

1 Optimizing taxonomic classification of marker gene
2 sequences

3
4 Nicholas A. Bokulich^{1#*}, Benjamin D. Kaehler^{2#*}, Jai Ram Rideout¹, Matthew Dillon¹, Evan
5 Bolyen¹, Rob Knight³, Gavin A. Huttley^{2#}, J. Gregory Caporaso^{1,4,#}

6
7 ¹The Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, AZ, USA

8 ²Research School of Biology, Australian National University, Canberra, Australia

9 ³Departments of Pediatrics and Computer Science & Engineering, and Center for
10 Microbiome Innovation, University of California San Diego, La Jolla, CA, USA

11 ⁴Department of Biological Sciences, Northern Arizona University, Flagstaff, AZ, USA

12

13 *These authors contributed equally

14

15 #Corresponding authors

16 Gregory Caporaso
17 Department of Biological Sciences
18 1298 S Knoles Drive
19 Building 56, 3rd Floor
20 Northern Arizona University
21 Flagstaff, AZ, USA
22 (303) 523-5485
23 (303) 523-4015 (fax)
24 Email: gregcaporaso@gmail.com

25

26 Nicholas Bokulich
27 The Pathogen & Microbiome Institute
28 PO Box 4073
29 Flagstaff, Arizona 86011-4073, USA
30 Email: nicholas.bokulich@nau.edu

31
32 Benjamin Kaehler
33 Research School of Biology
34 46 Sullivans Creek Road,
35 The Australian National University,
36 Acton ACT 2601, Australia
37 Email: benjamin.kaehler@anu.edu.au

38
39 Gavin Huttley
40 Research School of Biology
41 46 Sullivans Creek Road,
42 The Australian National University,
43 Acton ACT 2601, Australia
44 Email: gavin.huttley@anu.edu.au
45

46

47 **Abstract**

48 **Background:** Taxonomic classification of marker-gene sequences is an important step in
49 microbiome analysis. **Results:** We present q2-feature-classifier
50 (<https://github.com/qiime2/q2-feature-classifier>), a QIIME 2 plugin containing several
51 novel machine-learning and alignment-based taxonomy classifiers that meet or exceed
52 classification accuracy of existing methods. We evaluated and optimized several commonly
53 used taxonomic classification methods (RDP, BLAST, BLAST+, UCLUST) and several new
54 methods (a scikit-learn naive Bayes machine-learning classifier, and VSEARCH and
55 SortMeRNA alignment-based methods). **Conclusions:** Our results illustrate the importance
56 of parameter tuning for optimizing classifier performance, and we make explicit

57 recommendations regarding parameter choices for a range of standard operating
58 conditions. q2-feature-classifier and our evaluation framework, tax-credit, are both free,
59 open-source, BSD-licensed packages available on GitHub.

60

61 **Background**

62 High-throughput sequencing technologies have transformed our ability to explore
63 complex microbial communities, offering insight into microbial impacts on human health
64 [1] and global ecosystems [2]. This is achieved most commonly by sequencing short,
65 conserved marker genes amplified with ‘universal’ PCR primers, such as 16S rRNA for
66 bacteria and archaea, or internal transcribed spacer (ITS) regions for fungi. Targeted
67 marker-gene primers can also be used to profile specific taxa or functional groups, such as
68 nifH genes [3]. These sequences often are compared against an annotated reference
69 sequence database to determine the likely taxonomic origin of each sequence with as much
70 specificity as possible. Accurate and specific taxonomic information is a crucial component
71 of many experimental designs.

72 Challenges in this process include the short length of typical sequencing reads with
73 current technology, sequencing and PCR errors [4], selection of appropriate marker genes
74 that contain sufficient heterogeneity to differentiate target species but that are
75 homogeneous enough in some regions to design broad-spectrum primers, quality of

76 reference sequence annotations [5], and selection of a method that accurately predicts the
77 taxonomic affiliation of millions of sequences at minimal computational cost. Numerous
78 methods have been developed for taxonomy classification of DNA sequences, but few have
79 been directly compared in the specific case of short marker-gene sequences.

80 We introduce q2-feature-classifier, a QIIME 2 (<https://qiime2.org/>) plugin for
81 taxonomy classification of marker-gene sequences. QIIME 2 is the successor to the QIIME
82 [6] microbiome analysis package. The q2-feature-classifier plugin allows users to use any of
83 the numerous machine-learning classifiers available in scikit-learn [7][8] for marker gene
84 taxonomy classification, and currently provides two alignment-based taxonomy consensus
85 classifiers based on BLAST+ [9] and vsearch [10]. We evaluate the latter two methods and
86 the scikit-learn multinomial naive Bayes classifier (labelled “Naive Bayes” in the Results
87 section) for the first time. We show that the classifiers provided in q2-feature-classifier
88 match or outperform the classification accuracy of several widely-used methods for
89 sequence classification, and that performance of the naive Bayes classifier can be
90 significantly increased by providing it with information regarding expected taxonomic
91 composition.

92 We also developed tax-credit (<https://github.com/caporaso-lab/tax-credit/>), an
93 extensible computational framework for evaluating taxonomy classification accuracy. This
94 framework streamlines the process of methods benchmarking by compiling multiple
95 different test data sets, including mock communities [11] and simulated sequence reads. It
96 additionally stores pre-computed results from previously evaluated methods, including the
97 results presented here, and provides a framework for parameter sweeps and method

98 optimization. tax-credit could be used as an evaluation framework by other research
99 groups in the future, or its raw data could be easily extracted for integration in another
100 evaluation framework.

101

102 **Results**

103 We used tax-credit to optimize and compare multiple taxonomy classifiers. We
104 evaluated two commonly used, pre-existing classifiers that are wrapped in QIIME 1 (RDP
105 Classifier (version 2.2) [12], legacy BLAST (version 2.2.22) [13]), two QIIME 1 alignment-
106 based consensus taxonomy classifiers (the default UCLUST classifier available in QIIME 1
107 (based on version 1.2.22q) [14], and SortMeRNA (version 2.0 29/11/2014) [15]), two
108 alignment-based consensus taxonomy classifiers newly released in q2-feature-classifier
109 (based on BLAST+ (version 2.6.0) [9] and vsearch (version 2.0.3) [10]), and a new
110 multinomial naive Bayes machine-learning classifier in q2-feature-classifier (see materials
111 and methods for information about q2-feature-classifier methods and source code
112 availability). We performed parameter sweeps to determine optimal parameter
113 configurations for each method.

114 **Mock community evaluations**

115 We first benchmarked classifier performance on mock communities, which are
116 artificially constructed mixtures of microbial cells or DNA combined at known ratios [11].

117 We utilized 15 bacterial 16S rRNA mock communities and 4 fungal internal transcribed
118 spacer (ITS) mock communities (Table 1) sourced from mockrobiota [11], a public
119 repository for mock community data. Mock communities are useful for method
120 benchmarking because: 1) unlike for simulated communities, they allow quantitative
121 assessments of method performance under actual operating conditions, i.e., incorporating
122 real sequencing errors that can be difficult to model accurately; and 2) unlike for natural
123 community samples, the actual composition of a mock community is known in advance,
124 allowing quantitative assessments of community profiling accuracy.

125 An additional priority was to test the effect of setting class weights on classification
126 accuracy for the naive Bayes classifier implemented in q2-feature-classifier. In machine
127 learning, class weights or prior probabilities are vectors of weights that specify the
128 frequency at which each class is expected to be observed (and should be distinguished
129 from the use of this term under Bayesian inference as a probability distribution of weights
130 vectors). An alternative to setting class weights is to assume that each query sequence is
131 equally likely to belong to any of the taxa that are present in the reference sequence
132 database. This assumption, known as uniform class priors in the context of a naive Bayes
133 classifier, is made by the RDP classifier [12], and its impact on marker-gene classification
134 accuracy has yet to be validated. Making either assumption, that the class weights are
135 uniform or known to some extent, will affect results and cannot be avoided. The mock
136 communities have taxonomic abundances that are far from uniform over the set of
137 reference taxonomies, as any real data set must. We can therefore use them to assess the
138 impact of making assumptions regarding class weights. Where we have set the class

139 weights to the known taxonomic composition of a sample, we have labelled the results
140 “bespoke”.

141 We evaluated classifier performance accuracy on mock community sequences
142 classified at taxonomic levels from class through species. Mock community sequences were
143 classified using the Greengenes 99% OTUs 16S rRNA or UNITE 99% OTUs ITS reference
144 sequences for bacterial and fungal mock communities, respectively. As expected,
145 classification accuracy decreased as classification depth increased, and all methods could
146 predict the taxonomic affiliation of mock community sequences down to genus level with
147 median F-measures exceeding 0.8 across all parameter sets (minimum: UCLUST F=0.81,
148 maximum: Naive Bayes Bespoke F=1.00) (Figure 1A). However, species affiliation was
149 predicted with much lower and more variable accuracy among method configurations
150 (median F-measure minimum: UCLUST F=0.42, maximum: Naive Bayes Bespoke F=0.95),
151 highlighting the importance of parameter optimization (discussed in more detail below).
152 Figure 1A illustrates line plots of mean F-measure at each taxonomic level, averaged across
153 all classifier configurations; hence, classifier performance is underestimated for some
154 classifiers that are strongly affected by parameter configurations or for which a wider
155 range of parameters were tested (e.g., Naive Bayes). Comparing only optimized methods
156 (i.e., the top-performing parameter configurations for each method), Naive Bayes Bespoke
157 achieved significantly higher F-measure (paired t-test $P < 0.05$) (Figure 1B), recall, taxon
158 detection rate, and taxon accuracy rate scores (Figure 1C) and lower Bray-Curtis
159 dissimilarity than all other methods (Figure 1D).

160 Mock communities are necessarily simplistic, and cannot assess method
161 performance across a diverse range of taxa. Sequences matching the expected mock
162 community sequences are not removed from the reference database prior to classification,
163 in order to replicate normal operating conditions and assess recovery of expected
164 sequences. However, this approach may implicitly bias toward methods that find an exact
165 match to the query sequences, and does not approximate well natural microbial
166 communities in which few or no detected sequences exactly match the reference
167 sequences. Hence, we performed simulated sequence read classifications (described
168 below) to further test classifier performance.

169 **Cross-validated taxonomy classification**

170 Simulated sequence reads, derived from reference databases, allow us to assess
171 method performance across a greater diversity of sequences than a single mock community
172 generally encompasses. We first evaluated classifier performance using stratified k-fold
173 cross-validation of taxonomy classification to simulated reads. The k-fold cross-validation
174 strategy is modified slightly to account for the hierarchical nature of taxonomic
175 classifications, which all of the classifiers in this study (with the exception of legacy BLAST)
176 handle by assigning the lowest (i.e., most specific) taxonomic level where the classification
177 surpasses some user-defined “confidence” or “consensus” threshold (see materials and
178 methods). The modification is to truncate any expected taxonomy in each test set to the
179 maximum level at which an instance of that taxonomy exists in the training set. Simulated
180 reads were generated from Greengenes 99% OTUs 16S rRNA or UNITE 99% OTUs ITS

181 reference sequences with species-level annotations, and do not incorporate artificial
182 sequencing errors (see materials and methods for more details). In this set of tests and
183 below for novel taxa, the “bespoke” classifier had prior probabilities that were inferred
184 from the training set each time it was trained.

185 Classification of cross-validated reads yielded similar results to mock community
186 classification tests. For bacterial sequences, average classification accuracy for all methods
187 declined from near-perfect scores at family level (median F-measure minimum: BLAST+
188 $F=0.92$, maximum legacy BLAST $F=0.99$), but still retained accurate scores at species level
189 (median minimum: BLAST+ $F=0.76$, maximum SortMeRNA $F=0.84$), relative to some mock
190 community data sets (Figure 2A). Fungal sequences exhibited similar performance, with
191 the exception that mean BLAST+ and vsearch performance was markedly lower at all
192 taxonomic levels, indicating high sensitivity to parameter configurations, and species-level
193 F-measures were in general much lower (median minimum: BLAST+ $F=0.17$, maximum
194 UCLUST $F=0.45$) than those of bacterial sequence classifications (Figure 2A). At species
195 level, optimized UCLUST and SortMeRNA configurations achieved the highest F-measures
196 for 16S rRNA simulated sequences (Figure 2B). UCLUST achieved the highest F-measure
197 for ITS classification ($F = 0.51$). However, all optimized classifiers achieved similar F-
198 measure ranges, with the exception of legacy BLAST for ITS sequences (Figure 2B).

199

200

201 **Novel taxa evaluations**

202 Novel taxa classification offers a unique perspective on classifier behavior, assessing
203 how classifiers perform when challenged with a “novel” clade that is not represented in the
204 reference database. An ideal classifier should identify the nearest taxonomic lineage to
205 which this taxon belongs, but no further. In this evaluation, a reference database is
206 subsampled k times to generate query and reference sequence sets, as for cross-validated
207 classification, but two important distinctions exist: 1) the reference database used for
208 classification excludes any sequence that matches the taxonomic affiliation of the query
209 sequences at taxonomic level L , the taxonomic rank at which classification is being
210 attempted; and 2) this is performed at each taxonomic level, in order to assess
211 classification performance when each method encounters a “novel” species, genus, family,
212 et cetera.

213 Due to these differences, interpretation of novel taxa evaluation results is different
214 from that of mock community and cross-validated classifications. For the latter,
215 classification accuracy may be assessed at each taxonomic level for each classification
216 result: mean classification accuracy at family level and species level evaluate the same
217 results but focus on different taxonomic levels of classification. For novel taxa, however,
218 different query and reference sequences are compiled for classification at each taxonomic
219 level and separate classifications are performed for each. Hence, classifications at family
220 and species level are independent events — one assesses how accurately each method

221 performs when it encounters a “novel” family that is not represented in the reference
222 database, the other when a “novel” species is encountered.

223 Novel taxa evaluations employ a suite of modified metrics, to provide more
224 information on what types of classification errors occur. Precision, recall, and F-measure
225 calculations at each taxonomic level L assess whether an accurate taxonomy classification
226 was made at level $L-1$: for example, a “novel” species should be assigned a genus, because
227 the correct species class is not represented within the reference database. Any species-
228 level classification in this scenario is an *overclassification* (affecting both recall and
229 precision). Overclassification is one of the key metrics for novel taxa evaluation, indicating
230 the degree to which novel sequences will be interpreted as known organisms. This
231 overclassification is often highly undesirable because it leads, for example, to the incorrect
232 classification of unknown but harmless environmental sequences as known pathogens.
233 Novel sequences that are classified within the correct clade, but to a less specific level than
234 L , are *underclassified* (affecting recall but not precision). Sequences that are classified into a
235 completely different clade are *misclassified* (affecting both recall and precision).

236 Precision, recall, and F-measure all gradually increase from average scores near 0.0
237 at class level, reaching peak scores at genus level for bacteria and species level for fungi
238 (Figure 3A-C). These trends are paired with gradual decreases in underclassification and
239 misclassification rates for all classification methods, indicating that all classifiers perform
240 poorly when they encounter sequences with no known match at the class, order, or family
241 levels (Figure 3D-F). At species level, UCLUST, BLAST+, and vsearch achieved significantly
242 better F-measures than all other methods for 16S rRNA classifications ($P < 0.05$) (Figure

243 3G). UCLUST achieved significantly better F-measures than all other methods for ITS
244 classifications (Figure 3G). Over-, under-, and misclassification scores are less informative
245 for optimizing classifiers for real use cases, as most methods could be optimized to yield
246 near-zero scores for each of these metrics separately, but only through extreme
247 configurations, leading to F-measures that would be unacceptable under any scenario. Note
248 that all comparisons were made between methods optimized to maximize (or minimize) a
249 single metric, and hence the configurations that maximize precision are frequently
250 different from those that maximize recall or other metrics. This trade-off between different
251 metrics is discussed in more detail below.

252 The novel taxa evaluation provides an estimate of classifier performance given a
253 specific reference database, but its generalization is limited by the quality of the reference
254 databases available and by the label-based approach used for partitioning and evaluation.
255 Mislabeled and polyphyletic clades in the database, e.g. Clostridium group, increase the
256 probability of (potentially erroneous) misclassification. A complementary analysis based
257 on sequence similarity between a novel query and top reference hit could mitigate this
258 issue. However, we choose to apply a label-based approach, as it better reflects the
259 biological problem that users can expect to encounter; i.e., using a particular reference
260 sequence database (which will contain some quantity of mislabeled and polyphyletic taxa
261 inherent to currently available resources), how likely is a classifier to misclassify a
262 taxonomic label?

263

264 **Multi-evaluation method optimization**

265 The mock community and cross-validation classification evaluations yielded similar trends
266 in configuration performance, but optimizing parameters choices for the novel taxa
267 generally lead to suboptimal choices for the mock community and cross-validation tests
268 (Figure 4). We sought to determine the relationship between method configuration
269 performance for each evaluation, and use this information to select configurations that
270 perform best across all evaluations. For 16S rRNA sequence species-level classification,
271 method configurations that achieve maximum F-measures for mock and cross-validated
272 sequences perform poorly for novel taxa classification (Figure 4B). Optimization is more
273 straightforward for genus-level classification of 16S rRNA sequences (Figure 4A) and for
274 fungal sequences (Figure 4C-D), for which configuration performance (measured as mean
275 F-measure) is maximized by similar configurations among all three evaluations.

276 To identify optimal method configurations, we set accuracy score minimum
277 thresholds for each evaluation by identifying natural breaks in the range of quality scores,
278 selecting methods and parameter ranges that meet these criteria. Table 2 lists method
279 configurations that maximize species-level classification accuracy scores for mock
280 community, cross-validated, and novel taxa evaluations under several common operating
281 conditions. “Balanced” configurations are recommended for general use, and are methods
282 that maximize F-measure scores. “Precision” and “Recall” configurations maximize
283 precision and recall scores, respectively, for mock, cross-validated, and novel-taxa
284 classifications (Table 2). “Novel” configurations optimize F-measure scores for novel taxa

285 classification, and secondarily for mock and cross-validated performance (Table 2). These
286 configurations are recommended for use with sample types that are expected to contain
287 large proportions of unidentified species, for which overclassification is undesirable.
288 However, these configurations may not perform optimally for classification of known
289 species (i.e., underclassification rates will be higher). For fungi, the same configurations
290 recommended for “Precision” perform well for novel taxa classification (Table 2). For 16S
291 rRNA sequences, BLAST+, UCLUST, and vsearch consensus classifiers perform best for
292 novel taxa classification (Table 2).

293

294 **Computational runtime**

295 High-throughput sequencing platforms (and experiments) continue to yield increasing
296 sequence counts, which — even after quality filtering and dereplication or operational
297 taxonomic unit clustering steps common to most microbiome analysis pipelines — may
298 exceed thousands of unique sequences that need classification. Increasing numbers of
299 query sequences and references sequences may lead to unacceptable runtimes, and under
300 some experimental conditions the top-performing method (based on precision, recall, or
301 some other metric) may be insufficient to handle large numbers of sequences within an
302 acceptable time frame. For example, quick turnarounds may be vital under clinical
303 scenarios as microbiome evaluation becomes common clinical practice, or commercial
304 scenarios, when large sample volumes and client expectations may constrain turnaround
305 times and method selection.

306 We assessed computational runtime as a linear function of 1) the number of query
307 sequences and 2) the number of reference sequences. Linear dependence is empirically
308 evident in Figure 5. For both of these metrics, the slope is the most important measure of
309 performance. The intercept indicates the amount of time taken to train the reference
310 sequences, load environmental variables, or other “setup” steps that will diminish in
311 significance as sequence counts grow, and hence are negligible.

312 UCLUST (0.000028 s/sequence), vsearch (0.000072 s/sequence), BLAST+
313 (0.000080 s/sequence), and legacy BLAST (0.000100 s/sequence) all exhibit shallow
314 slopes with increasing numbers of reference sequences. Naive Bayes (0.000483
315 s/sequence) and SortMeRNA (0.000543 s/sequence) yield moderately higher slopes, and
316 RDP (0.001696 s/sequence) demonstrates the steepest slope (Figure 5A). For runtime as a
317 function of query sequence count, UCLUST (0.002248 s/sequence), RDP (0.002920
318 s/sequence), and SortMeRNA (0.003819 s/sequence) have relatively shallow slopes
319 (Figure 5B). Naive Bayes (0.022984 s/sequence), BLAST+ (0.026222s/sequence) , and
320 vsearch (0.030190 s/sequence) exhibit greater slopes. Legacy BLAST (0.133292
321 s/sequence) yielded a slope magnitudes higher than other methods, rendering this method
322 impractical for large data sets.

323

324

325 Discussion

326 We have developed and validated several machine-learning and alignment-based
327 classifiers provided in q2-feature-classifier and benchmarked these classifiers, as well as
328 other common classification methods, to evaluate their strengths and weaknesses across a
329 range of parameter settings for each (Table 2).

330 Each classifier required some degree of optimization to define top-performing
331 parameter configurations, with the sole exception of QIIME 1's legacy BLAST wrapper,
332 which was unaffected by its only user-defined parameter, e-value, over a range of 10^{-10} to
333 1000. For all other methods, performance varied widely depending on parameter settings,
334 and a single method could achieve among the worst performance with one configuration
335 but among the best performance with another. Configurations greatly affected accuracy
336 with mock community, cross-validated, and novel taxa evaluations, indicating that
337 optimization is necessary under a variety of performance conditions, and optimization for
338 one condition may not necessarily translate to another. Mock community and cross-
339 validated evaluations exhibited similar results, but novel taxa evaluations selected different
340 optimal configurations for most methods (Figure 4), indicating that configurations
341 optimized to one condition, e.g., high-recall classification of known sequences, may be less
342 suited for other conditions, e.g., classification of novel sequences. Table 2 lists the top-
343 performing configuration for each method for several standard performance conditions.

344 Optimal configurations also varied among different evaluation metrics. Precision
345 and recall, in particular, exhibited some mutual opposition, such that methods increasing

346 precision reduced recall. For this reason, F-measure, the harmonic mean of precision and
347 recall, is a useful metric for choosing configurations that are well balanced for average
348 performance. “Balanced” method configurations — which maximize F-measure scores for
349 mock, cross-validated, and novel taxa evaluations (Table 2) — are best suited for a wide
350 range of user conditions. The naive Bayes classifier with kmer lengths of 6 or 7 and
351 confidence = 0.7 (or confidence \geq 0.9 if using bespoke class weights), RDP with confidence
352 = 0.6-0.7, and UCLUST (minimum consensus = 0.51, minimum similarity = 0.9, max accepts
353 = 3) perform best under these conditions (Table 2). Performance is dramatically improved
354 using bespoke class weights for 16S rRNA sequences (Figure 4A-B), though this approach is
355 developmental and only applicable when the expected composition of samples is known in
356 advance (a scenario that is becoming increasingly common with the increasing quantity of
357 public microbiome data, and which could be aided by microbiome data sharing resources
358 such as Qiita (<http://qiita.microbio.me>)). For ITS sequences, the naive Bayes classifier with
359 kmer lengths of 6 or 7 and confidence \geq 0.9, or RDP with confidence = 0.7-0.9, perform best,
360 and the effects of bespoke class weights are less pronounced (Figure 4C-D).

361 However, some users may require high-precision classifiers when false-positives
362 may be more damaging to the outcome, e.g., for detection of pathogens in a sample.
363 Precision scores are maximized by naive Bayes and RDP classifiers with high confidence
364 settings (Table 2). Optimizing for precision will significantly damage recall by yielding a
365 high number of false negatives.

366 Other users may require high-recall classifiers when false-negatives and
367 underclassification hinder interpretation, but false positives (mostly overclassification to a

368 closely related species) are less damaging. For example, in environments with high
369 numbers of unidentified species, a high-precision classifier may yield large numbers of
370 unclassified sequences; in such cases, a second pass with a high-recall configuration (Table
371 2) may provide useful inference of what taxa are most similar to these unclassified
372 sequences. When recall is optimized, precision tends to suffer slightly (leading to similar F-
373 measure scores to “balanced” configurations) but novel taxa classification accuracy is
374 minimized, as these configurations tend to overclassify (Table 2). Any user prioritizing
375 recall ought to be aware of and acknowledge these risks, e.g., when sharing or publishing
376 their results, and understand that many of the species-level classifications may be wrong,
377 particularly if the samples are expected to contain many uncharacterized species. For 16S
378 rRNA sequences, naive Bayes bespoke classifiers with kmer lengths between 12-32 and
379 confidence = 0.5 yield maximal recall scores, but RDP (confidence = 0.5) and naive Bayes
380 (uniform class weights, confidence = 0.5, kmer length = 11, 12, or 18) also perform well
381 (Table 2). Fungal recall scores are maximized by the same configurations recommended for
382 “Balanced” classification, i.e., naive Bayes classifiers with kmer lengths between 6-7 and
383 confidence between 0.92-0.98, or RDP with confidence between 0.7-0.9 (Table 2).

384 Runtime requirements may also be the chief concern dictating method selection for
385 some users. QIIME 1’s UCLUST wrapper provides the fastest runtime while still achieving
386 reasonably good performance for most evaluations; Naive Bayes, RDP, and BLAST+ also
387 delivered reasonably low runtime requirements, and outperform UCLUST on most other
388 evaluation metrics.

389

390 **Conclusions**

391 The classification methods provided in q2-feature-classifier will support improved
392 taxonomy classification of marker-gene sequences, and are released as a free, open-source
393 plugin for use with QIIME 2. We demonstrate that these methods perform as well as or
394 better than other leading taxonomy classification methods on a number of performance
395 metrics. The naive Bayes, vsearch, and BLAST+ consensus classifiers described here are
396 released for the first time in QIIME 2, with optimized “balanced” configurations (Table 2)
397 set as defaults.

398 We also present the results of a benchmark of several widely used taxonomy
399 classifiers, and recommend the top-performing methods and configurations for the most
400 common user scenarios. Our recommendations for “balanced” methods (Table 2) will be
401 appropriate for most users who are classifying 16S rRNA or fungal ITS sequences, but other
402 users may prioritize high-precision (low false-positive) or high-recall (low false-negative)
403 methods.

404 We have also shown that great potential exists for improving the accuracy of
405 taxonomy classifications by appropriately setting class weights for the machine learning
406 classifiers. Currently, no tools exist that allow users to generate appropriate values for
407 these class weights in real applications. Compiling appropriate class weights for different
408 sample types could be a promising approach to further improve taxonomic classification of
409 marker gene sequence reads.

410

411 **Methods**

412 **Mock communities**

413 All mock communities were sourced from mockrobiota [11]. Raw fastq files were
414 demultiplexed and processed using tools available in QIIME 2 (version 2017.4)
415 (<https://qiime2.org/>). Reads were demultiplexed with q2-demux
416 (<https://github.com/qiime2/q2-demux>) and quality filtered and dereplicated with q2-
417 dada2 [4]. Representative sequence sets for each dada2 sequence variant were used for
418 taxonomy classification with each classification method.

419 The inclusion of multiple mock community samples is important to avoid overfitting;
420 optimizing method performance to a small set of data could result in overfitting to the
421 specific community compositions or conditions under which those data were generated,
422 which reduces the robustness of the classifier.

423 **Cross-validated simulated reads**

424 The simulated reads used here were derived from the reference databases using the
425 “Cross-validated classification performance” notebooks in our project repository. The
426 reference databases were either Greengenes or UNITE (99% OTUs) that were cleaned
427 according to taxonomic label to remove sequences with ambiguous or null labels.

20

428 Reference sequences were trimmed to simulate amplification using standard PCR primers
429 and slice out the first 250 bases downstream (3') of the forward primer. The exact
430 sequences were used for cross validation, and were not altered to simulate any sequencing
431 error. The bacterial primers used were 515F/806R [17], and the fungal primers used were
432 BITSf/B58S3r [18]. Each database was stratified by taxonomy and 10-fold randomised
433 cross-validation data sets were generated using scikit-learn's library functions. Where a
434 taxonomic label had less than 10 instances, taxonomies were amalgamated to make
435 sufficiently large strata. If, as a result, a taxonomy in any test set was not present in the
436 corresponding training set, the expected taxonomy label was truncated to the nearest
437 common taxonomic rank observed in the training set (e.g., *Lactobacillus casei* would
438 become *Lactobacillus*). The notebook detailing simulated read generation (for both cross-
439 validated and novel taxa reads) prior to taxonomy classification is available at
440 [https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/novel-taxa/dataset-](https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/novel-taxa/dataset-generation.ipynb)
441 [generation.ipynb](https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/novel-taxa/dataset-generation.ipynb).
442 Classification performance was also slightly modified from a standard machine-learning
443 scenario as the classifiers in this study are able to refuse classification if they are not
444 confident above a taxonomic level for a given sample. This also accommodates the
445 taxonomy truncation that we performed for this test. The methodology was consistent with
446 that used below for novel taxa, but we defer this description to the next section.

447 **“Novel taxa” simulation analysis**

448 “Novel taxa” analysis was performed to test the performance of classifiers when assigning
449 taxonomy to sequences that are not represented in a reference database, e.g., as a
450 simulation of what occurs when a method encounters an undocumented species. In this
451 analysis, simulated amplicons were filtered from those used for the cross-validation
452 analysis. For all sequences present in each test set, sequences sharing taxonomic affiliation
453 at a given taxonomic level L (e.g., to species level) in the corresponding training set were
454 removed. Taxa are stratified among query and test sets such that for each query taxonomy
455 at level L, no reference sequences match that taxonomy, but at least one reference
456 sequence will match the taxonomic lineage at level L-1 (e.g., same genus but different
457 species). An ideal classifier would assign taxonomy to the nearest common taxonomic
458 lineage (e.g., genus), but would not “overclassify” to near neighbors (e.g., assign species-
459 level taxonomy when species X is removed from the reference database). For example, a
460 “novel” sequence representing the species *Lactobacillus brevis* should be classified as
461 “*Lactobacillus*”, without species-level annotation, in order to be considered a true positive
462 in this analysis. As described above for cross-validated reads, these novel taxa simulated
463 communities were also tested in both bacterial (B) and fungal (F) databases on simulated
464 amplicons trimmed to simulate 250-nt sequencing reads.

465 Novel taxa classification performance is evaluated using precision, recall, F-
466 measure, overclassification rates, underclassification rates, and misclassification rates for
467 each taxonomic level (phylum to species), computed with the following modified

468 definitions (see below, *Performance analyses using simulated reads*, for full description of
469 precision, recall, and F-measure calculations; these calculations use the modified
470 definitions of true positive, false positive, and false negative as described here):

471 1) A true positive is considered the nearest correct lineage contained in the reference
472 database. For example, if *Lactobacillus brevis* is removed from the reference
473 database and used as a query sequence, the only correct taxonomy classification
474 would be “*Lactobacillus*”, without species-level classification.

475 2) A false positive would be either an classification to a different *Lactobacillus* species
476 (*Overclassification*), or any genus other than *Lactobacillus* (*Misclassification*).

477 3) A false negative occurs if an expected taxonomy classification (e.g., “*Lactobacillus*”)
478 is not observed in the results. Note that this will be the modified taxonomy expected
479 when using a naive reference database, and is not the same as the true taxonomic
480 affiliation of a query sequence in the novel taxa analysis. A false negative results
481 from misclassification, overclassification, or when the classification contains the
482 correct basal lineage, but does not assign a taxonomy label at level L
483 (*Underclassification*). E.g., classification as “*Lactobacillaceae*”, but no genus-level
484 classification.

485 **Taxonomy classification**

486 Representative sequences for all analyses (mock community, cross-validated, and novel
487 taxa) were classified taxonomically using the following taxonomy classifiers and setting
488 sweeps:

489 1. q2-feature-classifier multinomial naive Bayes classifier. Varied k-mer length
490 in {4, 6, 7, 8, 9, 10, 11, 12, 14, 16, 18, 32} and confidence threshold in {0, 0.5, 0.7, 0.9,
491 0.92, 0.94, 0.96, 0.98, 1}.

492 2. BLAST+ [9] local sequence alignment, followed by consensus taxonomy
493 classification implemented in q2-feature-classifier. Varied max accepts from 1 to 100;
494 percent identity from 0.80 to 0.99; and minimum consensus from 0.51 to 0.99. See
495 description below.

496 3. vsearch [10] global sequence alignment, followed by consensus taxonomy
497 classification implemented in q2-feature-classifier. Varied max accepts from 1 to
498 100; percent identity from 0.80 to 0.99; and minimum consensus from 0.51 to 0.99.
499 See description below.

500 4. Ribosomal Database Project (RDP) naïve Bayesian classifier [12] (QIIME1
501 wrapper), with confidence thresholds between 0.0 to 1.0 in steps of 0.1.

502 5. Legacy BLAST [13] (QIIME1 wrapper) varying e-value thresholds from 1e-9
503 to 1000.

504 6. SortMeRNA [15] (QIIME1 wrapper) varying minimum consensus fraction
505 from 0.51 to 0.99; similarity from 0.8 to 0.9; max accepts from 1 to 10; and coverage
506 from 0.8 to 0.9.

507 7. UCLUST [14] (QIIME1 wrapper) varying minimum consensus fraction from
508 0.51 to 0.99; similarity from 0.8 to 0.9; and max accepts from 1 to 10.

509

510 With the exception of the UCLUST classifier, we have only benchmarked the performance of
511 open-source, free, marker-gene-agnostic classifiers, i.e., those that can be trained/aligned
512 on a reference database of *any* marker gene. Hence, we excluded classifiers that can only
513 assign taxonomy to a particular marker gene (e.g., only bacterial 16S rRNA genes) and
514 those that rely on specialized or unavailable reference databases and cannot be trained on
515 other databases, effectively restricting their use for other marker genes and custom
516 databases.

517 Classification of bacterial/archaeal 16S rRNA sequences was made using the Greengenes
518 reference sequence database (13_8 release) [5] preclustered at 99% ID, with V4 domain
519 amplicons extracted using primers 515f/806r with q2-feature-classifier's `extract_reads`
520 method. Classification of fungal ITS sequences was made using the UNITE database
521 (version 7.1 QIIME developer release) [19] preclustered at 99% ID. For the cross
522 validation and novel taxa tests we prefiltered to remove sequences with incomplete or
523 ambiguous taxonomies (containing the substrings 'unknown', 'unidentified', or '_sp' or
524 terminating at any level with '_').

525
526 The notebooks detailing taxonomy classification sweeps of mock communities are available
527 at <https://github.com/caporaso-lab/tax-credit/tree/0.2.2/ipynb/mock-community>. Cross-
528 validated read classification sweeps are available at <https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/cross-validated/taxonomy-assignment.ipynb>. Novel taxa
529 classification sweeps are available at <https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/novel-taxa/taxonomy-assignment.ipynb>.

532

533 **Runtime analyses**

534 The tax-credit framework employs two different runtime metrics: as a function of 1) the
535 number of query sequences or 2) the number of reference sequences. Taxonomy classifier
536 runtimes were logged while performing classifications of pseudorandom subsets of 1,
537 2,000, 4,000, 6,000, 8,000, and 10,000 sequences from the Greengenes 99% OTU database.
538 Each subset was drawn once then used for all of the tests as appropriate. All runtimes were
539 computed on the same Linux workstation (Ubuntu 16.04.2 LTS, Intel Xeon CPU E7-4850 v3
540 @ 2.20GHz, 1TB memory). The exact commands used for runtime analysis are presented in
541 the “Runtime analyses” notebook in the project repository ([https://github.com/caporaso-](https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/runtime/analysis.ipynb)
542 [lab/tax-credit/blob/0.2.2/ipynb/runtime/analysis.ipynb](https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/runtime/analysis.ipynb)).

543 **Performance analyses using simulated reads**

544 Cross-validated and novel taxa reads are evaluated using the classic precision, recall, and F-
545 measure metrics [5] (novel taxa use the standard calculations as described below, but
546 modified definitions for true positive (*TP*), false positive (*FP*), and false negative (*FN*), as
547 described above for novel taxa).

548 Precision, recall, and F-measure are calculated as follows:

- 549 ○ *Precision* = $TP/(TP+FP)$ or the fraction of sequences that were classified correctly at
550 level L.
- 551 ○ *Recall* = $TP/(TP+FN)$ or the fraction of expected taxonomic labels that were
552 predicted at level L.

553 ○ *F-measure* = $2 \times \text{Precision} \times \text{Recall} / (\text{Precision} + \text{Recall})$, or the harmonic mean of
554 precision and recall.

555 The Jupyter notebook detailing commands used for evaluation of cross-validated read
556 classifications is available at <https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/cross-validated/evaluate-classification.ipynb>. The notebook for
557 evaluation of novel taxa classifications is available at <https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/novel-taxa/evaluate-classification.ipynb>.

560 **Performance analyses using mock communities**

561 The Jupyter notebook detailing commands used for evaluation of mock communities,
562 including the three evaluation types described below, is available at
563 [https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/mock-](https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/mock-community/evaluate-classification-accuracy.ipynb)
564 community/evaluate-classification-accuracy.ipynb.

565 **Precision and Recall**

566 Classic precision, recall, and F-measure are used to calculate mock community
567 classification accuracy, using the definitions given above for simulated reads. These metrics
568 require knowing the expected classification of each sequence, which we determine by
569 performing a gapless alignment between each representative sequence in the mock
570 community and the marker-gene sequences of each microbial strain added to the mock
571 community. These “expected sequences” are provided for the mock communities in
572 mockrobiota [11]. Representative sequences are assigned the taxonomy of the best

573 alignment, and any representative sequence with more than 3 mismatches to the expected
574 sequences are excluded from precision/recall calculations. If a representative sequence
575 aligns to more than one expected sequence equally well, all top hits are accepted as the
576 “correct” classification. This scenario is rare and typically only occurred when different
577 strains of the same species were added to the same mock community to intentionally
578 produce this challenge (e.g., for mock-12 as described by [4]). Precision, recall, and F-
579 measure are then calculated by comparing the “expected” classification for each mock
580 community sequence to the classifications predicted by each taxonomy classifier using the
581 full reference databases, as described above.

582 **Taxon accuracy rate and taxon detection rate**

583 Taxon accuracy rate (TAR) and taxon detection rate (TDR) are used for qualitative
584 compositional analyses of mock communities. As the true taxonomy labels for each
585 sequence in a mock community are not known with absolute certainty, TAR and TDR are
586 useful alternatives to precision and recall that instead rely on the presence/absence of
587 expected taxa, or microbiota that are intentionally added to the mock community. In
588 practice, TAR/TDR are complementary metrics to precision/recall and should provide
589 similar results if the expected classifications for mock community representative
590 sequences are accurate.

591 At a given taxonomic level, a classification is a:

- 592 ○ true positive (*TP*), if that taxon is both observed and expected.

593 ○ false positive (*FP*), if that taxon is observed but not expected.

594 ○ false negative (*FN*), if a taxon is expected but not observed.

595 These are used to calculate TAR and TDR as:

596 ○ $TAR = TP/(TP+FP)$ or the fraction of observed taxa that were expected at level L.

597 ○ $TDR = TP/(TP+FN)$ or the fraction of expected taxa that are observed at level L.

598

599 **Bray-Curtis Dissimilarity**

600 Bray-Curtis dissimilarity [20] is used to measure the degree of dissimilarity between two
601 samples as a function of the abundance of each species label present in each sample,
602 treating each species as equally related. This is a useful metric for evaluating classifier
603 performance by assessing the relative distance between each predicted mock community
604 composition (abundance of taxa in a sample based on results of a single classifier) and the
605 expected composition of that sample. For each classifier, Bray-Curtis distances between the
606 expected and observed taxonomic compositions are calculated for each sample in each
607 mock community dataset; this yields a single expected-observed distance for each
608 individual observation. The distance distributions for each method are then compared
609 statistically using paired or unpaired t-tests to assess whether one method (or
610 configuration) performs consistently better than another.

611 **New taxonomy classifiers**

612 We describe q2-feature-classifier (<https://github.com/qiime2/q2-feature-classifier>), a
613 plugin for QIIME 2 (<https://qiime2.org/>) that performs multi-class taxonomy classification
614 of marker-gene sequence reads. In this work we compare the consensus BLAST+ and
615 vsearch methods and the naive Bayes scikit-learn classifier. The software is free and open-
616 source.

617 **Machine learning taxonomy classifiers**

618 The q2-feature-classifier plugin allows users to apply any of the suite of machine learning
619 classifiers available in scikit-learn (<http://scikit-learn.org>) to the problem of taxonomy
620 classification of marker-gene sequences. It functions as a lightweight wrapper that
621 transforms the problem into a standard document classification problem. Advanced users
622 can input any appropriate scikit-learn classifier pipeline, which can include a range of
623 feature extraction and transformation steps as well as specifying a machine learning
624 algorithm.

625

626 The plugin provides a default method which is to extract k-mer counts from reference
627 sequences and train the scikit-learn multinomial naive Bayes classifier, and it is this
628 method that we test extensively here. Specifically, the pipeline consists of a
629 `sklearn.feature_extraction.text.HashingVectorizer` feature extraction step followed by a
630 `sklearn.naive_bayes.MultinomialNB` classification step. The use of a hashing feature
631 extractor allows the use of significantly longer k-mers than the 8-mers that are used by

632 RDP Classifier, and we tested up to 32-mers. Like most scikit-learn classifiers, we are able
633 to set class weights when training the multinomial naive Bayes classifiers. In the naive
634 Bayes setting, setting class weights means that class priors are not derived from the
635 training data or set to be uniform, as they are for the RDP Classifier. For more detail on how
636 class weights enter the calculations please refer to the scikit-learn User Guide
637 (<http://scikit-learn.org>).

638
639 In most settings, it is highly unlikely that the assumption of uniform weights is correct. That
640 assumption is that each of the taxa in the reference database is equally likely to appear in
641 each sample. Setting class weights to more realistic values can greatly aid the classifier in
642 making more accurate predictions, as we show in this work. When testing the mock
643 communities we made use of the fact that the sequence compositions were known *a priori*
644 for the bespoke classifier. For the simulated reads studies, we allowed the classifier to set
645 the class weights from the class frequencies observed in each training set for the bespoke
646 classifier.

647
648 For this study, we performed two parameter sweeps on the mock communities: an initial
649 broad sweep to optimize feature extraction parameters and then a more focussed sweep to
650 optimise k-mer length and confidence parameter settings. These sweeps included varying
651 the assumptions regarding class weights. The focussed sweeps were also performed for the
652 cross-validated and novel taxa evaluations, but only for the assumption of uniform class

653 priors. The results for the focussed sweeps across all data sets are those which are
654 compared against the other classifiers in this work.

655

656 The broad sweeps used a modified scikit-learn pipeline which consisted of the
657 `sklearn.feature_extraction.text.HashingVectorizer`, followed by the
658 `sklearn.feature_extraction.text.TfidfTransformer`, then the
659 `sklearn.naive_bayes.MultinomialNB`. We performed a full grid search over the parameters
660 shown in Table 3. The conclusion from the initial sweep was that the `TfidfTransformer` step
661 did not significantly improve classification, that `n_features` should be set to 8192, feature
662 vectors should be normalised using L2 normalisation and that the alpha parameter for the
663 naive Bayes classifier should be set to 0.001. Please see [https://github.com/caporaso-](https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/mock-community/evaluate-classification-accuracy-nb-extra.ipynb)
664 [lab/tax-credit/blob/0.2.2/ipynb/mock-community/evaluate-classification-accuracy-nb-](https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/mock-community/evaluate-classification-accuracy-nb-extra.ipynb)
665 [extra.ipynb](https://github.com/caporaso-lab/tax-credit/blob/0.2.2/ipynb/mock-community/evaluate-classification-accuracy-nb-extra.ipynb) for details.

666 **Consensus taxonomy alignment-based classifiers**

667

668 Two new classifiers implemented in `q2-feature-classifier` perform consensus taxonomy
669 classification based on alignment of a query sequence to a reference sequence. The
670 methods `classify_consensus_vsearch` and `classify_consensus_blast` use the global aligner
671 `vsearch` [10] or the local aligner `BLAST+` [9], respectively, to return up to `maxaccepts`
672 reference sequences that align to the query with at least `perc_identity` similarity. A
673 consensus taxonomy is then assigned to the query sequence by determining the taxonomic

674 lineage on which at least `min_consensus` of the aligned sequences agree. This consensus
675 taxonomy is truncated at the taxonomic level at which less than `min_consensus` of
676 taxonomies agree. For example, if a query sequence is classified with `maxaccepts=3`,
677 `min_consensus=0.51`, and the following top hits:

678

679 `k__Bacteria; p__Firmicutes; c__Bacilli; o__Lactobacillales; f__Lactobacillaceae;`
680 `g__Lactobacillus; s__brevis`

681 `k__Bacteria; p__Firmicutes; c__Bacilli; o__Lactobacillales; f__Lactobacillaceae;`
682 `g__Lactobacillus; s__brevis`

683 `k__Bacteria; p__Firmicutes; c__Bacilli; o__Lactobacillales; f__Lactobacillaceae;`
684 `g__Lactobacillus; s__delbrueckii`

685

686 The taxonomy label assigned will be `k__Bacteria; p__Firmicutes; c__Bacilli;`

687 `o__Lactobacillales; f__Lactobacillaceae; g__Lactobacillus; s__brevis`. However, if

688 `min_consensus=0.99`, the taxonomy label assigned will be `k__Bacteria; p__Firmicutes;`

689 `c__Bacilli; o__Lactobacillales; f__Lactobacillaceae; g__Lactobacillus`.

690

691

692 **Declarations**

693 **Ethics approval and consent to participate**

694 Not applicable

695 **Consent for publication**

696 Not applicable

697 **Availability of data and materials**

698 Mock community sequence data used in this study are publicly available in mockrobiota
699 [11] under the study identities listed in Table 1. All other data generated in this study, and
700 all new software, is available in our GitHub repositories under the BSD license. The tax-
701 credit repository can be found at: <https://github.com/caporaso-lab/tax-credit>, and static
702 versions of all analysis notebooks, which contain all code and analysis results, can be
703 viewed there. The q2-feature-classifier repository can be accessed at
704 <https://github.com/qiime2/q2-feature-classifier>; as a QIIME2 core plugin, it is
705 automatically installed any time QIIME2 (<https://qiime2.org/>) is installed.

706

707 **Project name:** q2-feature-classifier

708 **Project home page:** <https://github.com/qiime2/q2-feature-classifier>

709 **Operating system(s):** macOS X, Linux

710 **Programming language:** Python

711 **Other requirements:** QIIME2

712 **License:** BSD-3-Clause

713 **Any restrictions to use by non-academics:** None

714

715 **Project name:** tax-credit

716 **Project home page:** <https://github.com/caporaso-lab/tax-credit>
717 **Operating system(s):** macOS X, Linux
718 **Programming language:** Python
719 **Other requirements:** None (QIIME2 required for some optional functions)
720 **License:** BSD-3-Clause
721 **Any restrictions to use by non-academics:** None
722
723

724 **Funding**

725 This work was funded in part by National Science Foundation award 1565100 to JGC and
726 RK, awards from the Alfred P. Sloan Foundation to JGC and RK, and National Health and
727 Medical Research Council of Australia award APP1085372 to GAH and JGC. These funding
728 bodies had no role in the design of the study, the collection, analysis, or interpretation of
729 data, or in writing the manuscript.

730 **Acknowledgments**

731 The authors thank Stephen Gould and Cheng Soon Ong for advice on machine learning
732 optimisation.

733 **Authors' Contributions**

734 NAB, RK, and JGC conceived and designed tax-credit. NAB, BDK, JGC, and JRR contributed
735 to tax-credit. BDK, MD, JGC, and NAB contributed to q2-feature-classifier. BDK, JGC, MD,
736 JRR, and EB provided QIIME 2 integration with q2-feature-classifier. JGC and GAH provided
737 materials and support. NAB, BDK, JGC, and GAH wrote the manuscript with input from all
738 co-authors.

739 **Competing Interests**

740 The authors declare that they have no competing interests.

741

742 **Tables and Figures**

743 Table 1. Mock communities currently integrated in tax-credit.

Study ID*	target-gene	Species	Strains	Citation
mock-1	16S	46	48	[21]
mock-2	16S	46	48	[21]
mock-3	16S	21	21	[21]
mock-4	16S	21	21	[21]
mock-5	16S	21	21	[21]
mock-7	16S	67	67	[22]
mock-8	16S	67	67	[11]
mock-9	ITS	13	16	[11]
mock-10	ITS	13	16	[11]
mock-12	16S	26	27	[4]
mock-16	16S	56	59	[23]
mock-18	16S	15	15	[24]
mock-19	16S	15	27	[24]
mock-20	16S	20	20	[25]
mock-21	16S	20	20	[25]
mock-22	16S	20	20	[25]
mock-23	16S	20	20	[25]
mock-24	ITS	8	8	[26]
mock-26	ITS	11	11	[27]

744 *All studies are available on mockrobiota [11] at <https://github.com/caporaso->

745 [lab/mockrobiota/tree/master/data/\[studyID\]](https://github.com/caporaso-lab/mockrobiota/tree/master/data/[studyID])

746

747 Table 2. Optimized methods configurations for standard operating conditions.

Target	Condition	Method	Parameters	Mock			Cross-validated			Novel taxa			Threshold
				F	P	R	F	P	R	F	P	R	
16S rRNA	Balanced	NB-bespoke	[6,6]:0.9	0.705	0.98	0.582	0.827	0.931	0.744	0.165	0.243	0.125	F = (0.49, 0.8, 0.1)
			[6,6]:0.92	0.705	0.98	0.581	0.825	0.936	0.737	0.165	0.251	0.123	F = (0.7, 0.8, 0.15)
			[6,6]:0.94	0.703	0.98	0.579	0.822	0.942	0.729	0.162	0.259	0.118	
			[7,7]:0.92	0.712	0.978	0.592	0.831	0.931	0.751	0.151	0.221	0.115	
			[7,7]:0.94	0.708	0.978	0.586	0.829	0.936	0.743	0.157	0.239	0.117	
	naive-bayes	[7,7]:0.7	0.495	0.797	0.38	0.819	0.886	0.761	0.115	0.138	0.099		
		rdp	0.6	0.564	0.798	0.457	0.815	0.868	0.768	0.102	0.128	0.084	
			0.7	0.55	0.799	0.438	0.812	0.892	0.746	0.124	0.173	0.096	
	uclust	0.51:0.9:3	0.498	0.746	0.392	0.846	0.876	0.817	0.154	0.201	0.126		
	Precision	NB-bespoke	[6,6]:0.98	0.676	0.987	0.537	0.803	0.956	0.692	0.163	0.303	0.111	P = (0.94, 0.95, 0.25)
			[7,7]:0.98	0.687	0.98	0.551	0.815	0.951	0.713	0.164	0.283	0.115	
			rdp	1	0.239	0.941	0.16	0.632	0.968	0.469	0.12	0.457	0.069
	Recall	NB-bespoke	[12,12]:0.5	0.754	0.8	0.721	0.815	0.83	0.801	0.053	0.058	0.049	R = (0.47, 0.75, 0.04)
			[14,14]:0.5	0.758	0.802	0.726	0.811	0.826	0.797	0.052	0.057	0.048	R = (0.7, 0.75, 0.04)
			[16,16]:0.5	0.755	0.785	0.732	0.808	0.825	0.792	0.052	0.058	0.047	
			[18,18]:0.5	0.772	0.803	0.748	0.805	0.823	0.789	0.055	0.061	0.05	
			[32,32]:0.5	0.937	0.966	0.913	0.788	0.818	0.76	0.054	0.067	0.045	
		naive-bayes	[11,11]:0.5	0.567	0.77	0.479	0.793	0.82	0.768	0.059	0.065	0.055	
			[12,12]:0.5	0.567	0.769	0.479	0.79	0.816	0.765	0.059	0.064	0.055	
		[18,18]:0.5	0.564	0.764	0.477	0.779	0.807	0.753	0.057	0.063	0.051		
	rdp	0.5	0.577	0.791	0.48	0.816	0.848	0.787	0.068	0.079	0.06		
Novel	blast+	10:0.51:0.8	0.436	0.723	0.325	0.816	0.896	0.749	0.225	0.332	0.171	F = (0.4, 0.8, 0.2)	
		uclust	0.76:0.9:5	0.467	0.775	0.348	0.84	0.938	0.76	0.219	0.358	0.158	
		vsearch	10:0.51:0.8	0.45	0.74	0.342	0.814	0.891	0.75	0.226	0.333	0.171	
			10:0.51:0.9	0.45	0.74	0.342	0.82	0.896	0.755	0.219	0.338	0.162	
Fungi	Balanced	naive-bayes	[6,6]:0.94	0.874	0.935	0.827	0.481	0.57	0.416	0.374	0.438	0.327	F = (0.85, 0.45, 0.37)

			[6,6]:0.96	0.874	0.935	0.827	0.495	0.597	0.423	0.399	0.473	0.344	
			[6,6]:0.98	0.874	0.935	0.827	0.505	0.629	0.423	0.426	0.52	0.361	
			[7,7]:0.98	0.874	0.935	0.827	0.485	0.596	0.409	0.388	0.47	0.33	
		NB-bespoke	[6,6]:0.94	0.928	0.968	0.915	0.48	0.567	0.416	0.371	0.433	0.325	
			[6,6]:0.96	0.928	0.968	0.915	0.491	0.59	0.42	0.393	0.466	0.34	
			[6,6]:0.98	0.927	0.97	0.913	0.504	0.624	0.422	0.421	0.512	0.358	
			[7,7]:0.98	0.935	0.97	0.921	0.487	0.596	0.412	0.386	0.466	0.329	
		rdp	0.7	0.929	0.939	0.922	0.479	0.572	0.413	0.382	0.451	0.332	
			0.8	0.924	0.939	0.915	0.507	0.633	0.422	0.434	0.534	0.366	
			0.9	0.922	0.937	0.913	0.517	0.698	0.411	0.47	0.617	0.379	
	Precision	naive-bayes	[6,6]:0.98	0.874	0.935	0.827	0.505	0.629	0.423	0.426	0.52	0.361	P = (0.92, 0.6, 0.3)
		NB-bespoke	[6,6]:0.98	0.927	0.97	0.913	0.504	0.624	0.422	0.421	0.512	0.358	
		rdp	0.8	0.924	0.939	0.915	0.507	0.633	0.422	0.434	0.534	0.366	
			0.9	0.922	0.937	0.913	0.517	0.698	0.411	0.47	0.617	0.379	
			1	0.821	0.943	0.742	0.461	0.81	0.322	0.459	0.774	0.327	
	Recall	NB-bespoke	[6,6]:0.92	0.938	0.971	0.924	0.467	0.544	0.409	0.353	0.407	0.312	R = (0.9, 0.4, 0.3)
			[6,6]:0.94	0.928	0.968	0.915	0.48	0.567	0.416	0.371	0.433	0.325	
			[6,6]:0.96	0.928	0.968	0.915	0.491	0.59	0.42	0.393	0.466	0.34	
			[6,6]:0.98	0.927	0.97	0.913	0.504	0.624	0.422	0.421	0.512	0.358	
			[7,7]:0.96	0.935	0.969	0.921	0.47	0.56	0.404	0.357	0.422	0.31	
			[7,7]:0.98	0.935	0.97	0.921	0.487	0.596	0.412	0.386	0.466	0.329	
		rdp	0.7	0.929	0.939	0.922	0.479	0.572	0.413	0.382	0.451	0.332	
			0.8	0.924	0.939	0.915	0.507	0.633	0.422	0.434	0.534	0.366	
			0.9	0.922	0.937	0.913	0.517	0.698	0.411	0.47	0.617	0.379	
	Novel	naive-bayes	[6,6]:0.98	0.874	0.935	0.827	0.505	0.629	0.423	0.426	0.52	0.361	F = (0.85, 0.45, 0.4)
		NB-bespoke	[6,6]:0.98	0.927	0.97	0.913	0.504	0.624	0.422	0.421	0.512	0.358	
		rdp	0.8	0.923	0.939	0.915	0.507	0.633	0.422	0.434	0.534	0.366	
			0.9	0.921	0.937	0.913	0.517	0.698	0.411	0.47	0.617	0.379	

748

749 ^aF = F-measure, P = precision, R = recall750 ^bNaive Bayes parameters: k-mer range, confidence751 ^cRDP parameters: confidence752 ^dBLAST+/vsearch parameters: max accepts, minimum consensus, minimum percent identity753 ^eUCLUST parameters: minimum consensus, similarity, max accepts

754

38

755 Threshold describes the score cutoffs used to define optimal method ranges, in the format:
 756 [metric = (mock score, cross-validated score, novel-taxa score)]. If two cutoffs are given,
 757 the second indicates a higher cutoff used to select parameters for the developmental NB-
 758 bespoke method, and the configurations listed are the union of the two cutoffs: the second
 759 cutoff for selecting NB-bespoke, the first for selecting all other methods.
 760

761

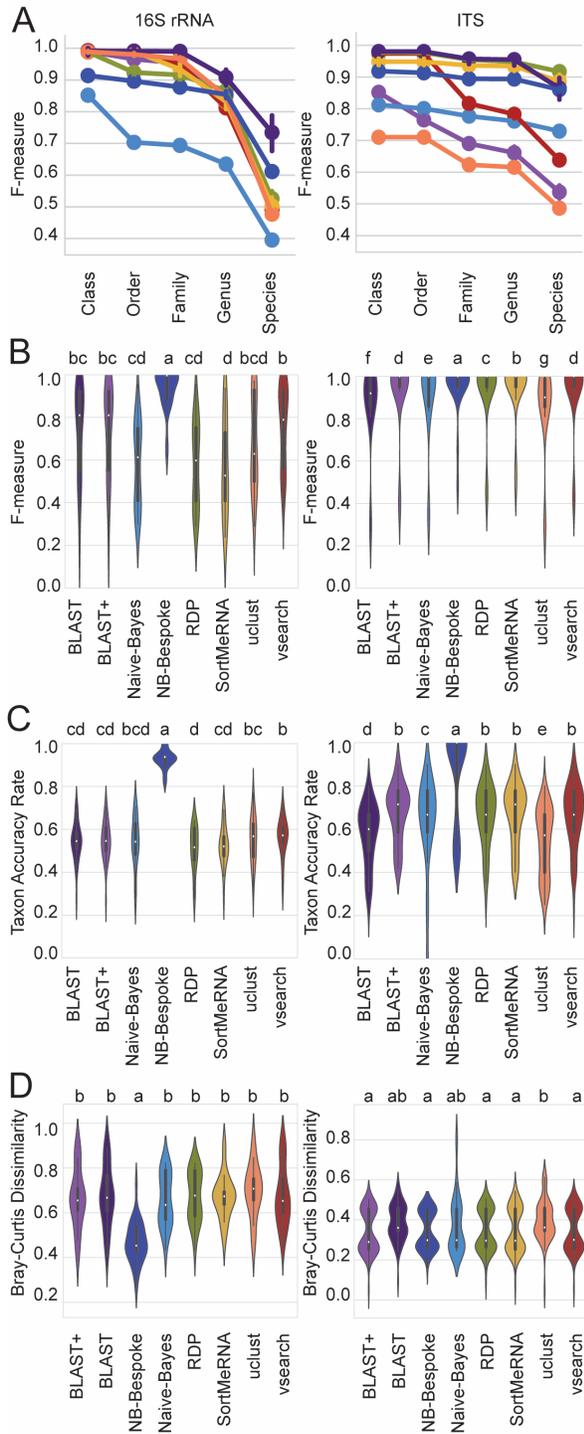
762 Table 3. Naive Bayes broad grid search parameters

Step	Parameter	Values
sklearn.feature_extraction.text.HashingVectorizer	n_features	1024, 8192, 65536
	ngram_range	[4,4], [8, 8], [16, 16], [4,16]
sklearn.feature_extraction.text.TfidfTransformer	norm	l1', 'l2', None
	use_idf	True, False
sklearn.naive_bayes.MultinomialNB	alpha	0.001, 0.01, 0.1
	class_prior	None, array of class weights
post processing	confidence	0, 0.2, 0.4, 0.6, 0.8

763

764

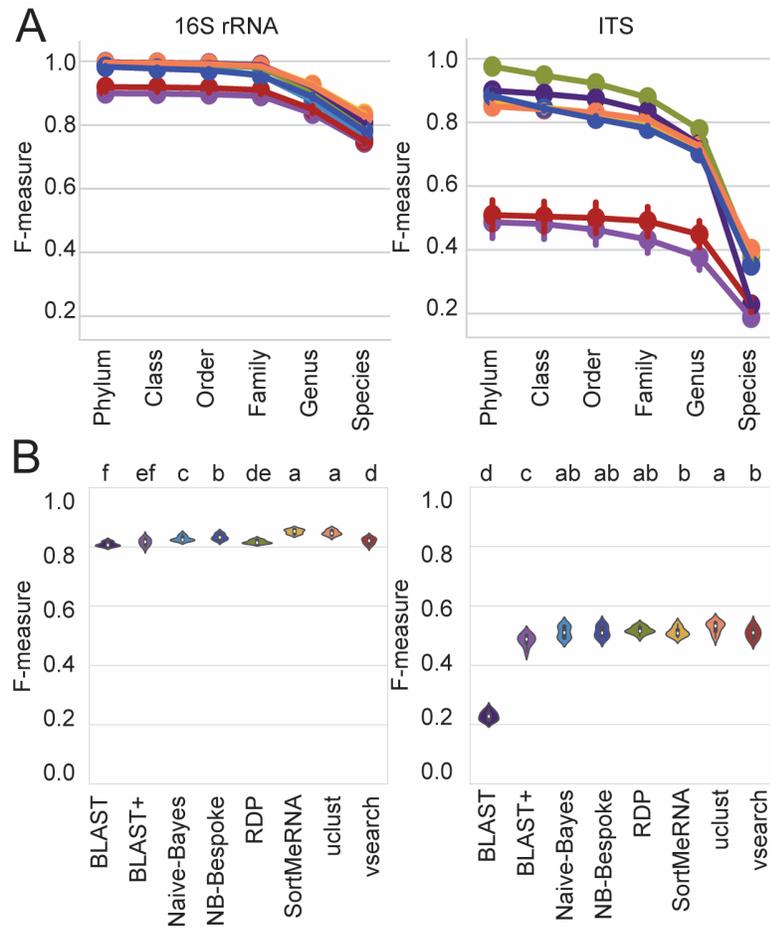
765



766

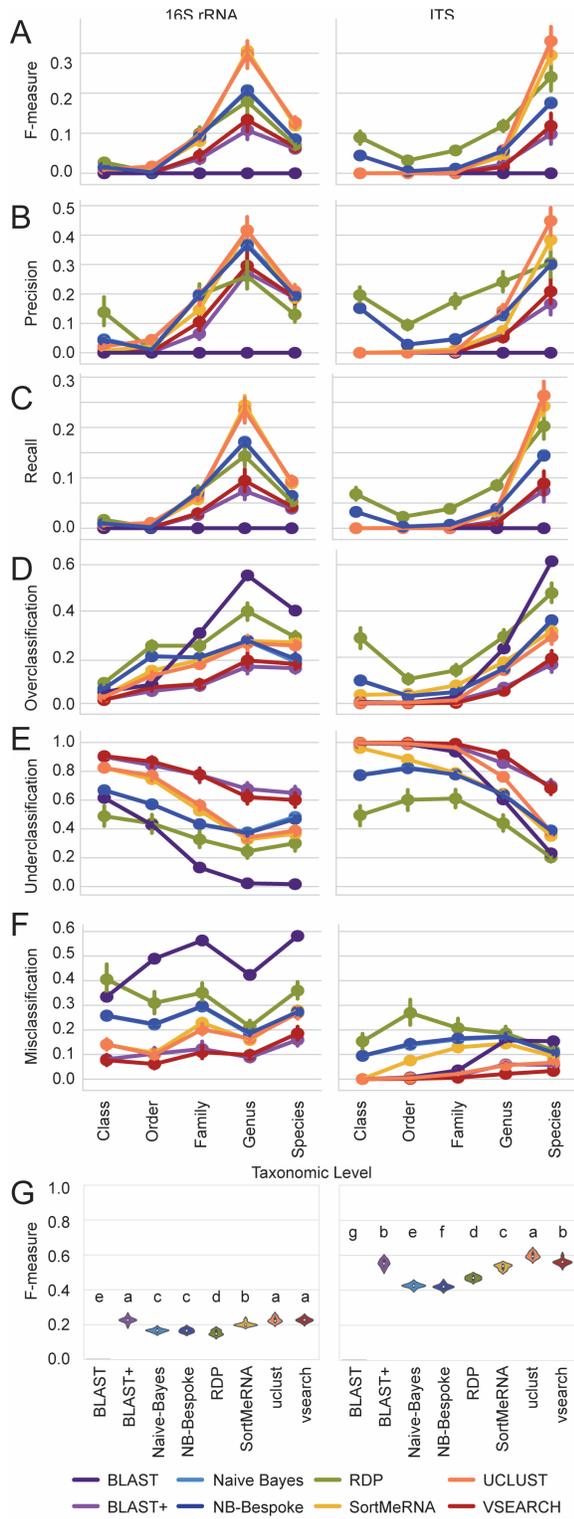
767

768 Figure 1. Classifier performance on mock community datasets for 16S rRNA sequences (left
769 column) and fungal ITS sequences (right column). A, Average F-measure for each taxonomy
770 classification method (averaged across all configurations and all mock community
771 datasets) from class to species level. Error bars = 95% confidence intervals. B, Average F-
772 measure for each optimized classifier (averaged across all mock communities) at species
773 level. C, Average taxon accuracy rate for each optimized classifier (averaged across all mock
774 communities) at species level. D, Average Bray-Curtis distance between the expected mock
775 community composition and its composition as predicted by each optimized classifier
776 (averaged across all mock communities) at species level. Violin plots show median (white
777 point), quartiles (black bars), and kernel density estimation (violin) for each score
778 distribution. Violins with different lower-case letters have significantly different means
779 (paired t-test false detection rate-corrected $P < 0.05$).



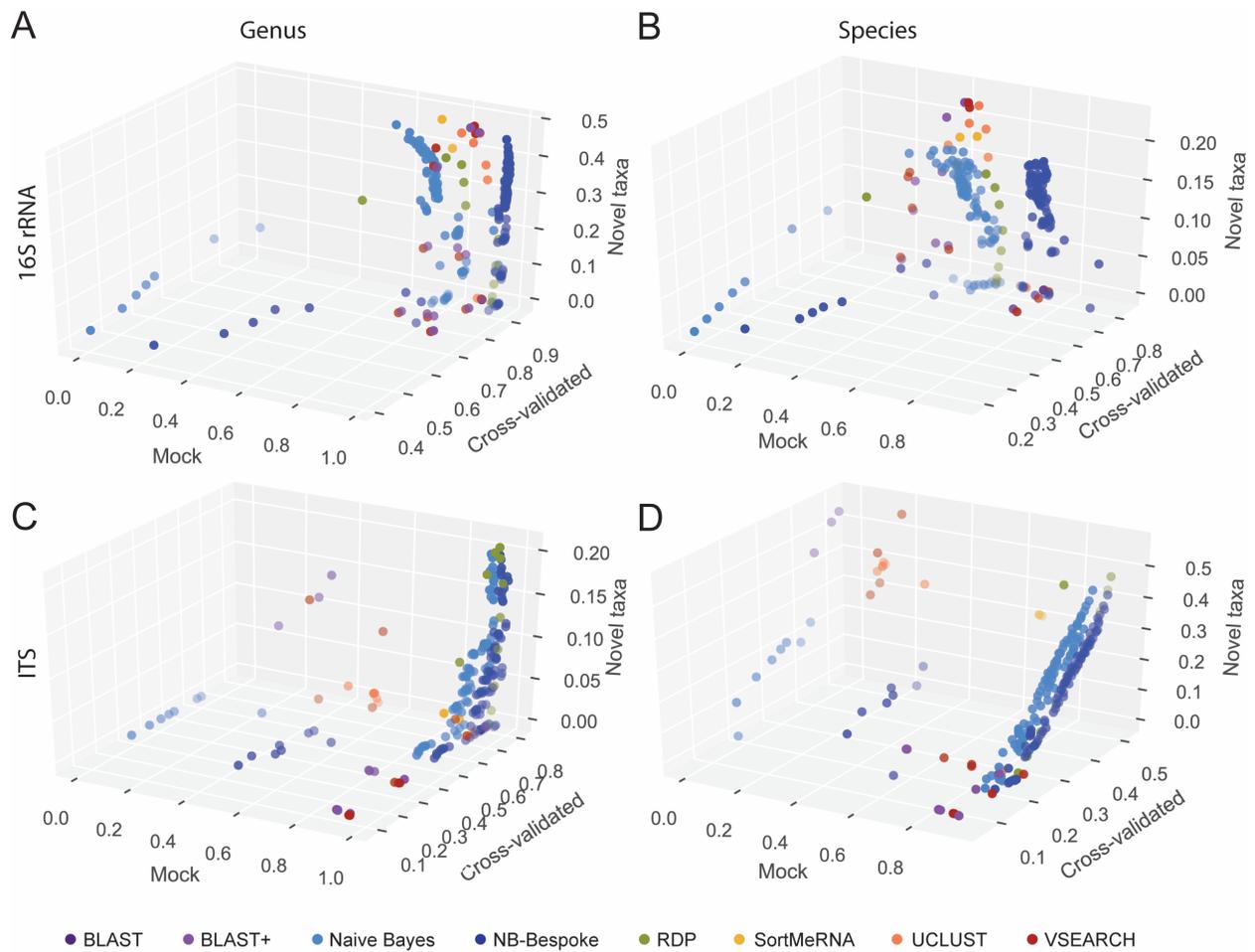
780

781 Figure 2. Classifier performance on cross-validated sequence datasets for 16S rRNA
 782 sequences (left column) and fungal ITS sequences (right column). A, Average F-measure for
 783 each taxonomy classification method (averaged across all configurations and all cross-
 784 validated sequence datasets) from class to species level. Error bars = 95% confidence
 785 intervals. B, Average F-measure for each optimized classifier (averaged across all cross-
 786 validated sequence datasets) at species level. Violins with different lower-case letters have
 787 significantly different means (paired t-test false detection rate-corrected $P < 0.05$).



788

789 Figure 3. Classifier performance on novel-taxa simulated sequence datasets for 16S rRNA
790 sequences (left column) and fungal ITS sequences (right column). A-F, Average F-measure
791 (A), precision (B), recall (C), overclassification (D), underclassification (E), and
792 misclassification (F) for each taxonomy classification method (averaged across all
793 configurations and all novel taxa sequence datasets) from phylum to species level. Error
794 bars = 95% confidence intervals. B, Average F-measure for each optimized classifier
795 (averaged across all novel taxa sequence datasets) at species level. Violins with different
796 lower-case letters have significantly different means (paired t-test false detection rate-
797 corrected $P < 0.05$).



798

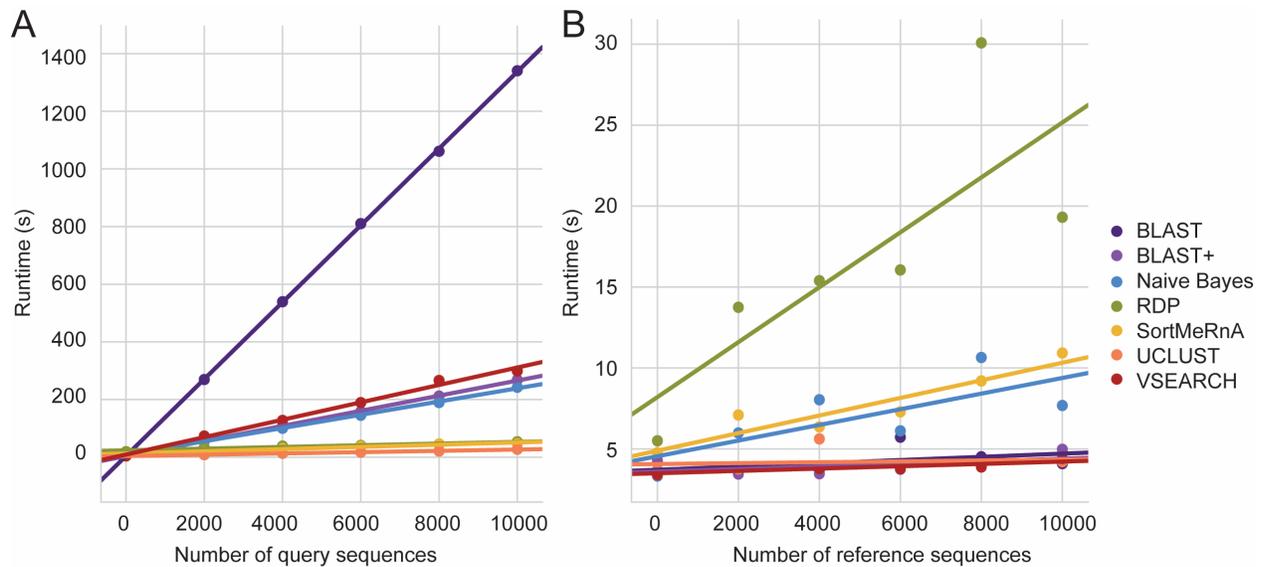
799

800

801

802

Figure 4. Classification accuracy comparison between mock community, cross-validated, and novel taxa evaluations. Scatterplots show mean F-measure scores for each method configuration, averaged across all samples, for classification of 16S rRNA at genus level (A) and species level (B), and fungal ITS sequences at genus level (C) and species level (D).



803

804 Figure 5. Runtime performance comparison of taxonomy classifiers. Runtime (s) for each
 805 taxonomy classifier either varying the number of query sequences and keeping a constant
 806 10000 reference sequences (A) or varying the number of reference sequences and keeping
 807 a constant 1 query sequence (B).

808

809

810 References

- 811 1. Human Microbiome Project Consortium. A framework for human microbiome research.
 812 Nature. 2012 Jun 13;486(7402):215–21.
- 813 2. Gilbert JA, Jansson JK, Knight R. The Earth Microbiome project: successes and aspirations. BMC
 814 Biol [Internet]. 2014;12(1). Available from: <http://dx.doi.org/10.1186/s12915-014-0069-1>
- 815 3. Wang Q, Quensen JF 3rd, Fish JA, Lee TK, Sun Y, Tiedje JM, et al. Ecological patterns of nifH
 816 genes in four terrestrial climatic zones explored with targeted metagenomics using FrameBot,

46

- 817 a new informatics tool. *MBio*. 2013 Sep 17;4(5):e00592–13.
- 818 4. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-resolution
819 sample inference from Illumina amplicon data. *Nat Methods*. 2016 Jul;13(7):581–3.
- 820 5. McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, et al. An improved
821 Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria
822 and archaea. *ISME J*. 2012 Mar;6(3):610–8.
- 823 6. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al. QIIME
824 allows analysis of high-throughput community sequencing data. *Nat Methods*. 2010
825 May;7(5):335–6.
- 826 7. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B, Grisel, O., Blondel, M.,
827 Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M.,
828 Perrot, M., Duchesnay, E. Scikit-learn: Machine Learning in Python. *J Mach Learn Res*.
829 2011;12:2825–30.
- 830 8. Buitinck, L., Louppe, G., Blondel, M., Pedregosa, F., Mueller, A., Grisel, O., Niculae, V.,
831 Prettenhofer, P., Gramfort, A., Grobler, J., Layton, R., VanderPlas, J., Joly, A., Holt, B., Varoquaux
832 G. API design for machine learning software: experiences from the scikit-learn project. In:
833 ECML PKDD Workshop: Languages for Data Mining and Machine Learning. 2013. p. 108–22.
- 834 9. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST : architecture
835 and applications. *BMC Bioinformatics*. 2009;10(1):421.
- 836 10. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool for
837 metagenomics. *PeerJ*. 2016 Oct 18;4:e2584.
- 838 11. Bokulich NA, Rideout JR, Mercurio WG, Shiffer A, Wolfe B, Maurice CF, et al. mockrobiota: a
839 Public Resource for Microbiome Bioinformatics Benchmarking. *mSystems* [Internet]. 2016
840 Sep;1(5). Available from: <http://dx.doi.org/10.1128/mSystems.00062-16>
- 841 12. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian classifier for rapid assignment of rRNA
842 sequences into the new bacterial taxonomy. *Appl Environ Microbiol*. 2007 Aug;73(16):5261–7.
- 843 13. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol*
844 *Biol*. 1990 Oct 5;215(3):403–10.
- 845 14. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*. 2010
846 Oct 1;26(19):2460–1.
- 847 15. Kopylova E, Noé L, Touzet H. SortMeRNA: fast and accurate filtering of ribosomal RNAs in
848 metatranscriptomic data. *Bioinformatics*. 2012 Dec 15;28(24):3211–7.
- 849 16. Müller AC, Behnke S. pystruct - Learning Structured Prediction in Python. *J Mach Learn Res*.
850 2014;15:2055–60.
- 851 17. Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, et al. Ultra-high-

- 852 throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. *ISME J.*
853 2012 Aug;6(8):1621–4.
- 854 18. Bokulich NA, Mills DA. Improved Selection of Internal Transcribed Spacer-Specific Primers
855 Enables Quantitative, Ultra-High-Throughput Profiling of Fungal Communities. *Appl Environ*
856 *Microbiol.* 2013;79(8):2519–26.
- 857 19. Kõljalg U, Nilsson RH, Abarenkov K, Tedersoo L, Taylor AFS, Bahram M, et al. Towards a unified
858 paradigm for sequence-based identification of fungi. *Mol Ecol.* 2013 Nov;22(21):5271–7.
- 859 20. Bray JR, Curtis JT. An Ordination of the Upland Forest Communities of Southern Wisconsin.
860 *Ecol Monogr.* 1957;27(4):325–49.
- 861 21. Bokulich NA, Subramanian S, Faith JJ, Gevers D, Gordon JI, Knight R, et al. Quality-filtering
862 vastly improves diversity estimates from Illumina amplicon sequencing. *Nat Methods.* 2013
863 Jan;10(1):57–9.
- 864 22. Maurice CF, Haiser HJ, Turnbaugh PJ. Xenobiotics Shape the Physiology and Gene Expression of
865 the Active Human Gut Microbiome. *Cell.* 2013;152(1-2):39–50.
- 866 23. Schirmer M, Ijaz UZ, D'Amore R, Hall N, Sloan WT, Quince C. Insight into biases and sequencing
867 errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic Acids Res.* 2015 Mar
868 31;43(6):e37.
- 869 24. Tourlousse DM, Yoshiike S, Ohashi A, Matsukura S, Noda N, Sekiguchi Y. Synthetic spike-in
870 standards for high-throughput 16S rRNA gene amplicon sequencing. *Nucleic Acids Res.*
871 2016;gkw984.
- 872 25. Gohl DM, Vangay P, Garbe J, MacLean A, Hauge A, Becker A, et al. Systematic improvement of
873 amplicon marker gene methods for increased accuracy in microbiome studies. *Nat Biotechnol.*
874 2016 Sep;34(9):942–9.
- 875 26. Taylor DL, Walters WA, Lennon NJ, Bochicchio J, Krohn A, Caporaso JG, et al. Accurate
876 Estimation of Fungal Diversity and Abundance through Improved Lineage-Specific Primers
877 Optimized for Illumina Amplicon Sequencing. *Appl Environ Microbiol.* 2016 Dec
878 15;82(24):7217–26.
- 879 27. Ihrmark K, Bödeker ITM, Cruz-Martinez K, Friberg H, Kubartova A, Schenck J, et al. New
880 primers to amplify the fungal ITS2 region--evaluation by 454-sequencing of artificial and
881 natural communities. *FEMS Microbiol Ecol.* 2012 Dec;82(3):666–77.
- 882