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

2022-02-03

Kelly L. Sovacool<sup>1</sup>, Sarah L. Westcott<sup>2</sup>, M. Brodie Mumphrey<sup>1</sup>, Gabrielle A. Dotson<sup>1</sup>, Patrick D. Schloss<sup>2,3,†</sup>

1 Department of Computational Medicine and Bioinformatics, University of Michigan

2 Department of Microbiology and Immunology, University of Michigan

3 Center for Computational Medicine and Bioinformatics, University of Michigan

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

Sovacool et al.

## Abstract

Assigning amplicon sequences to operational taxonomic units (OTUs) is an important 2 step in characterizing microbial communities across large datasets. A notable difference 3 between de novo clustering and database-dependent reference clustering methods is that 4 OTU assignments from *de novo* methods may change when new sequences are added. 5 However, one may wish to incorporate new samples to previously clustered datasets 6 without clustering all sequences again, such as when comparing across datasets or 7 deploying machine learning models. Existing reference-based methods produce consistent 8 OTUs, but only consider the similarity of each guery sequence to a single reference 9 sequence in an OTU, resulting in assignments that are worse than those generated by 10 de novo methods. To provide an efficient method to fit sequences to existing OTUs, we 11 developed the OptiFit algorithm. Inspired by the *de novo* OptiClust algorithm, OptiFit 12 considers the similarity of all pairs of reference and guery sequences to produce OTUs 13 of the best possible quality. We tested OptiFit using four datasets with two strategies: 1) 14 clustering to a reference database or 2) splitting the dataset into a reference and query set, 15 clustering the references using OptiClust, then clustering the queries to the references. 16 The result is an improved implementation of reference-based clustering. OptiFit produces 17 similar quality OTUs as OptiClust at faster speeds when using the split dataset strategy. 18 OptiFit provides a suitable option for users requiring consistent OTU assignments at the 19 same guality afforded by *de novo* clustering methods. 20

#### 21 Importance

Advancements in DNA sequencing technology have allowed researchers to affordably
generate millions of sequence reads from microorganisms in diverse environments.
Efficient and robust software tools are needed to assign microbial sequences into
taxonomic groups for characterization and comparison of communities. The OptiClust
algorithm produces high quality groups by comparing sequences to each other, but the

assignments can change when new sequences are added to a dataset, making it difficult 27 to compare different studies. Other approaches assign sequences to groups by comparing 28 them to sequences in a reference database to produce consistent assignments, but the 29 quality of the groups produced is reduced compared to OptiClust. We developed OptiFit, a 30 new reference-based algorithm that produces consistent yet high quality assignments like 31 OptiClust. OptiFit allows researchers to compare microbial communities across different 32 studies or add new data to existing studies without sacrificing the quality of the group 33 assignments. 34

# 35 Introduction

Amplicon sequencing is a mainstay of microbial ecology. Researchers can affordably 36 generate millions of sequences to characterize the composition of hundreds of samples 37 from microbial communities without the need for culturing. In many analysis pipelines, 38 16S rRNA gene sequences are assigned to operational taxonomic units (OTUs) to 39 facilitate comparison of taxonomic composition between communities to avoid the need 40 for taxonomic classification. A distance threshold of 3% (or sequence similarity of 97%) is 41 commonly used to cluster sequences into OTUs based on pairwise comparisons of the 42 sequences within the dataset. The method chosen for clustering affects the quality of OTU 43 assignments and thus may impact downstream analyses of community composition (1-3). 44 OTU quality can be conceptualized as how well the OTU assignments match the definition 45 set by the distance threshold, i.e. whether sequence pairs that are at least as similar as 46 the distance threshold are assigned to the same OTU and sequence pairs that are more 47 dissimilar than the distance threshold are assigned to different OTUs. 48

There are two main categories of OTU clustering algorithms: *de novo* and reference-based. 49 OptiClust is a *de novo* clustering algorithm which uses the distance score between all 50 pairs of sequences in the dataset to cluster them into OTUs by maximizing the Matthews 51 Correlation Coefficient (MCC) (1). This approach takes into account the distances between 52 all pairs of sequences when assigning query sequences to OTUs, in contrast to other de 53 novo methods such as the greedy clustering algorithms implemented in USEARCH and 54 VSEARCH (4, 5). In methods employing greedy clustering algorithms, only the distance 55 between each sequence and a representative centroid sequence in the OTU is considered 56 while clustering. As a result, distances between pairs of sequences in the same OTU are 57 frequently larger than the specified threshold, i.e. they are false positives. In contrast, the 58 OptiClust algorithm takes into account the distance between all pairs of sequences when 59 considering how to cluster sequences into OTUs and is thus less willing to take on false 60

Sovacool et al.

61 positives.

A limitation of *de novo* clustering is that different OTU assignments will be produced 62 when new sequences are added to a dataset, making it difficult to use *de novo* clustering 63 to compare OTUs between different studies. Furthermore, since de novo clustering 64 requires calculating and comparing distances between all sequences in a dataset, the 65 execution time can be slow and memory requirements can be prohibitive for very large 66 datasets. Reference clustering attempts to overcome the limitations of *de novo* clustering 67 methods by using a representative set of sequences from a database, with each reference 68 sequence seeding an OTU. Commonly, the Greengenes set of representative full length 69 sequences clustered at 97% similarity is used as the reference with VSEARCH (5-7). 70 Query sequences are then clustered into OTUs based on their similarity to the reference 71 sequences. Any query sequences that are not within the distance threshold to any of 72 the reference sequences are either thrown out (closed reference clustering) or clustered 73 de novo to create additional OTUs (open reference clustering). While reference-based 74 clustering is generally fast, it is limited by the diversity of the reference database. Novel 75 sequences in the sample will be lost in closed reference mode if they are not represented 76 by a similar sequence in the database. We previously found that the OptiClust de novo 77 clustering algorithm created the highest quality OTU assignments of all clustering methods 78 (1). 79

To overcome the limitations of current reference-based and *de novo* clustering algorithms while maintaining OTU quality, we developed OptiFit, a reference-based clustering algorithm. While other tools represent reference OTUs with a single sequence, OptiFit uses all sequences in existing OTUs as the reference and fits new sequences to those reference OTUs. In contrast to other tools, OptiFit considers all pairwise distance scores between reference and query sequences when assigning sequences to OTUs in order to produce OTUs of the highest possible quality. Here, we tested the OptiFit algorithm with

the reference as a public database (e.g. Greengenes) or *de novo* OTUs generated using a
reference set from the full dataset and compared the performance to existing tools. To
evaluate the OptiFit algorithm and compare to existing methods, we used four published
datasets isolated from soil (8), marine (9), mouse gut (10), and human gut (11) samples.
OptiFit is available within the mothur software program.

## 92 **Results**

#### **The OptiFit algorithm**

OptiFit leverages the method employed by OptiClust of iteratively assigning sequences 94 to OTUs to produce the highest quality OTUs possible, and extends this method for 95 reference-based clustering. OptiClust first seeds each sequence into its own OTU as a 96 singleton. Then for each sequence, OptiClust considers whether the sequence should 97 move to a different OTU or remain in its current OTU, choosing the option that results in 98 a better MCC score (1). The MCC uses all values from a confusion matrix and ranges 99 from negative one to one, with a score of one occurring when all sequence pairs are true 100 positives and true negatives, a score of negative one occurring when all pairs are false 101 positives and false negatives, and a score of zero when there are equal numbers of true 102 and false assignments (i.e. no better than random guessing). Sequence pairs that are 103 similar to each other (i.e. within the distance threshold) are counted as true positives if 104 they are clustered into the same OTU, and false negatives if they are not in the the same 105 OTU. Sequence pairs that are not similar to each other are true negatives if they are not 106 clustered into the same OTU, and false positives if they are in the same OTU. Thus, a pair 107 of sequences is considered correctly assigned when their OTU assignment matches the 108 OTU definition set by the distance threshold. OptiClust iterations continue until the MCC 109 stabilizes or until a maximum number of iterations is reached. This process produces de 110 novo OTU assignments with the most optimal MCC given the input sequences. 11

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 guery 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 guery and reference sequence pairs are 118 considered when calculating the MCC. This process is repeated until the MCC changes by 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 guery 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

#### 129 Reference clustering with public databases

To test how OptiFit performs for reference-based clustering, we clustered each dataset 130 to three databases of reference OTUs: the Greengenes database v13 8 99 (6), the 131 SILVA non-redundant database v132 (12), and the Ribosomal Database Project (RDP) v16 132 (13). Reference OTUs for each database were created by performing *de novo* clustering 133 with OptiClust at a distance threshold of 3% using the V4 region of each sequence (see 134 Figure 2). After trimming to the V4 region, the databases contained 174,979, 16,192, and 135 173,648 unique sequences and produced de novo MCC scores of 0.72, 0.74, and 0.73 for 136 Greengenes, RDP, and SILVA, respectively. Clustering guery sequences with OptiFit to 137

Greengenes and SILVA in closed reference mode performed similarly, with median MCC 138 scores of 0.85 and 0.77 respectively, while the median MCC was 0.35 when clustering to 139 RDP (Figure 3; "db: Greengenes", "db: SILVA", and "db: RDP"). For comparison, clustering 140 datasets with OptiClust produced an average MCC score of 0.86 (Figure 3; "de novo"). 141 This gap in OTU quality mostly disappeared when clustering in open reference mode, 142 which produced median MCCs of 0.86 with Greengenes, 0.86 with SILVA, and 0.86 with 143 the RDP. Thus, open reference OptiFit produced OTUs of very similar quality as de novo 144 clustering with OptiClust, and closed reference OptiFit followed closely behind as long as a 145 suitable reference database was chosen. 146

Since closed reference clustering does not cluster guery sequences that could not be 147 clustered into reference OTUs, an additional measure of clustering performance to consider 148 is the fraction of query sequences that were able to be clustered. On average, more 149 sequences were clustered with Greengenes as the reference (59%) than with SILVA (50%) 150 or with the RDP (9.7%) (Figure 3). This mirrored the result reported above that Greengenes 151 produced better OTUs in terms of MCC score than either SILVA or RDP. Note that de novo 152 and open reference clustering methods always cluster 100% of sequences into OTUs. 153 The database chosen affects the final closed reference OTU assignments considerably in 154 terms of both MCC score and fraction of query sequences that could be clustered into the 155 reference OTUs. 156

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

To compare to the reference clustering methods used by QIIME2, we clustered each 166 dataset with VSEARCH against the Greengenes database of OTUs previously clustered 167 at 97% sequence similarity. Each reference OTU from the Greengenes 97% database 168 contains one reference sequence, and VSEARCH maps sequences to the reference 169 based on each individual guery sequence's similarity to the single reference sequence. 170 In contrast, OptiFit accepts reference OTUs which each may contain multiple sequences, 171 and the sequence similarity between all query and reference sequences is considered 172 when assigning sequences to OTUs. In closed reference mode, OptiFit produced 27.2% 173 higher quality OTUs than VSEARCH in terms of MCC score, but VSEARCH was able to 174 cluster 24.9% more query sequences than OptiFit to the Greengenes reference database 175 (Figure 3). This is because VSEARCH only considers the distances between each query 176 sequence to the single reference sequence, while OptiFit considers the distances between 177 all pairs of reference and query sequences in an OTU. When open reference clustering, 178 OptiFit produced higher quality OTUs than VSEARCH against the Greengenes database, 179 with median MCC scores of 0.86 and 0.56, respectively. In terms of run time, OptiFit 180 outperformed VSEARCH in both closed and open reference mode by 53.6% and 44.0% on 181 average, respectively. Thus, the more stringent OTU definition employed by OptiFit, which 182 prefers the query sequence to be similar to all other sequences in the OTU rather than to 183 only one sequence, resulted in fewer sequences being clustered to reference OTUs than 184 when using VSEARCH, but caused OptiFit to outperform VSEARCH in terms of both OTU 185 quality and execution time. 186

#### **187** Reference clustering with split datasets

When performing reference clustering against public databases, the database chosen 188 greatly affects the quality of OTUs produced. OTU quality may be poor when the reference 189 database consists of sequences that are too unrelated to the samples of interest, such as 190 when samples contain novel populations. While *de novo* clustering overcomes the quality 191 limitations of reference clustering to databases, OTU assignments are not consistent when 192 new sequences are added. Researchers may wish to cluster new sequences to existing 193 OTUs or to compare OTUs across studies. To determine how well OptiFit performs for 194 clustering new sequences to existing OTUs, we employed a split dataset strategy, where 195 each dataset was randomly split into a reference fraction and a query fraction. Reference 196 sequences were clustered de novo with OptiClust, then query sequences were clustered 197 to the *de novo* OTUs with OptiFit. 198

First, we tested whether OptiFit performed as well as *de novo* clustering when using the 199 split dataset strategy with half of the sequences selected for the reference by a simple 200 random sample (a 50% split) (Figure 3; "self-split"). OTU guality was similar to that from 201 OptiClust regardless of mode (0.031% difference in median MCC). In closed reference 202 mode, OptiFit was able to cluster 84.9% of guery sequences to reference OTUs with the 203 split strategy, a great improvement over the average 59% of sequences clustered to the 204 Greengenes database. In terms of run time, closed and open reference OptiFit performed 205 faster than OptiClust on whole datasets by 39.6% and 36.8%, respectively. Random 206 access memory (RAM) usage was similar, with OptiFit requiring slightly more RAM in 207 gigabytes than OptiClust. Open and closed reference OptiFit required 1.8% and 1.2% 208 more RAM than OptiClust, respectively (data not shown). The split dataset strategy also 209 performed 6.7% faster than the database strategy in closed reference mode and 65.5% 210 faster in open reference mode. Thus, reference clustering with the split dataset strategy 211 creates as high quality OTUs as *de novo* clustering yet at a faster run time, and fits far 212

<sup>213</sup> more query sequences than the database strategy.

While we initially tested this strategy using a 50% split of the data into reference and guery 214 fractions, we next investigated whether there was an optimal reference fraction size. To 215 identify the best reference size, reference sets with 10% to 90% of the sequences were 216 created, with the remaining sequences used for the query (Figure 4). OTU quality was 217 remarkably consistent across reference fraction sizes. For example, splitting the human 218 dataset 100 times yielded a coefficient of variation (i.e. the standard deviation divided by 219 the mean) of 0.0018 for the MCC score across all fractions. Run time generally decreased 220 as the reference fraction increased; for the human dataset, the median run time was 221 364.0 seconds with 10% of sequences in the reference and 290.8 seconds with 90% of 222 sequences in the reference. The RAM usage was virtually the same across reference 223 fraction sizes, with a coefficient of variation of 0.00089 for the human dataset (data not 224 shown). In closed reference mode, the fraction of sequences that mapped increased as 225 the reference size increased; for the human dataset, the median fraction mapped was 0.85 226 with 10% of sequences in the reference and 0.95 with 90% of sequences in the reference. 227 These trends held for the other datasets as well. Thus, the reference fraction did not affect 228 OTU quality in terms of MCC score nor the memory usage, but did affect the run time and 229 the fraction of sequences that mapped during the closed reference clustering. 230

After testing the split strategy using a simple random sample to select the reference 231 sequences, we then investigated other methods of splitting the data. We tested three 232 methods for selecting the fraction of sequences to be used as the reference at a size of 233 50%: a simple random sample, weighting sequences by relative abundance, and weighting 234 by similarity to other sequences in the dataset (Figure 4). OTU quality in terms of MCC 235 was similar across all three sampling methods (median MCC of 0.86). In closed-reference 236 clustering mode, the fraction of sequences that mapped were similar for simple and 237 abundance-weighted sampling (median fraction mapped of 0.85 and 0.84, respectively), 238

but worse for similarity-weighted sampling (median fraction mapped of 0.56). While simple 239 and abundance-weighted sampling produced better guality OTUs than similarity-weighted 240 sampling, OptiFit performed faster on similarity-weighted samples with a median runtime of 241 103.9 seconds compared to 135.4 and 134.8 seconds for simple and abundance-weighted 242 sampling, respectively. Thus, employing more complicated sampling strategies such as 243 abundance-weighted and similarity-weighted sampling did not confer any advantages over 244 selecting the reference via a simple random sample, and in fact decreased OTU quality in 245 the case of similarity-weighted sampling. 246

# 247 Discussion

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

When clustering to public databases, OTU quality dropped in closed reference mode to 255 different degrees depending on the database and dataset source, and no more than half 256 of query sequences were able to be clustered into OTUs across any dataset/database 257 combination. This may reflect limitations of reference databases, which are unlikely 258 to contain sequences from novel microbes. This drop in quality was most notable 259 with the RDP reference, which contained only 16,192 sequences compared to 173,648 260 sequences in SILVA and 174,979 in Greengenes. Note that Greengenes has not been 261 updated since 2013 at the time of this writing, while SILVA and the RDP are updated 262 regularly. We recommend that users who require an independent reference database 263

opt for large databases with regular updates and good coverage of microbial diversity for
 their environment. Since OptiClust still performs faster than open reference OptiFit and
 creates higher quality OTUs than closed reference OptiFit with the database strategy, we
 recommend using OptiClust rather than clustering to a database whenever consistent
 OTUs are not required.

The OptiClust and OptiFit algorithms produced higher quality OTUs than VSEARCH in 269 open reference, closed reference, or *de novo* modes. However, VSEARCH was able 270 to cluster more sequences to OTUs than OptiFit in closed reference mode. While both 271 OptiFit and VSEARCH use a distance or similarity threshold for determining how to cluster 272 sequences into OTUs, VSEARCH is more permissive than OptiFit regardless of mode. 273 The OptiFit and OptiClust algorithms use all of the sequences to define an OTU, preferring 274 that all pairs of sequences (including reference and query sequences) in an OTU are within 275 the distance threshold in order to maximize the MCC. In contrast, VSEARCH only requires 276 each query sequence to be similar to the single centroid sequence that seeded the OTU, 277 thus allowing pairs of query sequences to be less similar to each other than the threshold 278 specified. Because of this, VSEARCH sacrifices OTU quality by allowing more dissimilar 279 sequences to be clustered into the same OTUs. 280

When clustering with the split dataset strategy, OTU guality was remarkably similar when 281 reference sequences were selected by a simple random sample or weighted by abundance, 282 but quality was slightly worse when sequences were weighted by similarity. We recommend 283 using a simple random sample since the more sophisticated reference selection methods 284 do not offer any benefit. The similarity in OTU guality between OptiClust and OptiFit with 285 this strategy demonstrates the suitability of using OptiFit to cluster sequences to existing 286 OTUs, such as when comparing OTUs across studies. However, when consistent OTUs 287 are not required, we recommend using OptiClust for *de novo* clustering over the split 288 strategy with OptiFit since OptiClust is simpler to execute but performs similarly in terms of 289

Sovacool et al.

<sup>290</sup> both run time and OTU quality.

Unlike existing reference-based methods that cluster guery sequences to a single centroid 291 sequence in each reference OTU, OptiFit considers all sequences in each reference OTU 292 when clustering query sequences, resulting in OTUs of a similar high quality as those 293 produced by the *de novo* OptiClust algorithm. Potential applications include clustering 294 sequences to reference databases, comparing taxonomic composition of microbiomes 295 across different studies, or using OTU-based machine learning models to make predictions 296 on new data. OptiFit fills the missing option for clustering query sequences to existing 297 OTUs that does not sacrifice OTU quality for consistency of OTU assignments. 298

# 299 Materials and Methods

#### **Data processing steps**

We downloaded 16S rRNA gene amplicon sequences from four published datasets isolated 301 from soil (8), marine (9), mouse gut (10), and human gut (11) samples. These datasets 302 contain sequences from the V4 region of the 16S rRNA gene and represent a selection 303 of the broad types of natural communities that microbial ecologists study. We processed 304 the raw sequences using mothur according to the Schloss Lab MiSeg SOP (15) and 305 accompanying study by Kozich et al. (16). These steps included trimming and filtering 306 for quality, aligning to the SILVA reference alignment (12), discarding sequences that 307 aligned outside the V4 region, removing chimeric reads with UCHIME (17), and calculating 308 distances between all pairs of sequences within each dataset prior to clustering. 309

#### **Reference database clustering**

To generate reference OTUs from public databases, we downloaded sequences from the Greengenes database (v13\_8\_99) (6), SILVA non-redundant database (v132) (12), and the Ribosomal Database Project (v16) (13). These sequences were processed using the same steps outlined above followed by clustering sequences into *de novo* OTUs with OptiClust.
Processed reads from each of the four datasets were clustered with OptiFit to the reference
OTUs generated from each of the three databases. When reference clustering with
VSEARCH, processed datasets were clustered directly to the unprocessed Greengenes
97% OTU reference alignment, since this method is how VSEARCH is typically used by
the QIIME2 software for reference-based clustering (7, 18).

#### 320 Split dataset clustering

For each dataset, half of the sequences were selected to be clustered de novo into 321 reference OTUs with OptiClust. We used three methods for selecting the subset of 322 sequences to be used as the reference: a simple random sample, weighting sequences by 323 relative abundance, and weighting by similarity to other sequences in the dataset. Dataset 324 splitting was repeated with 100 random seeds. With the simple random sampling method, 325 dataset splitting was also repeated with reference fractions ranging from 10% to 90% of 326 the dataset. For each dataset split, the remaining guery sequences were clustered into the 327 reference OTUs with OptiFit. 328

#### 329 Benchmarking

OptiClust and OptiFit randomize the order of guery sequences prior to clustering and 330 employ a random number generator to break ties when OTU assignments are of equal 331 quality. As a result, they produce slightly different OTU assignments when repeated 332 with different random seeds. To capture any variation in OTU guality or execution time, 333 clustering was repeated with 100 random seeds for each combination of parameters and 334 input datasets. We used the benchmark feature provided by Snakemake to measure the 335 run time of every clustering job. We calculated the MCC on each set of OTUs to quantify 336 the quality of clustering, as described by Westcott et al. (1). 337

#### **Data and code availability**

We implemented the analysis workflow in Snakemake (19) and wrote scripts in R (20), Python (21), and GNU bash (22). Software used includes mothur v1.47.0 (23), VSEARCH v2.15.2 (5), the tidyverse metapackage (24), R Markdown (25), ggraph (26), ggtext (27), numpy (28), the SRA toolkit (29), and conda (30). The complete workflow and supporting files required to reproduce this manuscript are available at https://github.com/SchlossLab/ Sovacool\_OptiFit\_mSphere\_2022.

# 345 Acknowledgements

<sup>346</sup> We thank members of the Schloss Lab for their feedback on the figures.

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

# **351** Author Contributions

KLS wrote the analysis code, evaluated the algorithm, and wrote the original draft of the manuscript. SLW designed and implemented the OptiFit algorithm and assisted in debugging the analysis code. MBM and GAD contributed analysis code. PDS conceived the study, supervised the project, and assisted in debugging the analysis code. All authors reviewed and edited the manuscript.

# 357 **References**

- Westcott SL, Schloss PD. 2017. OptiClust, an Improved Method for Assigning Amplicon-Based Sequence Data to Operational Taxonomic Units. mSphere 2:e00073–17. doi:10.1128/mSphereDirect.00073-17.
- Schloss PD. 2016. Application of a Database-Independent Approach To Assess the Quality of Operational Taxonomic Unit Picking Methods. mSystems 1:e00027–16.
   doi:10.1128/mSystems.00027-16.
- Westcott SL, Schloss PD. 2015. De novo clustering methods outperform reference-based methods for assigning 16S rRNA gene sequences to operational taxonomic units. PeerJ 3:e1487. doi:10.7717/peerj.1487.
- 4. Edgar RC. 2010. Search and clustering orders of magnitude faster than BLAST.
   Bioinformatics 26:2460–2461. doi:10.1093/bioinformatics/btq461.
- <sup>366</sup> 5. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. 2016. VSEARCH: A versatile
   <sup>367</sup> open source tool for metagenomics. PeerJ 4:e2584. doi:10.7717/peerj.2584.
- 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. doi:10.1128/AEM.03006-05.
- 7. Clustering sequences into OTUs using Q2-vsearch QIIME 2 2021.2.0
   documentation. https://docs.qiime2.org/2021.2/tutorials/otu-clustering/.

- 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.
- 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.
- Schloss PD, Schubert AM, Zackular JP, Iverson KD, Young VB, Petrosino JF.
   2012. Stabilization of the murine gut microbiome following weaning. Gut Microbes
   377
   3333–393. doi:10.4161/gmic.21008.
- 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.
- <sup>380</sup> 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.
- 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.

14. Navas-Molina JA, Peralta-Sánchez JM, González A, McMurdie PJ, 384 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. 385 15. Schloss PD, Westcott SL. MiSeq SOP. https://mothur.org/MiSeq SOP. 386 387 16. Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD. 2013. 388 Development of a Dual-Index Sequencing Strategy and Curation Pipeline for Analyzing Amplicon Sequence Data on the MiSeg Illumina Sequencing Platform. Appl Environ Microbiol **79**:5112–5120. doi:10.1128/AEM.01043-13. 389 17. Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R. 2011. UCHIME 390 improves sensitivity and speed of chimera detection. Bioinformatics 27:2194–2200. doi:10.1093/bioinformatics/btr381. 391

Sovacool et al.

- 18. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, 392 Alexander H, Alm EJ, Arumugam M, Asnicar F, Bai Y, Bisanz JE, Bittinger K, Breinrod 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, Fouguier 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, Kosciolek T, Kreps J, Langille MGI, Lee J, Ley R, Liu Y-X, Loftfield 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, Ul-Hasan S, van der Hooft JJJ, Vargas F, Vázguez-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. 393
- <sup>394</sup> 19. Köster J, Rahmann S. 2012. Snakemake a scalable bioinformatics workflow
   <sup>395</sup> engine. Bioinformatics 28:2520–2522. doi:10.1093/bioinformatics/bts480.
- <sup>396</sup> 20. **R Core Team**. 2020. R: A language and environment for statistical computing.
   <sup>397</sup> Manual, R Foundation for Statistical Computing, Vienna, Austria.
- <sup>398</sup> 21. Van Rossum G, Drake FL. 2009. Python 3 Reference Manual | Guide books.
- 399

# 400 22. Bash Reference Manual. https://www.gnu.org/software/bash/manual/bash.html. 401

- 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.
- Wickham H, Averick M, Bryan J, Chang W, McGowan LD, François R, Grolemund 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.
- <sup>406</sup> 25. Xie Y, Allaire JJ, Grolemund G. 2018. R Markdown: The Definitive Guide. Taylor
   <sup>407</sup> & Francis, CRC Press.
- Pedersen TL. 2021. Ggraph: An implementation of grammar of graphics for graphs
   and networks.
- 410 27. Wilke CO. 2020. Ggtext: Improved text rendering support for 'Ggplot2'. Manual.

411

<sup>412</sup> 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.

29. SRA-Tools - NCBI. http://ncbi.github.io/sra-tools/.
30. 2016. Anaconda Software Distribution. Anaconda Documentation. Anaconda Inc.



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 of 3%. Previously, 50 reference sequences were clustered de novo with OptiClust (see the OptiClust supplemental text (1)). Reference sequences A through Q (colored orange) were within the distance threshold to at least one other reference sequence; the remaining reference sequences formed additional singleton OTUs (not shown). The goal of OptiFit is to assign the guery sequences W through Z (colored green) to the reference OTUs. Here, there are 50 reference sequences and 4 guery sequences which make 1,431 sequence pairs, of which 23 pairs are within the 3% distance threshold. Initially (step 1), OptiFit places each query sequence in its own OTU, resulting in 14 true positives, 9 false negatives, 0 false positives, and 1,408 true negatives for an MCC score of 0.78. Then, for each guery 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 (steps 2-4). 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 was needed. Note that sequence Z was dissimilar from all other sequences and thus it remained a singleton. The final MCC score is 0.91 with 20 true positives, 3 false negatives, 1 false positive, and 1407 true negatives.



**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.



**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 three clustering strategies; 1) *de novo* clustering the whole dataset using OptiClust, 2) splitting the dataset with 50% of the sequences as a reference set and the other 50% as a query set, clustering the references using OptiClust, then clustering the query sequences to the reference OTUs with OptiFit, and 3) clustering 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.



**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 with OptiClust is also shown for comparison.