

Machine learning classification by fitting amplicon sequences to existing OTUs

Running title: self-reference-based OTU clustering for ML classification

Courtney R. Armour¹, Kelly L. Sovacool², William L. Close^{1,*}, Begüm D. Topçuoğlu^{1,#}, Jenna Wiens³, Patrick D. Schloss^{1,†}

¹ Department of Microbiology and Immunology, University of Michigan, Ann Arbor, Michigan, USA

² Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, Michigan, USA

³ Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, Michigan, USA

* Current Affiliation: Bio-Rad Laboratories, Hercules, California, USA

Current Affiliation: Bristol Myers Squibb, Summit, New Jersey, USA

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

Observation format

1 **Abstract**

2 Machine learning classification using the gut microbiome relies on assigning 16S rRNA gene
3 sequences into operational taxonomic units (OTUs) to quantify microbial composition. OTU
4 abundances are then used to train a classification model that can be applied to classify new
5 samples. The standard approaches to clustering sequences include reference-based and *de*
6 *novo* clustering. Reference-based clustering requires a well-curated reference database that
7 may not exist for all systems. *De novo* clustering tends to produce higher quality OTU
8 assignments than reference-based, but clusters depend on the sequences in the dataset and
9 therefore OTU assignments will change when new samples are sequenced. This lack of stability
10 complicates machine learning classification since new sequences must be reclustered with the
11 old data and the model retrained with the new OTU assignments. The OptiFit algorithm
12 addresses these issues by fitting new sequences into existing OTUs. While OptiFit produces
13 high quality OTU clusters, it is unclear whether this method for fitting new sequence data into
14 existing OTUs will impact the performance of classification models trained with the older data.
15 We used OptiFit to cluster sequences into existing OTUs and evaluated model performance in
16 classifying a dataset containing samples from patients with and without colonic screen relevant
17 neoplasia (SRN). We compared the performance of this model to standard methods including
18 *de novo* and database-reference-based clustering. We found that using OptiFit performed as
19 well or better in classifying SRNs. OptiFit can streamline the process of classifying new samples
20 by avoiding the need to retrain models using reclustered sequences.

21 **Importance**

22 There is great potential for using microbiome data to aid in diagnosis. A challenge with OTU-
23 based classification models is that 16S rRNA gene sequences are often assigned to OTUs
24 based on similarity to other sequences in the dataset. If data are generated from new patients,
25 the old and new sequences must all be reassigned to OTUs and the classification model
26 retrained. Yet there is a desire to have a single, validated model that can be widely deployed.
27 To overcome this obstacle, we applied the OptiFit clustering algorithm to fit new sequence data
28 to existing OTUs allowing for reuse of the model. A random forest model implemented using
29 OptiFit performed as well as the traditional reassign and retrain approach. This result shows that
30 it is possible to train and apply machine learning models based on OTU relative abundance data
31 that do not require retraining or the use of a reference database.

32 There is increasing evidence for an association between the composition of the gut microbiome
33 and a variety of diseases, such as crohn's disease and colorectal cancer (1, 2). There is great
34 potential to diagnose disease with gut microbiome sequence data and machine learning.
35 Taxonomic composition of microbial communities can be assessed using amplicon sequencing
36 of the 16S rRNA gene, which is the input to classification models. Analysis of 16S rRNA gene
37 sequence data generally relies on assigning sequences into operational taxonomic units
38 (OTUs). The process of OTU clustering can either be reference-based or *de novo*. The quality
39 of OTUs generated with reference-based clustering is generally poor compared to those
40 generated with *de novo* clustering (3). While *de novo* clustering produces high-quality OTU
41 clusters where sequences are accurately grouped based on similarity thresholds, the resulting
42 OTU clusters depend on the sequences within the dataset and the addition of new data has the
43 potential to redefine OTU cluster composition. The unstable nature of *de novo* OTU clustering
44 complicates deployment of machine learning models since integration of additional data
45 requires reclustering all the data and retraining the model. The ability to integrate new data into
46 a validated model without reclustering and retraining could allow for the application of a single
47 model that can continually classify new data. Recently, Sovacool *et al.* introduced OptiFit, a
48 method for fitting new sequence data into existing OTUs (4). While OptiFit can effectively fit new
49 sequence data to existing OTU clusters, it is unknown if the use of OptiFit will have an impact
50 on classification performance. Here, we tested the ability of OptiFit to cluster new sequence
51 data into existing OTU clusters for the purpose of classifying disease based on gut microbiome
52 composition.

53 We compared the ability of several approaches for assigning 16S rRNA gene sequences to
54 OTUs including, *de novo* and reference-based clustering. For reference-based clustering, we
55 used closed-reference clustering to a public database (database-reference-based) and to OTUs
56 generated from a subset of the samples (self-reference-based). To test how the model

57 performance compared between these approaches, we used a publicly available dataset of 16S
58 rRNA gene sequences from stool samples of healthy subjects (n = 226) as well as subjects with
59 screen-relevant neoplasia (SRN) consisting of advanced adenoma and carcinoma (n = 229) (5).
60 For the *de novo* workflows, all the 16S rRNA sequence data was clustered into OTUs. The OTU
61 clustering was conducted using two common algorithms: 1) the OptiClust algorithm in mothur
62 (6) and 2) the VSEARCH algorithm used in QIIME2 (7, 8). For both algorithms, the resulting
63 abundance data was then split into training and testing sets, where the training set was used to
64 tune hyperparameters and ultimately train and select the model. The model was applied to the
65 testing set and performance evaluated (Figure 1A). We also conducted reference-based OTU
66 clustering using OptiFit to fit the sequence data into OTUs based on the greengenes reference
67 database. To compare with another commonly used method, we also used the VSEARCH
68 algorithm to fit the sequence data to the greengenes reference (Figure 1B). In the OptiFit self-
69 reference workflow, the data was split into a training and a testing set. The training set was
70 clustered into OTUs and used to train a classification model. The OptiFit algorithm was used to
71 fit sequence data of samples not part of the original dataset into the existing OTUs, and used
72 the same model to classify the samples (Figure 1C). For each of the workflows the process was
73 repeated for 100 random splits of the data to account for variation caused by the choice of the
74 random number generator seed.

75 We first examined the quality of the resulting OTU clusters from each method using the
76 Matthews correlation coefficient (MCC). MCC is a metric used to measure OTU cluster quality
77 based on the similarity of all pairs of sequences and whether they are appropriately clustered or
78 not (3). MCC scores range between negative one and one, and measure how well clustering
79 assignment correlates with the distance between sequences. To ensure that OptiFit
80 appropriately integrated new sequence data into the existing OTUs, we expected the MCC
81 scores produced by the OptiFit workflow to be similar to that of *de novo* clustering using the

82 OptiClust algorithm. In the OptiFit workflow the test data was fit to the clustered training data for
83 each of the 100 data splits resulting in an MCC score for each split of the data. In the remaining
84 workflows, the data was only clustered once and then split into the training and testing sets
85 resulting in a single MCC score for each method. Indeed, the MCC scores were similar between
86 the OptiClust *de novo* (MCC = 0.884) and OptiFit self-reference workflows (average MCC =
87 0.879, standard deviation = 0.002). Consistent with prior findings, the reference-based methods
88 produced lower MCC scores (OptiFit Greengenes MCC = 0.786; VSEARCH Greengenes MCC
89 = 0.531) than the *de novo* methods (OptiClust *de novo* MCC = 0.884; VSEARCH *de novo* MCC
90 = 0.641) (4). Another metric we examined for the OptiFit workflow was the fraction of sequences
91 from the test set that mapped to the reference OTUs. Since sequences that did not map to
92 reference OTUs were eliminated, if a high percentage of reads did not map to an OTU we
93 expected this loss of data to negatively impact classification performance. We found that loss of
94 data was not an issue since on average 99.8% (standard deviation = 0.68%) of sequences in
95 the subsampled test set mapped to the reference OTUs. This number is higher than the
96 average fraction of reads mapped in the OptiFit Greengenes workflow (96.8% +/- 3.5). These
97 results indicate that the OptiFit self-reference method performed as well as the OptiClust *de*
98 *novovo* method and better than using an external database.

99 We next assessed model performance using OTU relative abundances from the training data
100 from the workflows to train a model to predict SRNs and used the model on the held out data.
101 Using the predicted and actual diagnosis classification, we calculated the area under the
102 receiver operating characteristic curve (AUROC) for each data split. During cross-validation
103 (CV) training, the performance of the OptiFit self-reference and OptiClust *de novo* models were
104 not significantly different (p-value = 0.066; Figure 2A), while performance for both VSEARCH
105 methods was significantly lower than the OptiClust *de novo*, OptiFit self, and OptiFit
106 Greengenes methods (p-values < 0.05). The trained model was then applied to the test data

107 classifying samples as either control or SRN. The VSEARCH Greengenes method performed
108 slightly worse than the OptiClust *de novo* method (p-value = 0.030). However the performance
109 on the test data for the OptiClust *de novo*, OptiFit Greengenes, OptiFit self-reference, and
110 VSEARCH *de novo* approaches were not significantly different (p-values > 0.05; Figures 2B and
111 2C). These results indicate that new data could be fit to existing OTU clusters using OptiFit
112 without impacting model performance.

113 We tested the ability of OptiFit to integrate new data into existing OTUs for the purpose of
114 machine learning classification using OTU relative abundance. A potential problem with using
115 OptiFit is that any sequences from the new samples that do not map to the existing OTU
116 clusters will be discarded, resulting in a possible loss of information. However, we demonstrated
117 that OptiFit can be used to fit new sequence data into existing OTU clusters and it could perform
118 as well in predicting SRN compared to *de novo* clustering all the sequence data together. In this
119 instance, the performance of OptiFit was equivalent to using a database-reference-based
120 method despite the lower quality of the OTU clusters in the database-reference-based
121 approach. This likely indicates that the sequences that are important to the model are well
122 characterized by the reference database. However, a less well studied system may not be as
123 well characterized by a reference-database which would make the ability to utilize one's own
124 data a reference an exciting possibility. The ability to integrate data from new samples into
125 existing OTUs enables the implementation of a single machine learning model. This is important
126 for model implementation because not all of the data needs to be available or known at the time
127 of model generation. A robust machine learning model can be implemented as part of a non-
128 invasive and low-cost diagnostic for SRN and other diseases.

129 **Materials and Methods**

130 **Dataset.** Raw 16S rRNA gene sequence data from the V4 region were previously generated
131 from human stool samples. Sequences were downloaded from the NCBI Sequence Read

132 Archive (accession no. SRP062005) (5, 9). This dataset contains stool samples from 490
133 subjects. For this analysis, samples from subjects identified in the metadata as normal, high risk
134 normal, or adenoma were categorized as “normal”, while samples from subjects identified as
135 advanced adenoma or carcinoma were categorized as “screen relevant neoplasia” (SRN). The
136 resulting dataset consisted of 261 normal samples and 229 SRN samples.

137 **Data processing.** The full dataset was preprocessed with mothur (v1.47) (10) to join forward
138 and reverse reads, merge duplicate reads, align to the SILVA reference database (v132) (11),
139 precluster, remove chimeras with UCHIME (9), assign taxonomy, and remove non-bacterial
140 reads following the Schloss Lab MiSeq standard operating procedure described on the mothur
141 website (https://mothur.org/wiki/miseq_sop/). 100 splits of the 490 samples were generated
142 where 80% of the samples (392 samples) were randomly assigned to the training set and the
143 remaining 20% (98 samples) were assigned to the test set. Using 100 splits of the data
144 accounts for the variation that may be observed depending on the samples that are in the
145 training or test sets. Each sample was in the training set an average of 80 times (standard
146 deviation = 4.1) and the test set an average of 20 times (standard deviation = 4.1).

147 ***Reference-based workflows.***

148 1. OptiFit Self: The preprocess data was split into the training and testing sets. The training
149 set was clustered into OTUs using OptiClust, then the test set was fit to the OTUs of the
150 training set using the OptiFit algorithm (4). The OptiFit algorithm was run with method
151 open so that any sequences that did not map to the existing OTU clusters would form
152 new OTUs. The data was then subsampled to 10,000 reads and any novel OTUs from
153 the test set were removed. This process was repeated for each of the 100 splits resulting
154 in 100 training and testing datasets.

155 2. OptiFit Greengenes: Reference sequences from the Greengenes database v13_8_99
156 (12) were downloaded and processed with mothur by trimming to the V4 region and
157 clustered *de novo* with OptiClust (6). The preprocessed data was fit to the clustered
158 reference data using OptiFit with the method open to allow any sequences that did not
159 map to the existing reference clusters would form new OTUs. The data was then
160 subsampled to 10,000 reads and any novel OTUs from the test set were removed. The
161 dataset was then split into two sets where 80% of the samples were assigned to the
162 training set and 20% to the testing set. This process was repeated for each of the 100
163 splits resulting in 100 training and testing datasets.

164 3. VSEARCH Greengenes: Preprocessed data was clustered using VSEARCH v2.15.2 (7)
165 directly to unprocessed Greengenes 97% OTU reference alignment consistent with how
166 VSEARCH is typically used by the QIIME2 software for reference-based clustering (8).
167 The data was then subsampled to 10,000 reads and any novel OTUs from the test set
168 were removed. The dataset was then split into two sets where 80% of the samples were
169 assigned to the training set and 20% to the testing set. This process was repeated for
170 each of the 100 splits resulting in 100 training and testing datasets.

171 ***De novo workflows.***

172 1. OptiClust *de novo*: All the preprocessed data was clustered together with OptiClust (6) to
173 generate OTUs. The data was subsampled to 10,000 reads per sample and the resulting
174 abundance tables were split into the training and testing sets. The process was repeated
175 for each of the 100 splits resulting in 100 training and testing datasets.

176 2. VSEARCH *de novo*: All the preprocessed data was clustered using VSEARCH v2.15.2
177 (7) with 97% identity and then subsampled to 10,000 reads per sample. The process

178 was repeated for each of the 100 splits resulting in 100 training and testing datasets for
179 both workflows.

180 **Machine Learning.** A random forest model was trained with the R package mikrompl (v 1.2.0)
181 (13) to predict the diagnosis (SRN or normal) for the samples in the test set for each data split.
182 The training set was preprocessed to normalize OTU counts (scale and center), collapse
183 correlated OTUs, and remove OTUs with zero variance. The preprocessing from the training set
184 was then applied to the test set. Any OTUs in the test set that were not in the training set were
185 removed. P-values comparing model performance were calculated as previously described (14).
186 The averaged ROC curves were plotted by taking the average and standard deviation of the
187 sensitivity at each specificity value.

188 **Code Availability.**

189 The analysis workflow was implemented in Snakemake (15). Scripts for analysis were written in
190 R (16) and GNU bash (17). The software used includes mothur v1.47.0 (10), VSEARCH v2.15.2
191 (7), RStudio (18), the Tidyverse metapackage (19), R Markdown (20), the SRA toolkit (21), and
192 conda (22). The complete workflow and supporting files required to reproduce this study are
193 available at: https://github.com/SchlossLab/Armour_OptiFitGLNE_mBio_2023

194 **Acknowledgments**

195 This work was supported through a grant from the NIH (R01CA215574).

196 **References**

- 197 1. **Morgan XC, Tickle TL, Sokol H, Gevers D, Devaney KL, Ward DV, Reyes JA, Shah SA,**
198 **LeLeiko N, Snapper SB, Bousvaros A, Korzenik J, Sands BE, Xavier RJ, Huttenhower C.**
199 2012. Dysfunction of the intestinal microbiome in inflammatory bowel disease and treatment.
200 *Genome Biol* **13**:R79.
- 201 2. **Sobhani I, Tap J, Roudot-Thoraval F, Roperch JP, Letulle S, Langella P, Corthier G,**
202 **Tran Van Nhieu J, Furet JP.** 2011. Microbial dysbiosis in colorectal cancer (CRC) patients.
203 *PLoS One* **6**:e16393.
- 204 3. **Westcott SL, Schloss PD.** 2015. De novo clustering methods outperform reference-based
205 methods for assigning 16S rRNA gene sequences to operational taxonomic units. *PeerJ*
206 **3**:e1487. doi:10.7717/peerj.1487.
- 207 4. **Sovacool KL, Westcott SL, Mumphrey MB, Dotson GA, Schloss PD.** 2022. OptiFit: An
208 improved method for fitting amplicon sequences to existing OTUs. *mSphere* **7**:e00916–21.
209 doi:10.1128/msphere.00916-21.
- 210 5. **Baxter NT, Ruffin MT, Rogers MAM, Schloss PD.** 2016. Microbiota-based model improves
211 the sensitivity of fecal immunochemical test for detecting colonic lesions. *Genome Medicine*
212 **8**:37. doi:10.1186/s13073-016-0290-3.
- 213 6. **Westcott SL, Schloss PD.** 2017. OptiClust, an improved method for assigning amplicon-
214 based sequence data to operational taxonomic units. *mSphere* **2**:e00073–17.
215 doi:10.1128/mSphereDirect.00073-17.
- 216 7. **Rognes T, Flouri T, Nichols B, Quince C, Mahé F.** 2016. VSEARCH: a versatile open
217 source tool for metagenomics. *PeerJ* **4**:e2584. doi:10.7717/peerj.2584.

- 218 8. **Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, Alexander H,**
219 **Alm EJ, Arumugam M, Asnicar F, Bai Y, Bisanz JE, Bittinger K, Brejnrod A, Brislawn CJ,**
220 **Brown CT, Callahan BJ, Caraballo-Rodríguez AM, Chase J, Cope EK, Da Silva R, Diener**
221 **C, Dorrestein PC, Douglas GM, Durall DM, Duvallet C, Edwardson CF, Ernst M, Estaki M,**
222 **Fouquier J, Gauglitz JM, Gibbons SM, Gibson DL, Gonzalez A, Gorlick K, Guo J, Hillmann**
223 **B, Holmes S, Holste H, Huttenhower C, Huttley GA, Janssen S, Jarmusch AK, Jiang L,**
224 **Kaehler BD, Kang KB, Keefe CR, Keim P, Kelley ST, Knights D, Koester I, Kosciulek T,**
225 **Kreps J, Langille MGI, Lee J, Ley R, Liu Y-X, Lofffield E, Lozupone C, Maher M, Marotz C,**
226 **Martin BD, McDonald D, McIver LJ, Melnik AV, Metcalf JL, Morgan SC, Morton JT, Naimey**
227 **AT, Navas-Molina JA, Nothias LF, Orchanian SB, Pearson T, Peoples SL, Petras D,**
228 **Preuss ML, Pruesse E, Rasmussen LB, Rivers A, Robeson MS, Rosenthal P, Segata N,**
229 **Shaffer M, Shiffer A, Sinha R, Song SJ, Spear JR, Swafford AD, Thompson LR, Torres PJ,**
230 **Trinh P, Tripathi A, Turnbaugh PJ, Ul-Hasan S, Hooft JJJ van der, Vargas F, Vázquez-**
231 **Baeza Y, Vogtmann E, Hippel M von, Walters W, Wan Y, Wang M, Warren J, Weber KC,**
232 **Williamson CHD, Willis AD, Xu ZZ, Zaneveld JR, Zhang Y, Zhu Q, Knight R, Caporaso JG.**
233 2019. Reproducible, interactive, scalable and extensible microbiome data science using QIIME
234 2. *Nature Biotechnology* **37**:852–857. doi:10.1038/s41587-019-0209-9.
- 235 9. **Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R.** 2011. UCHIME improves
236 sensitivity and speed of chimera detection. *Bioinformatics* **27**:2194–2200.
237 doi:10.1093/bioinformatics/btr381.
- 238 10. **Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski**
239 **RA, Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn DJ,**
240 **Weber CF.** 2009. Introducing mothur: Open-source, platform-independent, community-
241 supported software for describing and comparing microbial communities. *Applied and*
242 *Environmental Microbiology* **75**:7537–7541. doi:10.1128/AEM.01541-09.

- 243 11. **Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO.**
244 2013. The SILVA ribosomal RNA gene database project: Improved data processing and web-
245 based tools. *Nucleic Acids Research* **41**:D590–D596. doi:10.1093/nar/gks1219.
- 246 12. **DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, Huber T, Dalevi D,**
247 **Hu P, Andersen GL.** 2006. Greengenes, a chimera-checked 16S rRNA gene database and
248 workbench compatible with ARB. *Applied and Environmental Microbiology* **72**:5069–5072.
249 doi:10.1128/AEM.03006-05.
- 250 13. **Topçuoğlu BD, Lapp Z, Sovacool KL, Snitkin E, Wiens J, Schloss PD.** 2021. mikropml:
251 User-Friendly R Package for Supervised Machine Learning Pipelines. *Journal of Open Source*
252 *Software* **6**:3073. doi:10.21105/joss.03073.
- 253 14. **Topçuoğlu BD, Lesniak NA, Ruffin MT, Wiens J, Schloss PD.** 2020. A framework for
254 effective application of machine learning to microbiome-based classification problems. *mBio*
255 **11**:e00434–20. doi:10.1128/mBio.00434-20.
- 256 15. **Koster J, Rahmann S.** 2012. Snakemake—a scalable bioinformatics workflow engine.
257 *Bioinformatics* **28**:2520–2522. doi:10.1093/bioinformatics/bts480.
- 258 16. **R Core Team.** 2020. R: A language and environment for statistical computing. R Foundation
259 for Statistical Computing, Vienna, Austria.
- 260 17. **GNU Project.** Bash reference manual.
- 261 18. **RStudio Team.** 2019. RStudio: Integrated development environment for r. RStudio, Inc.,
262 Boston, MA.
- 263 19. **Wickham H, Averick M, Bryan J, Chang W, McGowan LD, François R, Golemund G,**
264 **Hayes A, Henry L, Hester J, Kuhn M, Pedersen TL, Miller E, Bache SM, Müller K, Ooms J,**
265 **Robinson D, Seidel DP, Spinu V, Takahashi K, Vaughan D, Wilke C, Woo K, Yutani H.**

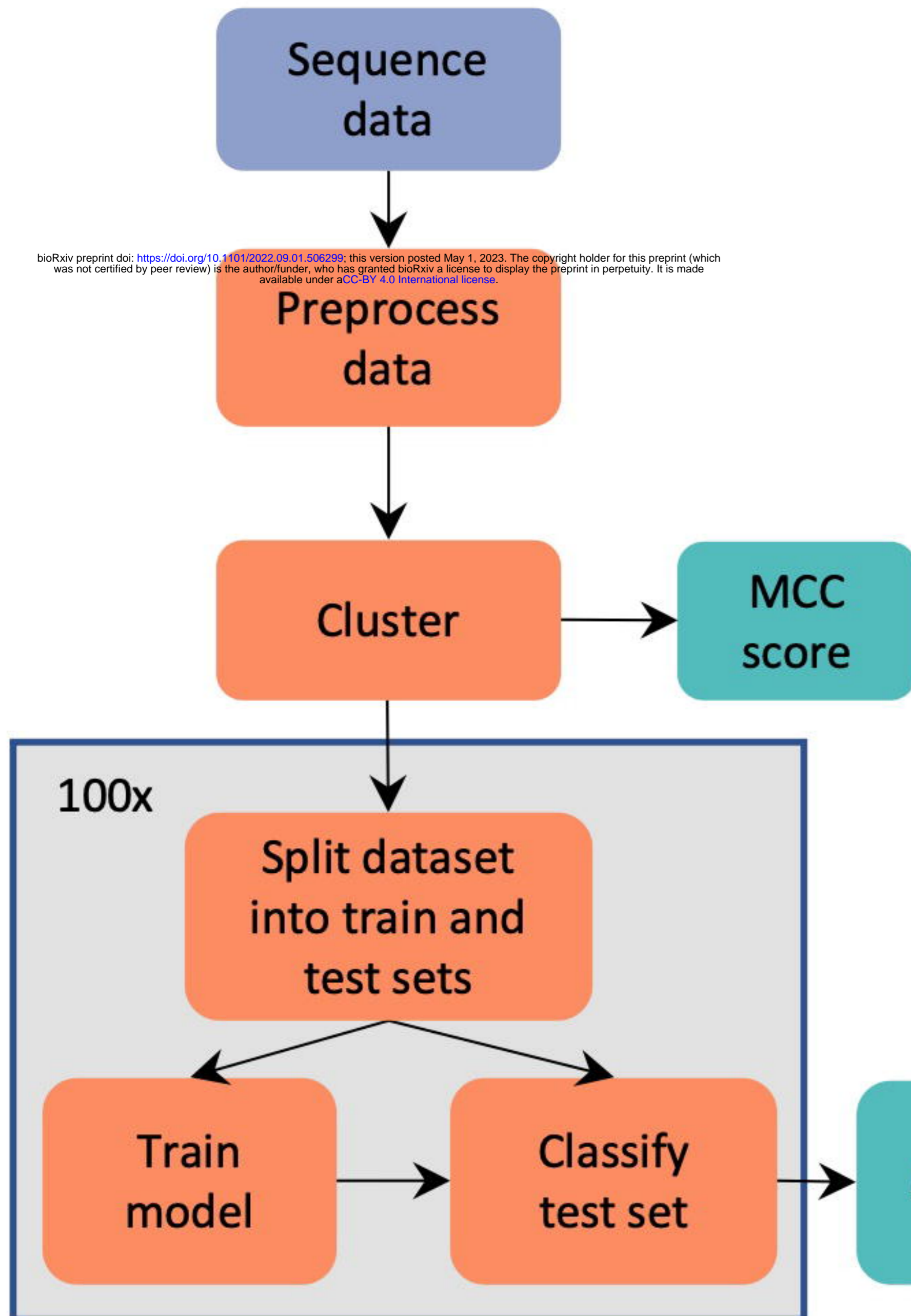
- 266 2019. Welcome to the Tidyverse. *Journal of Open Source Software* **4**:1686.
- 267 doi:10.21105/joss.01686.
- 268 20. **Xie Y, Allaire JJ, Golemund G**. 2018. *R Markdown: The Definitive Guide*. Taylor &
- 269 Francis, CRC Press.
- 270 21. SRA-Tools - NCBI. <http://ncbi.github.io/sra-tools/>.
- 271 22. 2016. *Anaconda Software Distribution. Anaconda Documentation*. Anaconda Inc.

272 **Figure Legends**

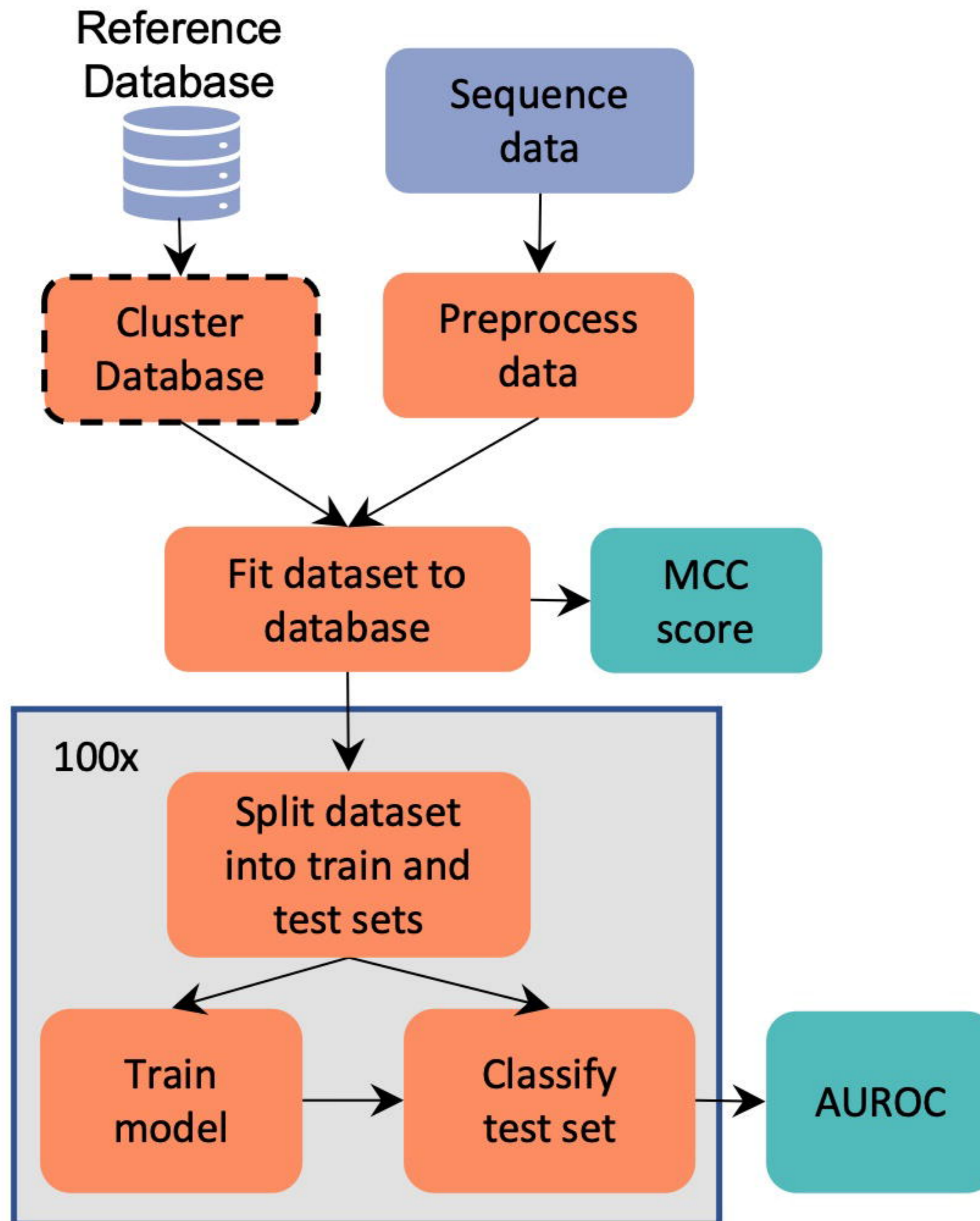
273 **Figure1: Overview of clustering workflows.** The *de novo* and database-reference-based
274 workflows were conducted using two approaches: OptiClust with mothur and VSEARCH as is
275 used in the QIIME pipeline.

276 **Figure 2: Model performance of OptiFit self-reference workflow is as good or better than**
277 **other methods. A)** Area under the receiver operating characteristic (AUROC) curve during
278 cross-validation (train) for the various workflows. **B)** AUROC on the test data for the various
279 workflows. The mean and standard deviation of the AUROC is represented by the black dot and
280 whiskers in panels A and B. The mean AUROC is printed below the points. **C)** Averaged
281 receiver operating characteristic (ROC) curves. Lines represent the average true positive rate
282 for the range of false positive rates.

A) *De novo*



B) Database reference-based



C) OptiFit self-reference-based

