

OptiFit: an improved method for fitting amplicon sequences to existing OTUs

2021-11-08

Kelly L. Sovacool¹, Sarah L. Westcott², M. Brodie Mumphrey¹, Gabrielle A. Dotson¹,
Patrick D. Schloss²†

¹ Department of Computational Medicine and Bioinformatics, University of Michigan

² Department of Microbiology and Immunology, University of Michigan

† To whom correspondence should be addressed: pschloss@umich.edu

1 **Abstract**

2 Assigning amplicon sequences to operational taxonomic units (OTUs) is often an important
3 step in characterizing the composition of microbial communities across large datasets.
4 OptiClust, a *de novo* OTU clustering method, has been shown to produce higher quality
5 OTU assignments than other methods and at comparable or faster speeds. A notable
6 difference between *de novo* clustering and database-dependent reference clustering
7 methods is that OTU assignments from *de novo* methods may change when new
8 sequences are added to a dataset. However, in some cases one may wish to incorporate
9 new samples into a previously clustered dataset without performing clustering again on
10 all sequences, such as when comparing across datasets or deploying machine learning
11 models where OTUs are features. Existing reference-based clustering methods produce
12 consistent OTUs, but they only consider the similarity of each query sequence to a single
13 reference sequence in an OTU, thus resulting in OTU assignments that are significantly
14 worse than those generated by *de novo* methods. To provide an efficient and robust
15 method to fit amplicon sequence data to existing OTUs, we developed the OptiFit algorithm.
16 Inspired by OptiClust, OptiFit considers the similarity of all pairs of reference and query
17 sequences in an OTU to produce OTUs of the best possible quality. We tested OptiFit
18 using four microbiome datasets with two different strategies: by clustering to an external
19 reference database or by splitting the dataset into a reference and query set and clustering
20 the query sequences to the reference set after clustering it using OptiClust. The result is
21 an improved implementation of closed and open-reference clustering. OptiFit produces
22 OTUs of similar quality as OptiClust and at faster speeds when using the split dataset
23 strategy, although the OTU quality and processing speed depends on the database chosen
24 when using the external database strategy. OptiFit provides a suitable option for users who
25 require consistent OTU assignments at the same quality afforded by *de novo* clustering
26 methods.

27 **Importance**

28 Advancements in DNA sequencing technology have allowed researchers to affordably
29 generate millions of sequence reads from microorganisms in diverse environments.
30 Efficient and robust software tools are needed to assign microbial sequences into
31 taxonomic groups for characterization and comparison of communities. The OptiClust
32 algorithm produces high quality groups by comparing sequences to each other, but the
33 assignments can change when new sequences are added to a dataset, making it difficult
34 to compare different studies. Other approaches assign sequences to groups by comparing
35 them to sequences in a reference database to produce consistent assignments, but the
36 quality of the groups produced is reduced compared to OptiClust. We developed OptiFit, a
37 new reference-based algorithm that produces consistent yet high quality assignments like
38 OptiClust. OptiFit allows researchers to compare microbial communities across different
39 studies or add new data to existing studies without sacrificing the quality of the group
40 assignments.

41 Introduction

42 Amplicon sequencing is a mainstay of microbial ecology. Researchers can affordably
43 generate millions of sequences to characterize the composition of hundreds of samples
44 from microbial communities without the need for culturing. In many analysis pipelines,
45 16S rRNA gene sequences are assigned to operational taxonomic units (OTUs) to
46 facilitate comparison of taxonomic composition between communities to avoid the need
47 for taxonomic classification. A distance threshold of 3% (or sequence similarity of 97%) is
48 commonly used to cluster sequences into OTUs based on pairwise comparisons of the
49 sequences within the dataset. The method chosen for clustering affects the quality of OTU
50 assignments and thus may impact downstream analyses of community composition (1–3).

51 There are two main categories of OTU clustering algorithms: *de novo* and reference-based.
52 OptiClust is a *de novo* clustering algorithm which uses the distance score between all
53 pairs of sequences in the dataset to cluster them into OTUs by maximizing the Matthews
54 Correlation Coefficient (MCC) (1). This approach takes into account the distances between
55 all pairs of sequences when assigning query sequences to OTUs, in contrast to other *de*
56 *novo* methods such as the greedy clustering algorithms implemented in USEARCH and
57 VSEARCH (4, 5). In methods employing greedy clustering algorithms, only the distance
58 between each sequence and a representative centroid sequence in the OTU is considered
59 while clustering. As a result, distances between pairs of sequences in the same OTU
60 are frequently larger than the specified threshold, i.e. they are false positives. In contrast,
61 the OptiClust algorithm takes into account the distance between all pairs of sequences
62 when considering how to cluster sequences into OTUs and is thus less willing to take
63 on false positives. A limitation of *de novo* clustering is that different OTU assignments
64 will be produced when new sequences are added to a dataset, making it difficult to use
65 *de novo* clustering to compare OTUs between different studies. Furthermore, since *de*
66 *novo* clustering requires calculating and comparing distances between all sequences in a

67 dataset, the execution time can be slow and memory requirements can be prohibitive for
68 very large datasets. Reference clustering attempts to overcome the limitations of *de novo*
69 clustering methods by using a representative set of sequences from a database, with each
70 reference sequence seeding an OTU. Commonly, the Greengenes set of representative full
71 length sequences clustered at 97% similarity is used as the reference with VSEARCH (5–7).
72 Query sequences are then clustered into OTUs based on their similarity to the reference
73 sequences. Any query sequences that are not within the distance threshold to any of
74 the reference sequences are either thrown out (closed reference clustering) or clustered
75 *de novo* to create additional OTUs (open reference clustering). While reference-based
76 clustering is generally fast, it is limited by the diversity of the reference database. Novel
77 sequences in the sample will be lost in closed reference mode if they are not represented
78 by a similar sequence in the database. Previous studies found that the OptiClust *de novo*
79 clustering algorithm created the highest quality OTU assignments of all clustering methods
80 (1).

81 To overcome the limitations of current reference-based and *de novo* clustering algorithms
82 while maintaining OTU quality, we developed OptiFit, a reference-based clustering
83 algorithm. While other tools represent reference OTUs with a single sequence, OptiFit
84 uses multiple sequences in existing OTUs as the reference and fits new sequences to
85 those reference OTUs. In contrast to other tools, OptiFit considers all pairwise distance
86 scores between reference and query sequences when assigning sequences to OTUs
87 in order to produce OTUs of the highest possible quality. Here, we tested the OptiFit
88 algorithm with the reference as a public database (e.g. Greengenes) or *de novo* OTUs
89 generated using a reference set from the full dataset and compared the performance to
90 existing tools. To evaluate the OptiFit algorithm and compare to existing methods, we used
91 four published datasets isolated from soil (8), marine (9), mouse gut (10), and human gut
92 (11) samples. OptiFit is available within the mothur software program.

93 **Results**

94 **The OptiFit algorithm**

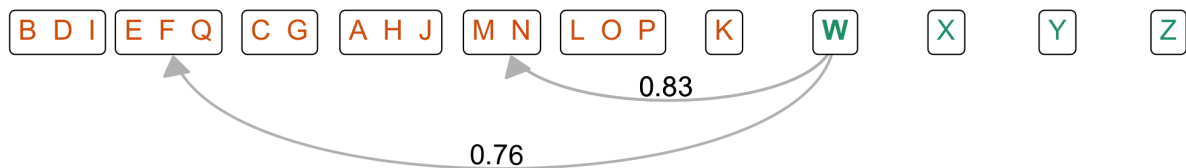
95 OptiFit leverages the method employed by OptiClust of iteratively assigning sequences
96 to OTUs to produce the highest quality OTUs possible, and extends this method for
97 reference-based clustering. OptiClust first seeds each sequence into its own OTU as a
98 singleton. Then for each sequence, OptiClust considers whether the sequence should
99 move to a different OTU or remain in its current OTU, choosing the option that results
100 in a better Matthews correlation coefficient (MCC) (1). The MCC uses all values from
101 a confusion matrix and ranges from negative one to one, with a score of one occurring
102 when all sequence pairs are true positives and true negatives and a score of negative one
103 occurring when all pairs are false positives and false negatives. Sequence pairs that are
104 similar to each other (i.e. within the distance threshold) are counted as true positives if
105 they are clustered into the same OTU, and false negatives if they are not in the the same
106 OTU. Sequence pairs that are not similar to each other are true negatives if they are not
107 clustered into the same OTU, and false positives if they are not in the same OTU. OptiClust
108 iterations continue until the MCC stabilizes or until a maximum number of iterations is
109 reached. This process produces *de novo* OTU assignments with the most optimal MCC
110 given the input sequences.

111 OptiFit begins where OptiClust ends, starting with a list of reference OTUs and their
112 sequences, a list of query sequences to cluster to the reference OTUs, and the sequence
113 pairs that are within the distance threshold (e.g. 0.03) (Figure 1). Initially, all query
114 sequences are placed into separate OTUs. Then, the algorithm iteratively reassigns the
115 query sequences to the reference OTUs to optimize the MCC. Alternatively, a sequence
116 will remain unassigned if the MCC value is maximized when the sequence is a singleton
117 rather than clustered into a reference OTU. All query and reference sequence pairs are
118 considered when calculating the MCC. This process is repeated until the MCC changes by

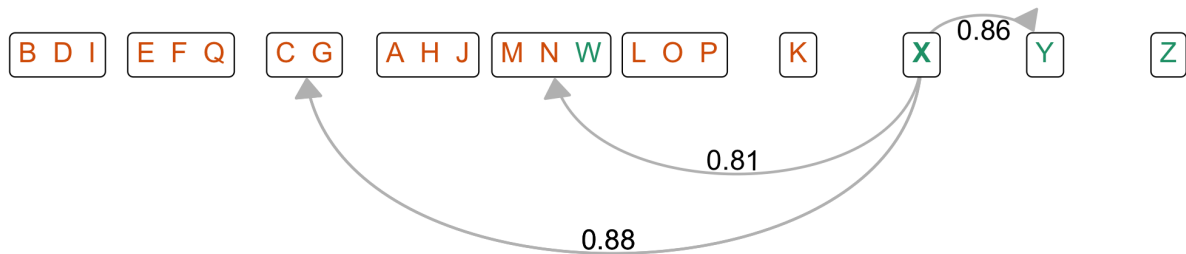
0. List of sequence pairs within the distance threshold

	D	F	G	H	I	I	J	J	N	O	P	P	P	Q	Q	W	W	W	X	X	X	X	Y
	B	E	C	A	B	D	A	H	M	L	K	L	O	E	F	F	M	N	C	G	N	Y	C
% distance	1.7	1.4	2.9	2.7	1.7	1.4	1.0	1.6	1.6	2.6	1.5	2.2	2.4	1.8	1.2	2.8	1.0	1.4	2.1	2.7	1.0	2.1	1.4

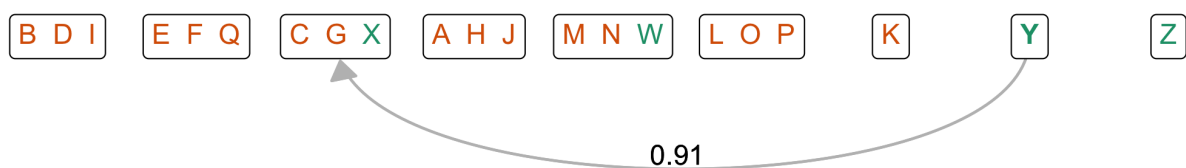
1. MCC = 0.78



2. MCC = 0.83



3. MCC = 0.88



4. MCC = 0.91



Figure 1: The OptiFit Algorithm. Here we present a toy example of the OptiFit algorithm fitting query sequences to existing OTUs, given the list of all sequence pairs that are within the distance threshold (here 3% is used). The goal of OptiFit is to assign the query sequences W through Z (colored **green**) to the reference OTUs created by clustering Sequences A through Q (colored **orange**) which were previously clustered *de novo* with OptiClust (see the OptiClust supplemental text (1)). Initially, OptiFit places each query sequence in its own OTU. Then, for each query sequence (**bolded**), OptiFit determines what the new MCC score would be if that sequence were moved to one of the OTUs containing at least one other similar sequence. The sequence is then moved to the OTU which would result in the best MCC score. OptiFit stops iterating over sequences once the MCC score stabilizes (in this example; only one iteration over each sequence is needed).

119 no more than 0.0001 (default) or until a maximum number of iterations is reached (default:
120 100). In the closed reference mode, any query sequences that cannot be clustered into
121 reference OTUs are discarded, and the results only contain OTUs that exist in the original
122 reference. In the open reference mode, unassigned query sequences are clustered *de*
123 *novo* using OptiClust to generate new OTUs. The final MCC is reported with the best
124 OTU assignments. There are two strategies for generating OTUs with OptiFit: 1) cluster
125 the query sequences to reference OTUs generated by *de novo* clustering an independent
126 database, or 2) split the dataset into a reference and query fraction, cluster the reference
127 sequences *de novo*, then cluster the query sequences to the reference OTUs.

128 **Reference clustering with public databases**

129 To test how OptiFit performs for reference-based clustering, we clustered each dataset to
130 three databases of reference OTUs: the Greengenes database, the SILVA non-redundant
131 database, and the Ribosomal Database Project (RDP) (6, 12, 13). Reference OTUs for
132 each database were created by performing *de novo* clustering with OptiClust at a distance
133 threshold of 3% using the V4 region of each sequence (see Figure 2). After trimming to
134 the V4 region, the databases contained 174,979, 16,192, and 173,648 unique sequences
135 and produced *de novo* MCC scores of 0.72, 0.74, and 0.73 for Greengenes, RDP, and
136 SILVA, respectively. Clustering sequences to Greengenes and SILVA in closed reference
137 mode performed similarly, with median MCC scores of 0.85 and 0.77 respectively, while
138 the median MCC was 0.35 when clustering to RDP (Figure 3). For comparison, clustering
139 datasets with OptiClust produced an average MCC score of 0.87. This gap in OTU quality
140 mostly disappeared when clustering in open reference mode, which produced median
141 MCCs of 0.86 with Greengenes, 0.85 with SILVA, and 0.86 with the RDP. Thus, open
142 reference OptiFit produced OTUs of very similar quality as *de novo* clustering, and closed
143 reference OptiFit followed closely behind as long as a suitable reference database was
144 chosen.

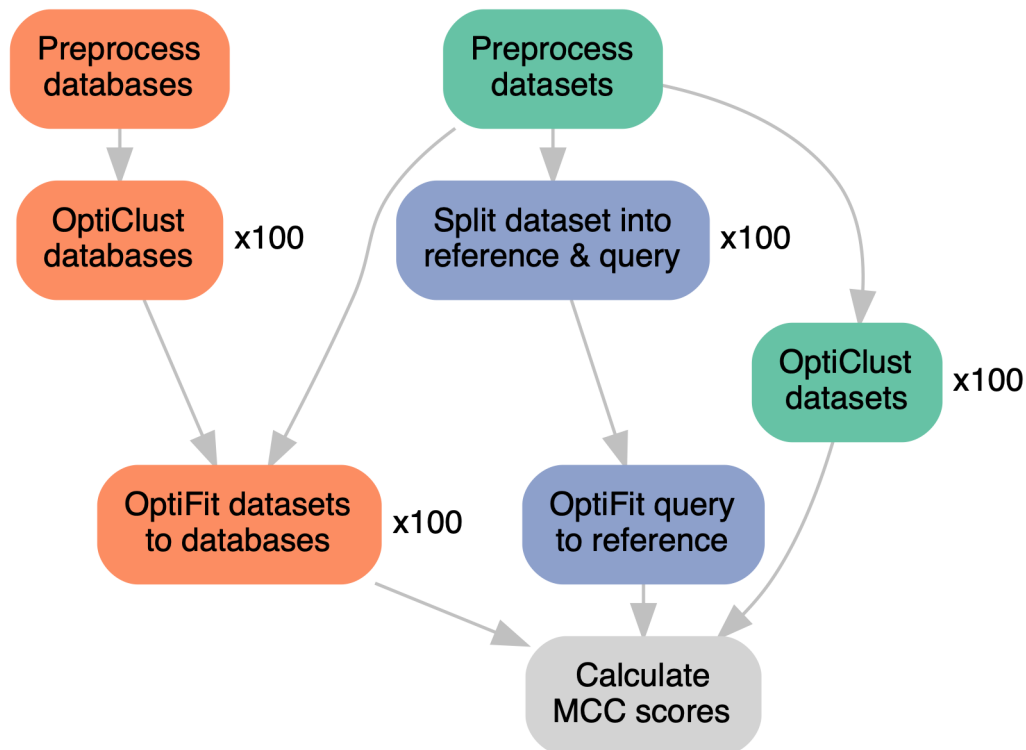


Figure 2: The Analysis Workflow. Reference sequences from Greengenes, the RDP, and SILVA were downloaded, preprocessed with mothur by trimming to the V4 region, and clustered *de novo* with OptiClust for 100 repetitions. Datasets from human, marine, mouse, and soil microbiomes were downloaded, preprocessed with mothur by aligning to the SILVA V4 reference alignment, then clustered *de novo* with OptiClust for 100 repetitions. Individual datasets were fit to reference databases with OptiFit; OptiFit was repeated 100 times for each dataset and database combination. Datasets were also randomly split into a reference and query fraction, and the query sequences were fit to the reference sequences with OptiFit for 100 repetitions. The final MCC score was reported for all OptiClust and OptiFit repetitions.

145 Since closed reference clustering does not cluster query sequences that could not be
146 clustered into reference OTUs, an additional measure of clustering performance to consider
147 is the fraction of query sequences that were able to be clustered. On average, more
148 sequences were clustered with Greengenes as the reference (59.1%) than with SILVA
149 (50.0%) or with the RDP (9.8%) (Figure 3). This mirrored the result reported above that
150 Greengenes produced better OTUs in terms of MCC score than either SILVA or RDP. Note
151 that *de novo* and open reference clustering methods always cluster 100% of sequences
152 into OTUs. The database chosen affects the final closed reference OTU assignments
153 considerably in terms of both MCC score and fraction of query sequences that could be
154 clustered into the reference OTUs.

155 Despite the drawbacks, closed reference methods have been used when fast execution
156 speed is required, such as when using very large datasets (14). To compare performance
157 in terms of speed, we repeated each OptiFit and OptiClust run 100 times and measured
158 the execution time. Across all dataset and database combinations, closed reference OptiFit
159 outperformed both OptiClust and open reference OptiFit (Figure 3). For example, with
160 the human dataset fit to SILVA reference OTUs, the average run times in seconds were
161 406.8 for closed reference OptiFit, 455.3 for *de novo* clustering the dataset, and 559.4 for
162 open reference OptiFit. Thus, the OptiFit algorithm continues the precedent that closed
163 reference clustering sacrifices OTU quality for execution speed.

164 To compare to the reference clustering methods used by QIIME2, we clustered each
165 dataset with VSEARCH against the Greengenes database of OTUs previously clustered
166 at 97% sequence similarity. Each reference OTU from the Greengenes 97% database
167 contains one reference sequence, and VSEARCH maps sequences to the reference
168 based on each individual query sequence's similarity to the single reference sequence.
169 In contrast, OptiFit accepts reference OTUs which each may contain multiple sequences,
170 and the sequence similarity between all query and reference sequences is considered

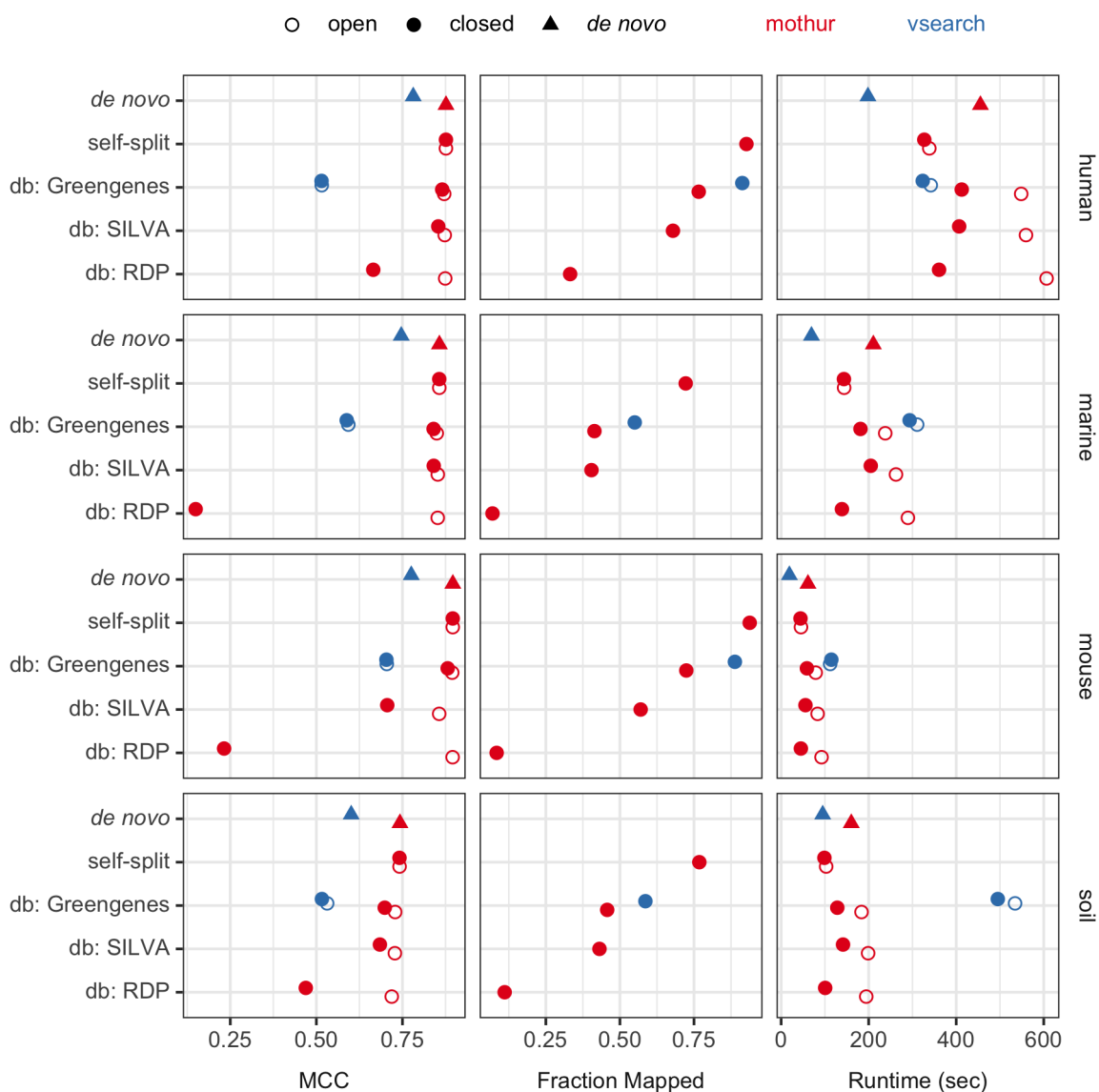


Figure 3: Benchmarking Results. The median MCC score, fraction of query sequences that mapped in closed-reference clustering, and runtime in seconds from repeating each clustering method 100 times. Each dataset underwent *de novo* clustering using OptiClust or reference-based clustering using OptiFit with one of two strategies; splitting the dataset and fitting 50% the sequences to the other 50%, or fitting the dataset to a reference database (Greengenes, SILVA, or RDP). Reference-based clustering was repeated with open and closed mode. For additional comparison, VSEARCH was used for *de novo* and reference-based clustering against the Greengenes database.

171 when assigning sequences to OTUs. In closed reference mode, OptiFit produced 27.2%
172 higher quality OTUs than VSEARCH, but VSEARCH was able to cluster 24.8% more query
173 sequences than OptiFit to the Greengenes reference database (Figure 3). This is because
174 VSEARCH only considers the distances between each query sequence to the single
175 reference sequence, while OptiFit considers the distances between all pairs of reference
176 and query sequences in an OTU. When open reference clustering, OptiFit produced higher
177 quality OTUs than VSEARCH against the Greengenes database, with median MCC scores
178 of 0.86 and 0.56, respectively. In terms of run time, OptiFit outperformed VSEARCH in
179 both closed and open reference mode by 54.6% and 49.5% on average, respectively. Thus,
180 the more stringent OTU definition employed by OptiFit, which prefers the query sequence
181 to be similar to all other sequences in the OTU rather than to only one sequence, resulted
182 in fewer sequences being clustered to reference OTUs than when using VSEARCH, but
183 caused OptiFit to outperform VSEARCH in terms of both OTU quality and execution time.

184 **Reference clustering with split datasets**

185 When performing reference clustering against public databases, the database chosen
186 greatly affects the quality of OTUs produced. OTU quality may be poor when the reference
187 database consists of sequences that are too unrelated to the samples of interest, such as
188 when samples contain novel populations. While *de novo* clustering overcomes the quality
189 limitations of reference clustering to databases, OTU assignments are not consistent when
190 new sequences are added. Researchers may wish to cluster new sequences to existing
191 OTUs or to compare OTUs across studies. To determine how well OptiFit performs for
192 clustering new sequences to existing OTUs, we employed a split dataset strategy, where
193 each dataset was randomly split into a reference fraction and a query fraction. Reference
194 sequences were clustered *de novo* with OptiClust, then query sequences were clustered
195 to the *de novo* OTUs with OptiFit.

196 First, we tested whether OptiFit performed as well as *de novo* clustering when using the
197 split dataset strategy with half of the sequences selected for the reference by a simple
198 random sample (a 50% split) (Figure 3; self-split). OTU quality was similar to that from
199 OptiClust regardless of mode (0.029% difference in median MCC). In closed reference
200 mode, OptiFit was able to cluster 84.8% of query sequences to reference OTUs with
201 the split strategy, a great improvement over the average 59.1% of sequences clustered
202 to the Greengenes database. In terms of run time, closed and open reference OptiFit
203 performed faster than OptiClust on whole datasets by 34.7% and 33.5%, respectively. The
204 split dataset strategy also performed 13.5% faster than the database strategy in closed
205 reference mode and 43.5% faster in open reference mode. Thus, reference clustering with
206 the split dataset strategy creates as high quality OTUs as *de novo* clustering yet at a faster
207 run time, and fits far more query sequences than the database strategy.

208 While we initially tested this strategy using a 50% split of the data into reference and query
209 fractions, we next investigated whether there was an optimal reference fraction size. To
210 identify the best reference size, reference sets with 10% to 90% of the sequences were
211 created, with the remaining sequences used for the query (Figure 4). OTU quality was
212 remarkably consistent across reference fraction sizes. For example, splitting the human
213 dataset 100 times yielded a coefficient of variation (i.e. the standard deviation divided by
214 the mean) of 0.00022 for the MCC score across all fractions. Run time generally decreased
215 as the reference fraction increased; for the human dataset, the median run time was
216 364.1 seconds with 10% of sequences in the reference and 291.3 seconds with 90% of
217 sequences in the reference. In closed reference mode, the fraction of sequences that
218 mapped increased as the reference size increased; for the human dataset, the median
219 fraction mapped was 0.85 with 10% of sequences in the reference and 0.95 with 90% of
220 sequences in the reference. These trends held for the other datasets as well. Thus, the
221 reference fraction did not affect OTU quality in terms of MCC score, but did affect the run
222 time and the fraction of sequences that mapped during the closed reference clustering.

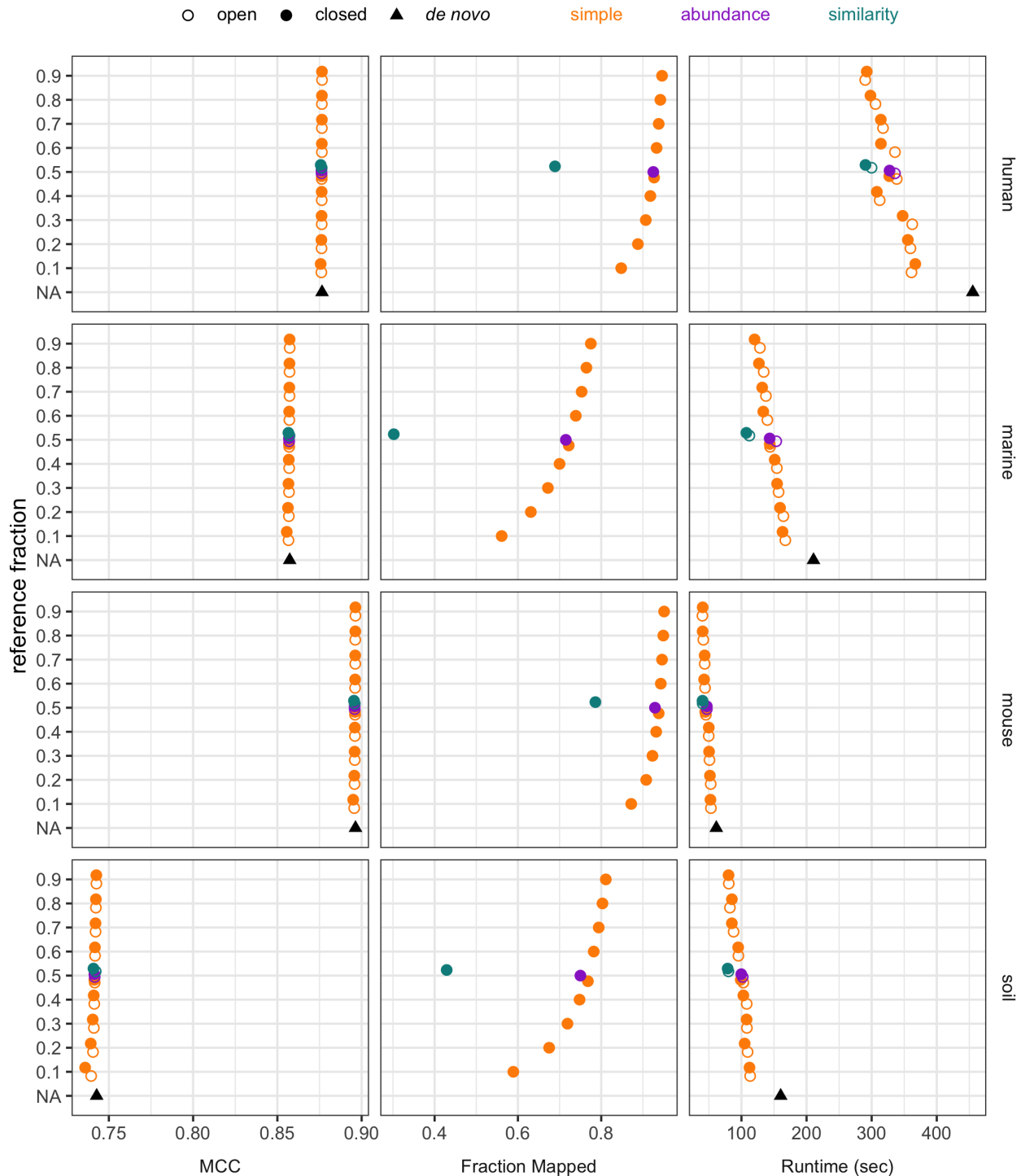


Figure 4: Split dataset strategy. The median MCC score, fraction of query sequences that mapped in closed-reference clustering, and runtime in seconds from repeating each clustering method 100 times. Each dataset was split into a reference and query fraction. Reference sequences were selected via a simple random sample, weighting sequences by relative abundance, or weighting by similarity to other sequences in the dataset. With the simple random sample method, dataset splitting was repeated with reference fractions ranging from 10% to 90% of the dataset and for 100 random seeds. *De novo* clustering each dataset is also shown for comparison.

223 After testing the split strategy using a simple random sample to select the reference
224 sequences, we then investigated other methods of splitting the data. We tested three
225 methods for selecting the fraction of sequences to be used as the reference at a size of
226 50%: a simple random sample, weighting sequences by relative abundance, and weighting
227 by similarity to other sequences in the dataset (Figure 4). OTU quality in terms of MCC
228 was similar across all three sampling methods (median MCC of 0.87). In closed-reference
229 clustering mode, the fraction of sequences that mapped were similar for simple and
230 abundance-weighted sampling (median fraction mapped of 0.85 and 0.84, respectively),
231 but worse for similarity-weighted sampling (median fraction mapped of 0.56). While simple
232 and abundance-weighted sampling produced better quality OTUs than similarity-weighted
233 sampling, OptiFit performed faster on similarity-weighted samples with a median runtime of
234 93.8 seconds compared to 123.2 and 122.6 seconds for simple and abundance-weighted
235 sampling, respectively. Thus, employing more complicated sampling strategies such as
236 abundance-weighted and similarity-weighted sampling did not confer any advantages over
237 selecting the reference via a simple random sample, and in fact decreased OTU quality in
238 the case of similarity-weighted sampling.

239 Discussion

240 We developed a new algorithm for clustering sequences to existing OTUs and have
241 demonstrated its suitability for reference-based clustering. OptiFit makes the iterative
242 method employed by OptiClust available for tasks where reference-based clustering is
243 required. We have shown that OTU quality is similar between OptiClust and OptiFit in open
244 reference mode, regardless of strategy employed. Open reference OptiFit performs slower
245 than OptiClust due to the additional *de novo* clustering step, so users may prefer OptiClust
246 for tasks that do not require reference OTUs.

247 When clustering to public databases, OTU quality dropped in closed reference mode to

248 different degrees depending on the database and dataset source, and no more than half
249 of query sequences were able to be clustered into OTUs across any dataset/database
250 combination. This may reflect limitations of reference databases, which are unlikely
251 to contain sequences from novel microbes. This drop in quality was most notable
252 with the RDP reference, which contained only 16,192 sequences compared to 173,648
253 sequences in SILVA and 174,979 in Greengenes. Note that Greengenes has not been
254 updated since 2013 at the time of this writing, while SILVA and the RDP are updated
255 regularly. We recommend that users who require an independent reference database
256 opt for large databases with regular updates and good coverage of microbial diversity for
257 their environment. Since OptiClust still performs faster than open reference OptiFit and
258 creates higher quality OTUs than closed reference OptiFit with the database strategy, we
259 recommend using OptiClust rather than clustering to a database whenever consistent
260 OTUs are not required.

261 The OptiClust and OptiFit algorithms produced higher quality OTUs than VSEARCH in
262 open reference, closed reference, or *de novo* modes. However, VSEARCH was able
263 to cluster more sequences to OTUs than OptiFit in closed reference mode. While both
264 OptiFit and VSEARCH use a distance or similarity threshold for determining how to cluster
265 sequences into OTUs, VSEARCH is more permissive than OptiFit regardless of mode.
266 The OptiFit and OptiClust algorithms use all of the sequences to define an OTU, preferring
267 that all pairs of sequences (including reference and query sequences) in an OTU are within
268 the distance threshold in order to maximize the MCC. In contrast, VSEARCH only requires
269 each query sequence to be similar to the single centroid sequence that seeded the OTU.
270 Because of this, VSEARCH sacrifices OTU quality by allowing more dissimilar sequences
271 to be clustered into OTUs.

272 When clustering with the split dataset strategy, OTU quality was remarkably similar when
273 reference sequences were selected by a simple random sample or weighted by abundance,

274 but quality was slightly worse when sequences were weighted by similarity. We recommend
275 using a simple random sample since the more sophisticated reference selection methods
276 do not offer any benefit. The similarity in OTU quality between OptiClust and OptiFit with
277 this strategy demonstrates the suitability of using OptiFit to cluster sequences to existing
278 OTUs, such as when comparing OTUs across studies. However, when consistent OTUs
279 are not required, we recommend using OptiClust for *de novo* clustering over the split
280 strategy with OptiFit since OptiClust is simpler to execute but performs similarly in terms of
281 both run time and OTU quality.

282 Unlike existing reference-based methods that cluster query sequences to a single centroid
283 sequence in each reference OTU, OptiFit considers all sequences in each reference OTU
284 when clustering query sequences, resulting in OTUs of a similar high quality as those
285 produced by the *de novo* OptiClust algorithm. Potential applications include clustering
286 sequences to reference databases, comparing taxonomic composition of microbiomes
287 across different studies, or using OTU-based machine learning models to make predictions
288 on new data. OptiFit fills the missing option for clustering query sequences to existing
289 OTUs that does not sacrifice OTU quality for consistency of OTU assignments.

290 **Materials and Methods**

291 **Data Processing Steps**

292 We downloaded 16S rRNA gene amplicon sequences from four published datasets isolated
293 from soil (8), marine (9), mouse gut (10), and human gut (11) samples. These datasets
294 contain sequences from the V4 region of the 16S rRNA gene and represent a selection
295 of the broad types of natural communities that microbial ecologists study. We processed
296 the raw sequences using mothur according to the Schloss Lab MiSeq SOP (15) and
297 accompanying study by Kozich *et al.* (16). These steps included trimming and filtering
298 for quality, aligning to the SILVA reference alignment (12), discarding sequences that

299 aligned outside the V4 region, removing chimeric reads with UCHIME (17), and calculating
300 distances between all pairs of sequences within each dataset prior to clustering.

301 **Reference database clustering**

302 To generate reference OTUs from public databases, we downloaded sequences from the
303 Greengenes database (v13_8_99) (6), SILVA non-redundant database (v132) (12), and the
304 Ribosomal Database Project (v16) (13). These sequences were processed using the same
305 steps outlined above followed by clustering sequences into *de novo* OTUs with OptiClust.
306 Processed reads from each of the four datasets were clustered with OptiFit to the reference
307 OTUs generated from each of the three databases. When reference clustering with
308 VSEARCH, processed datasets were clustered directly to the unprocessed Greengenes
309 97% OTU reference alignment, since this method is how VSEARCH is typically used by
310 the QIIME2 software for reference-based clustering (7, 18).

311 **Split dataset clustering**

312 For each dataset, half of the sequences were selected to be clustered *de novo* into
313 reference OTUs with OptiClust. We used three methods for selecting the subset of
314 sequences to be used as the reference: a simple random sample, weighting sequences by
315 relative abundance, and weighting by similarity to other sequences in the dataset. Dataset
316 splitting was repeated with 100 random seeds. With the simple random sampling method,
317 dataset splitting was also repeated with reference fractions ranging from 10% to 90% of
318 the dataset. For each dataset split, the remaining query sequences were clustered into the
319 reference OTUs with OptiFit.

320 **Benchmarking**

321 OptiClust and OptiFit randomize the order of query sequences prior to clustering and
322 employ a random number generator to break ties when OTU assignments are of equal

323 quality. As a result, they produce slightly different OTU assignments when repeated
324 with different random seeds. To capture any variation in OTU quality or execution time,
325 clustering was repeated with 100 random seeds for each combination of parameters and
326 input datasets. We used the benchmark feature provided by Snakemake to measure the
327 run time of every clustering job. We calculated the MCC on each set of OTUs to quantify
328 the quality of clustering, as described by Westcott *et al.* (1).

329 **Data and Code Availability**

330 We implemented the analysis workflow in Snakemake (19) and wrote scripts in R (20),
331 Python (21), and GNU bash (22). Software used includes mothur v1.47.0 (23), VSEARCH
332 v2.15.2 (5), the tidyverse metapackage (24), R Markdown (25), ggraph (26), ggtext (27),
333 numpy (28), the SRA toolkit (29), and conda (30). The complete workflow and supporting
334 files required to reproduce this manuscript are available at [https://github.com/SchlossLab/](https://github.com/SchlossLab/Sovacool_OptiFit_2021)
335 [Sovacool_OptiFit_2021](https://github.com/SchlossLab/Sovacool_OptiFit_2021).

336 **Acknowledgements**

337 We thank members of the Schloss Lab for their feedback on the figures.

338 KLS received support from the NIH Training Program in Bioinformatics (T32 GM070449).
339 Salary support for PDS came from NIH grants R01CA215574 and U01AI124255. The
340 funders had no role in study design, data collection and interpretation, or the decision to
341 submit the work for publication.

342 **Author Contributions**

343 KLS wrote the analysis code, evaluated the algorithm, and wrote the original draft of
344 the manuscript. SLW designed and implemented the OptiFit algorithm and assisted in
345 debugging the analysis code. MBM and GAD contributed analysis code. PDS conceived

346 the study, supervised the project, and assisted in debugging the analysis code. All authors
347 reviewed and edited the manuscript.

348 **References**

- 349 1. **Westcott SL, Schloss PD.** 2017. OptiClust, an Improved Method for Assigning
Amplicon-Based Sequence Data to Operational Taxonomic Units. *mSphere*
350 **2**:e00073–17. doi:10.1128/mSphereDirect.00073-17.
- 351 2. **Schloss PD.** 2016. Application of a Database-Independent Approach To Assess the
Quality of Operational Taxonomic Unit Picking Methods. *mSystems* **1**:e00027–16.
352 doi:10.1128/mSystems.00027-16.
- 353 3. **Westcott SL, Schloss PD.** 2015. De novo clustering methods outperform
reference-based methods for assigning 16S rRNA gene sequences to operational
354 taxonomic units. *PeerJ* **3**:e1487. doi:10.7717/peerj.1487.
- 355 4. **Edgar RC.** 2010. Search and clustering orders of magnitude faster than BLAST.
Bioinformatics **26**:2460–2461. doi:10.1093/bioinformatics/btq461.
356
- 357 5. **Rognes T, Flouri T, Nichols B, Quince C, Mahé F.** 2016. VSEARCH: A versatile
open source tool for metagenomics. *PeerJ* **4**:e2584. doi:10.7717/peerj.2584.
358
- 359 6. **DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, Huber
T, Dalevi D, Hu P, Andersen GL.** 2006. Greengenes, a Chimera-Checked 16S
rRNA Gene Database and Workbench Compatible with ARB. *AEM* **72**:5069–5072.
360 doi:10.1128/AEM.03006-05.
- 361 7. Clustering sequences into OTUs using Q2-vsearch — QIIME 2 2021.2.0
362 documentation. <https://docs.qiime2.org/2021.2/tutorials/otu-clustering/>.

- 363 8. **Johnston ER, Rodriguez-R LM, Luo C, Yuan MM, Wu L, He Z, Schuur EAG, Luo Y, Tiedje JM, Zhou J, Konstantinidis KT.** 2016. Metagenomics Reveals Pervasive Bacterial Populations and Reduced Community Diversity across the Alaska Tundra Ecosystem. *Front Microbiol* **7**. doi:10.3389/fmicb.2016.00579.
- 364
- 365 9. **Henson MW, Pitre DM, Weckhorst JL, Lanclos VC, Webber AT, Thrash JC.** 2016. Artificial Seawater Media Facilitate Cultivating Members of the Microbial Majority from the Gulf of Mexico. *mSphere* **1**. doi:10.1128/mSphere.00028-16.
- 366
- 367 10. **Schloss PD, Schubert AM, Zackular JP, Iverson KD, Young VB, Petrosino JF.** 2012. Stabilization of the murine gut microbiome following weaning. *Gut Microbes* **3**:383–393. doi:10.4161/gmic.21008.
- 368
- 369 11. **Baxter NT, Ruffin MT, Rogers MAM, Schloss PD.** 2016. Microbiota-based model improves the sensitivity of fecal immunochemical test for detecting colonic lesions. *Genome Med* **8**:37. doi:10.1186/s13073-016-0290-3.
- 370
- 371 12. **Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO.** 2013. The SILVA ribosomal RNA gene database project: Improved data processing and web-based tools. *Nucleic Acids Research* **41**:D590–D596. doi:10.1093/nar/gks1219.
- 372
- 373 13. **Cole JR, Wang Q, Fish JA, Chai B, McGarrell DM, Sun Y, Brown CT, Porras-Alfaro A, Kuske CR, Tiedje JM.** 2014. Ribosomal Database Project: Data and tools for high throughput rRNA analysis. *Nucl Acids Res* **42**:D633–D642. doi:10.1093/nar/gkt1244.
- 374

- 375 14. **Navas-Molina JA, Peralta-Sánchez JM, González A, McMurdie PJ, Vázquez-Baeza Y, Xu Z, Ursell LK, Lauber C, Zhou H, Song SJ, Huntley J, Ackermann GL, Berg-Lyons D, Holmes S, Caporaso JG, Knight R.** 2013. Chapter Nineteen - Advancing Our Understanding of the Human Microbiome Using QIIME, p. 371–444. *In* DeLong, EF (ed.), *Methods in Enzymology*. Academic Press.
- 376
- 377 15. **Schloss PD, Westcott SL.** MiSeq SOP. https://mothur.org/MiSeq_SOP.
- 378
- 379 16. **Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD.** 2013. Development of a Dual-Index Sequencing Strategy and Curation Pipeline for Analyzing Amplicon Sequence Data on the MiSeq Illumina Sequencing Platform. *Appl Environ Microbiol* **79**:5112–5120. doi:10.1128/AEM.01043-13.
- 380
- 381 17. **Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R.** 2011. UCHIME improves sensitivity and speed of chimera detection. *Bioinformatics* **27**:2194–2200. doi:10.1093/bioinformatics/btr381.
- 382

- 383 18. **Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, Alexander H, Alm EJ, Arumugam M, Asnicar F, Bai Y, Bisanz JE, Bittinger K, Brejnrod A, Brislawn CJ, Brown CT, Callahan BJ, Caraballo-Rodríguez AM, Chase J, Cope EK, Da Silva R, Diener C, Dorrestein PC, Douglas GM, Durall DM, Duvallet C, Edwardson CF, Ernst M, Estaki M, Fouquier J, Gauglitz JM, Gibbons SM, Gibson DL, Gonzalez A, Gorlick K, Guo J, Hillmann B, Holmes S, Holste H, Huttenhower C, Huttley GA, Janssen S, Jarmusch AK, Jiang L, Kaehler BD, Kang KB, Keefe CR, Keim P, Kelley ST, Knights D, Koester I, Kosciulek T, Kreps J, Langille MGI, Lee J, Ley R, Liu Y-X, Lofffield E, Lozupone C, Maher M, Marotz C, Martin BD, McDonald D, McIver LJ, Melnik AV, Metcalf JL, Morgan SC, Morton JT, Naimey AT, Navas-Molina JA, Nothias LF, Orchanian SB, Pearson T, Peoples SL, Petras D, Preuss ML, Pruesse E, Rasmussen LB, Rivers A, Robeson MS, Rosenthal P, Segata N, Shaffer M, Shiffer A, Sinha R, Song SJ, Spear JR, Swafford AD, Thompson LR, Torres PJ, Trinh P, Tripathi A, Turnbaugh PJ, UI-Hasan S, van der Hooft JJJ, Vargas F, Vázquez-Baeza Y, Vogtmann E, von Hippel M, Walters W, Wan Y, Wang M, Warren J, Weber KC, Williamson CHD, Willis AD, Xu ZZ, Zaneveld JR, Zhang Y, Zhu Q, Knight R, Caporaso JG.** 2019. Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nat Biotechnol* **37**:852–857. doi:10.1038/s41587-019-0209-9.
- 384
- 385 19. **Köster J, Rahmann S.** 2012. Snakemake — a scalable bioinformatics workflow engine. *Bioinformatics* **28**:2520–2522. doi:10.1093/bioinformatics/bts480.
- 386
- 387 20. **R Core Team.** 2020. R: A language and environment for statistical computing. Manual, R Foundation for Statistical Computing, Vienna, Austria.
- 388
- 389 21. **Van Rossum G, Drake FL.** 2009. Python 3 Reference Manual | Guide books.
- 390

- 391 22. Bash Reference Manual. <https://www.gnu.org/software/bash/manual/bash.html>.
392
- 393 23. **Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn DJ, Weber CF.** 2009. Introducing mothur: Open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Applied and Environmental Microbiology* **75**:7537–7541. doi:10.1128/AEM.01541-09.
394
- 395 24. **Wickham H, Averick M, Bryan J, Chang W, McGowan LD, François R, Grolemond G, Hayes A, Henry L, Hester J, Kuhn M, Pedersen TL, Miller E, Bache SM, Müller K, Ooms J, Robinson D, Seidel DP, Spinu V, Takahashi K, Vaughan D, Wilke C, Woo K, Yutani H.** 2019. Welcome to the Tidyverse. *Journal of Open Source Software* **4**:1686. doi:10.21105/joss.01686.
396
- 397 25. **Xie Y, Allaire JJ, Grolemond G.** 2018. *R Markdown: The Definitive Guide*. Taylor & Francis, CRC Press.
398
- 399 26. **Pedersen TL.** 2021. *Ggraph: An implementation of grammar of graphics for graphs and networks*.
400
- 401 27. **Wilke CO.** 2020. *Ggtext: Improved text rendering support for 'Ggplot2'*. Manual.
402
- 403 28. **Harris CR, Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, Wieser E, Taylor J, Berg S, Smith NJ, Kern R, Picus M, Hoyer S, van Kerkwijk MH, Brett M, Haldane A, del Río JF, Wiebe M, Peterson P, Gérard-Marchant P, Sheppard K, Reddy T, Weckesser W, Abbasi H, Gohlke C, Oliphant TE.** 2020. Array programming with NumPy. *Nature* **585**:357–362. doi:10.1038/s41586-020-2649-2.
404

405 29. SRA-Tools - NCBI. <http://ncbi.github.io/sra-tools/>.

406

407 30. 2016. Anaconda Software Distribution. Anaconda Documentation. Anaconda Inc.

408