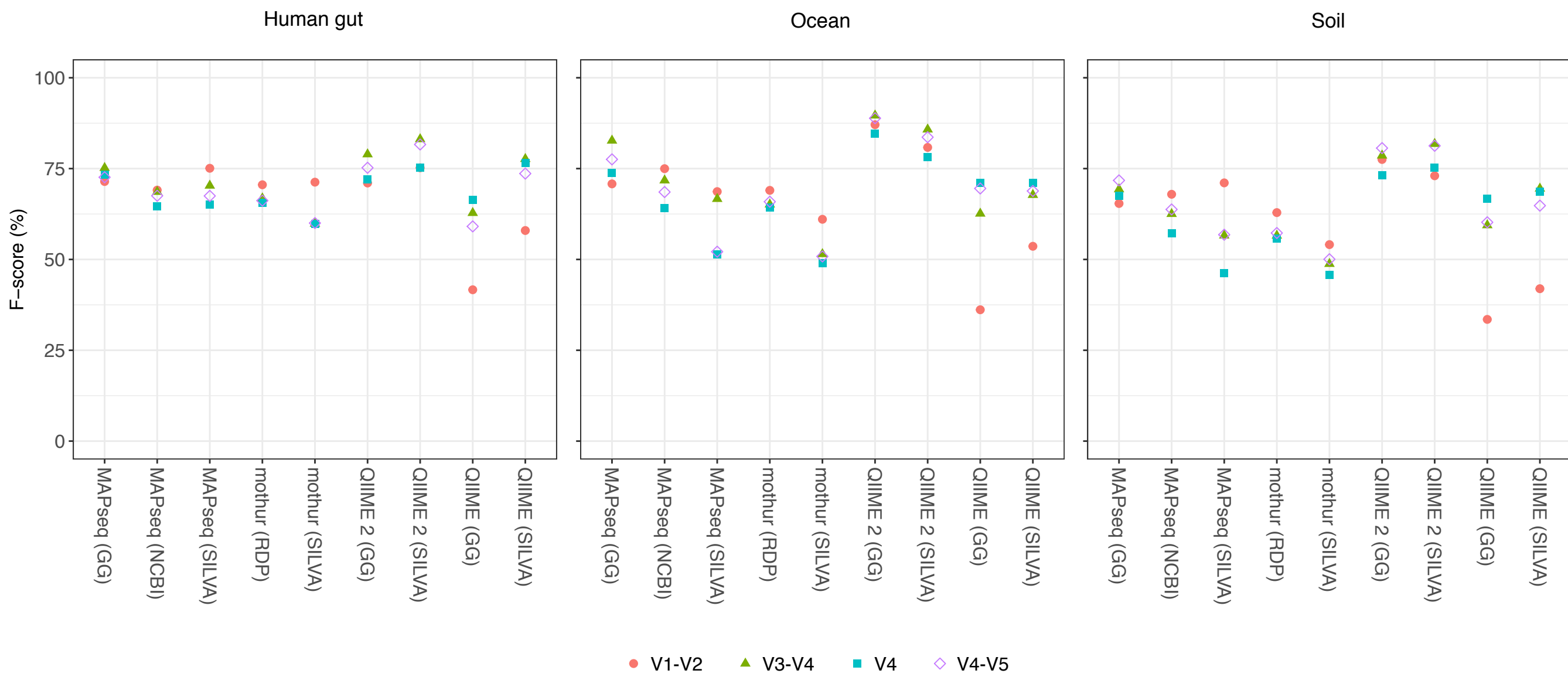


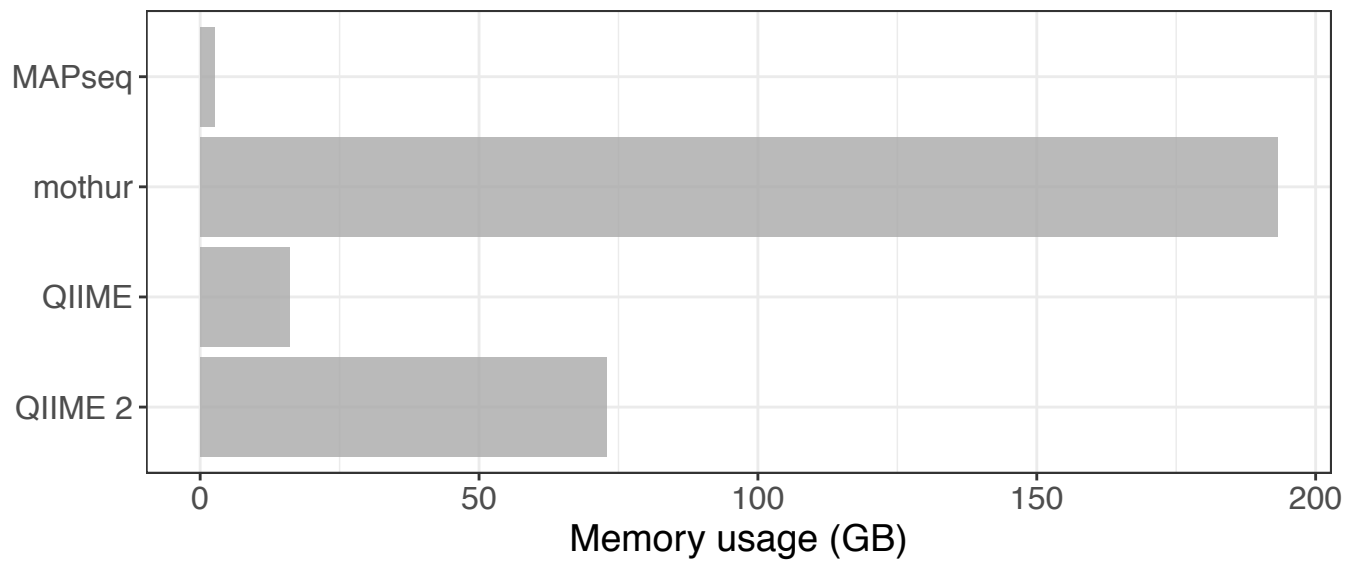
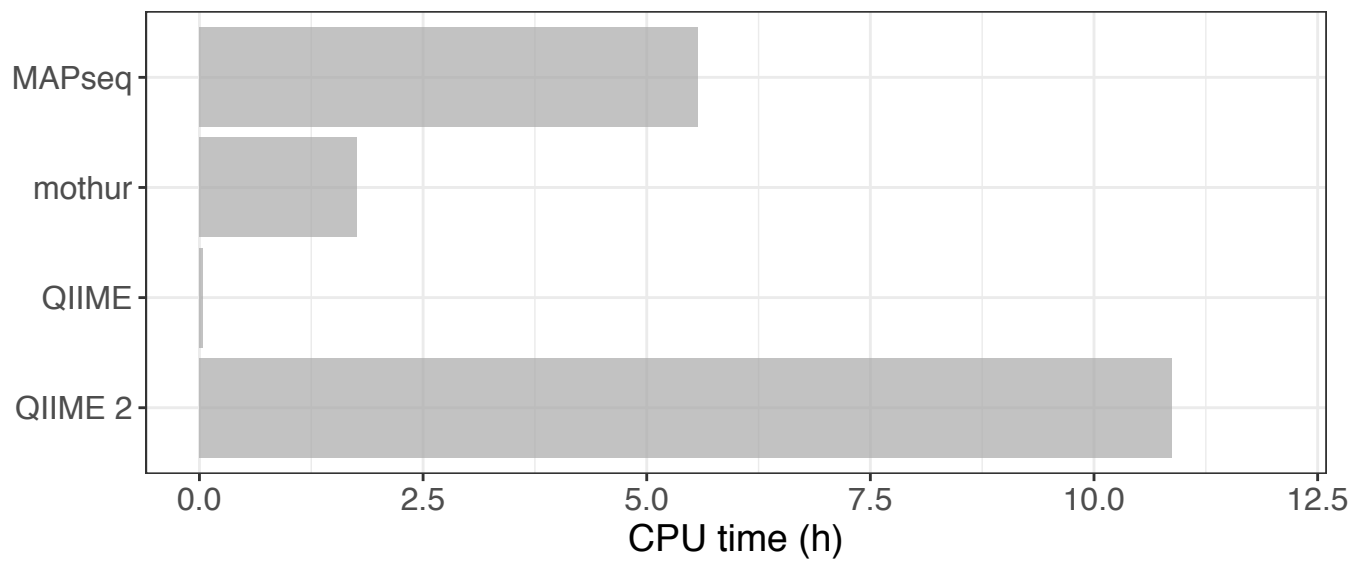
Ocean  
Genera 66 61 55 63 58 59 60 58 59 61

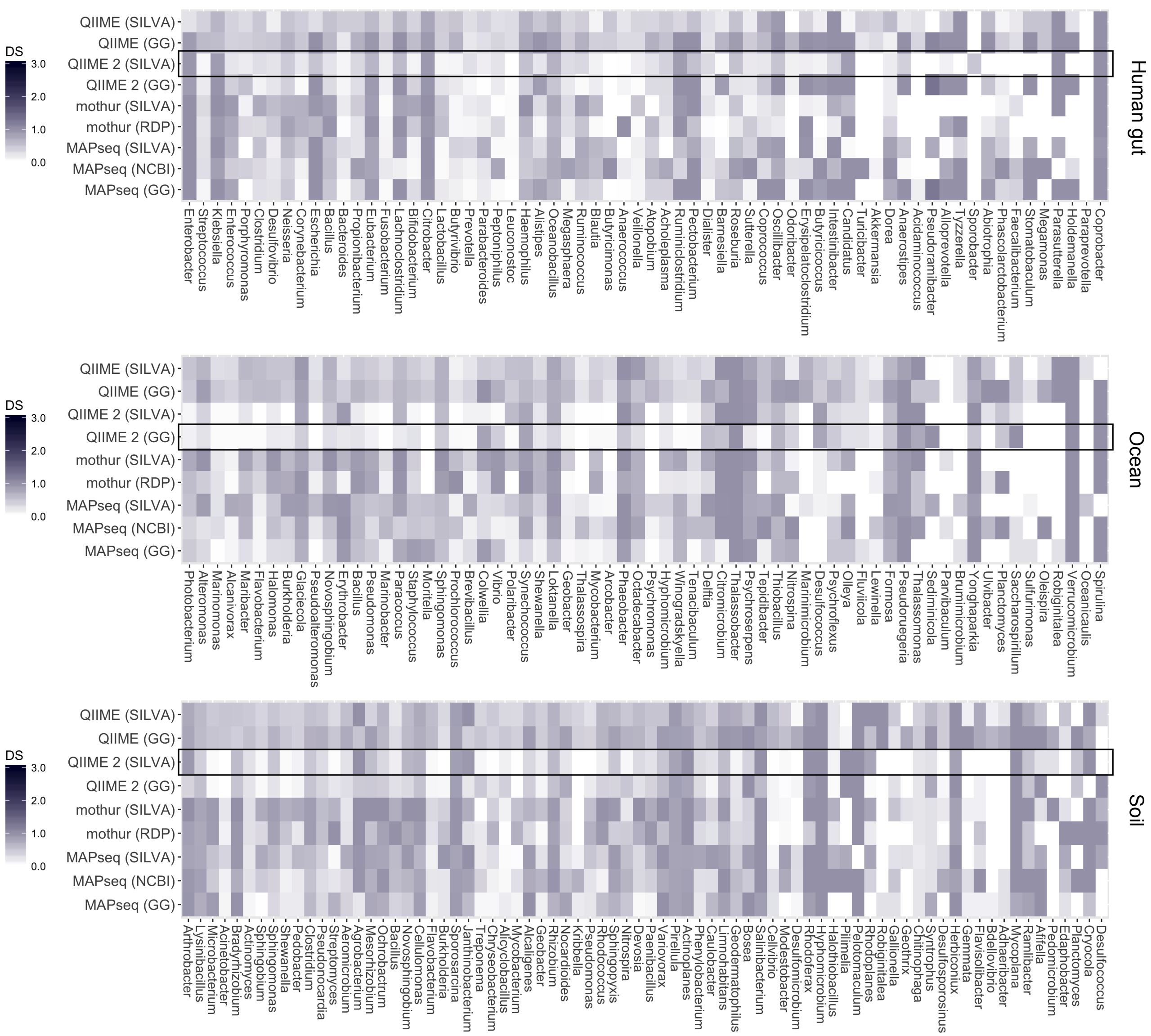


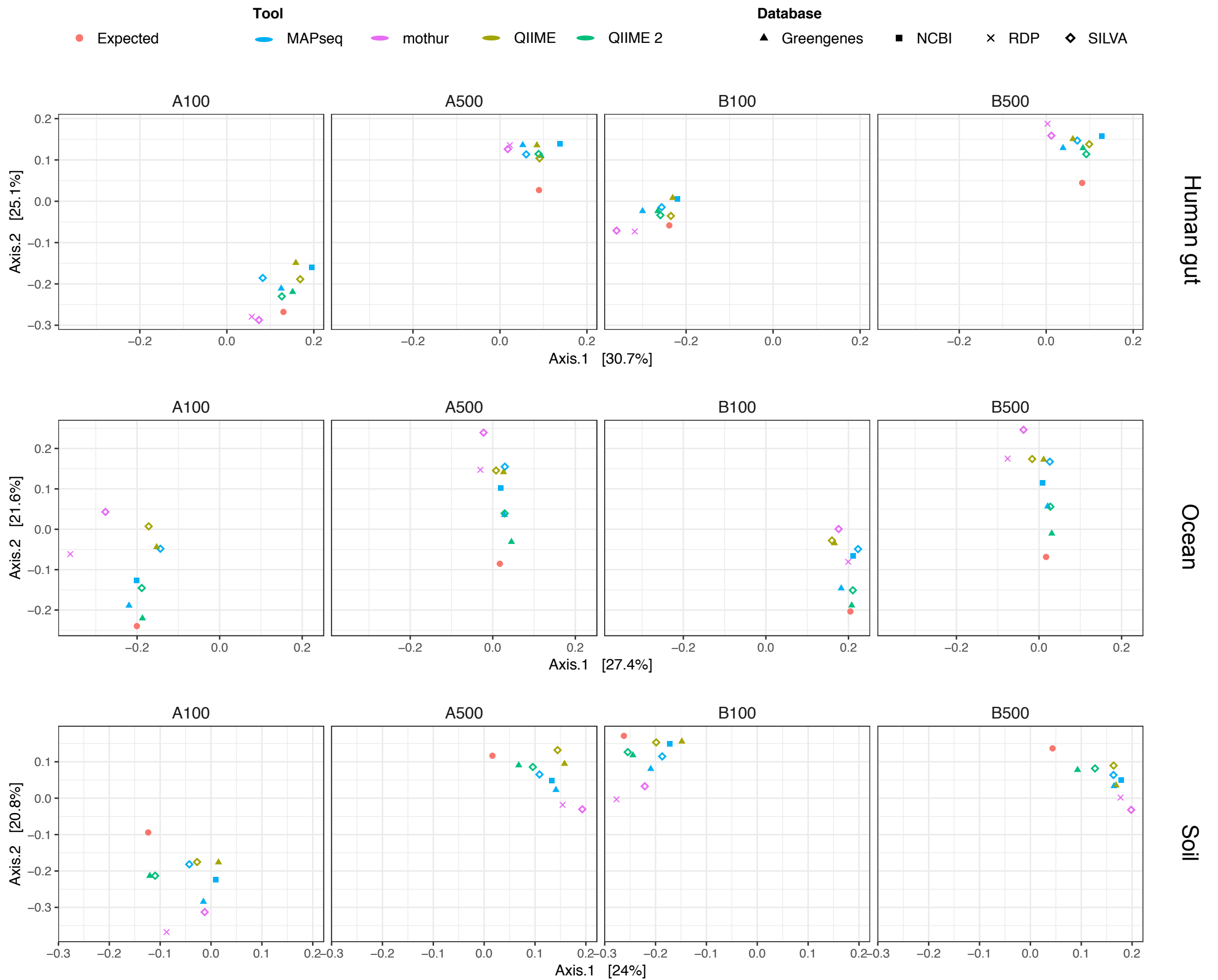
Soil  
Genera 76 71 64 71 66 62 70 67 65 67

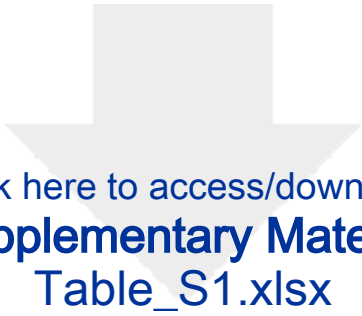


**A****B**


**A****B**



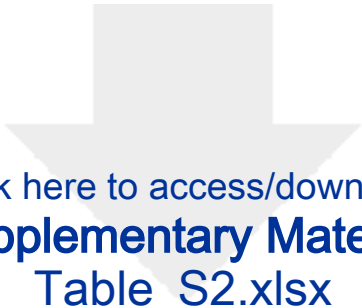





Click here to access/download  
**Supplementary Material**  
Table\_S1.xlsx

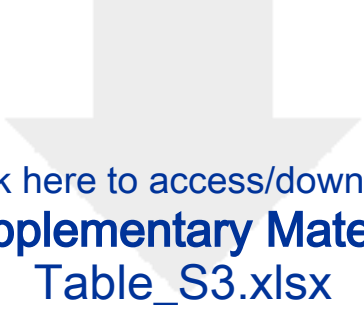







Click here to access/download  
**Supplementary Material**  
Table\_S2.xlsx







Click here to access/download  
**Supplementary Material**  
Table\_S3.xlsx




Click here to access/download  
**Supplementary Material**  
Table\_S4.xlsx




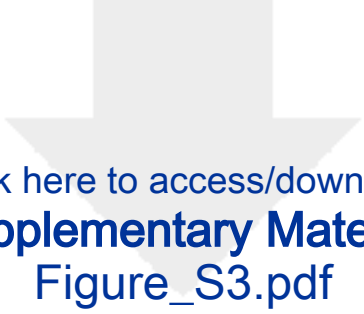
Click here to access/download  
**Supplementary Material**  
Figure\_S1.pdf







Click here to access/download  
**Supplementary Material**  
Figure\_S2.pdf





Click here to access/download  
**Supplementary Material**  
Figure\_S3.pdf





Click here to access/download  
**Supplementary Material**  
Figure\_S4.pdf



Click here to access/download  
**Supplementary Material**  
Figure\_S5.pdf