

1 Optimizing taxonomic classification of marker gene
2 amplicon 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 and 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 the
52 accuracy of existing methods for marker-gene amplicon sequence classification. We
53 evaluated and optimized several commonly used taxonomic classification methods (RDP,
54 BLAST, UCLUST) and several new methods (a scikit-learn naive Bayes machine-learning
55 classifier, and alignment-based taxonomy consensus methods of VSEARCH, BLAST+, and
56 SortMeRNA) for classification of marker-gene amplicon sequence data. **Conclusions:** Our

57 results illustrate the importance of parameter tuning for optimizing classifier performance,
58 and we make recommendations regarding parameter choices for a range of standard
59 operating conditions. q2-feature-classifier and our evaluation framework, tax-credit, are
60 both free, open-source, BSD-licensed packages available on GitHub.

61

62 **Background**

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

73 Challenges in this process include the short length of typical sequencing reads with
74 current technology, sequencing and PCR errors [4], selection of appropriate marker genes
75 that contain sufficient heterogeneity to differentiate target species but that are

76 homogeneous enough in some regions to design broad-spectrum primers, quality of
77 reference sequence annotations [5], and selection of a method that accurately predicts the
78 taxonomic affiliation of millions of sequences at low computational cost. Numerous
79 methods have been developed for taxonomy classification of DNA sequences, but few have
80 been directly compared in the specific case of short marker-gene sequences.

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

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

98 computed results from previously evaluated methods, including the results presented here,
99 and provides a framework for parameter sweeps and method optimization. tax-credit
100 could be used as an evaluation framework by other research groups in the future, or its raw
101 data could be easily extracted for integration in another evaluation framework.

102

103 **Results**

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

115 **Mock community evaluations**

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

126 An additional priority was to test the effect of setting class weights on classification
127 accuracy for the naive Bayes classifier implemented in q2-feature-classifier. In machine
128 learning, class weights or prior probabilities are vectors of weights that specify the
129 frequency at which each class is expected to be observed (and should be distinguished
130 from the use of this term under Bayesian inference as a probability distribution of weights
131 vectors). An alternative to setting class weights is to assume that each query sequence is
132 equally likely to belong to any of the taxa that are present in the reference sequence
133 database. This assumption, known as uniform class priors in the context of a naive Bayes
134 classifier, is made by the RDP classifier [12], and its impact on marker-gene classification
135 accuracy has yet to be validated. Making either assumption, that the class weights are

136 uniform or known to some extent, will affect results and cannot be avoided. The mock
137 communities have taxonomic abundances that are far from uniform over the set of
138 reference taxonomies, as any real data set must. We can therefore use them to assess the
139 impact of making assumptions regarding class weights. Where we have set the class
140 weights to the known taxonomic composition of a sample, we have labelled the results
141 “bespoke”.

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

158 achieved significantly higher F-measure (paired t-test $P < 0.05$) (Figure 1B), recall, taxon
159 detection rate, taxon accuracy rate (Figure 1C), and lower Bray-Curtis dissimilarity than all
160 other methods (Figure 1D).

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

170 **Cross-validated taxonomy classification**

171 Simulated sequence reads, derived from reference databases, allow us to assess
172 method performance across a greater diversity of sequences than a single mock community
173 generally encompasses. We first evaluated classifier performance using stratified k-fold
174 cross-validation of taxonomy classification to simulated reads. The k-fold cross-validation
175 strategy is modified slightly to account for the hierarchical nature of taxonomic
176 classifications, which all of the classifiers in this study (with the exception of legacy BLAST)
177 handle by assigning the lowest (i.e., most specific) taxonomic level where the classification
178 surpasses some user-defined “confidence” or “consensus” threshold (see materials and

179 methods). The modification is to truncate any expected taxonomy in each test set to the
180 maximum level at which an instance of that taxonomy exists in the training set. Simulated
181 reads were generated from Greengenes 99% OTUs 16S rRNA gene or UNITE 99% OTUs ITS
182 reference sequences with species-level annotations. Greengenes 16S rRNA gene simulated
183 reads were generated from full-length 16S rRNA genes (primers 27F/1492R) and V4
184 (primers 515F/806R) and V1-3 sub-domains (primers 27F/534R). The simulated reads do
185 not incorporate artificial sequencing errors (see materials and methods for more details).
186 In this set of tests and below for novel taxa, the “bespoke” classifier had prior probabilities
187 that were inferred from the training set each time it was trained.

188 Classification of cross-validated reads performed better at coarser levels of
189 classification (Figure 2A), similar to the trend observed in mock community results. For
190 bacterial sequences, average classification accuracy for all methods declined from near-
191 perfect scores at family level (V4 domain median F-measure minimum: BLAST+ F=0.92,
192 maximum: legacy BLAST F=0.99), but still retained accurate scores at species level (median
193 minimum: BLAST+ F=0.76, maximum: SortMeRNA F=0.84), relative to some mock
194 community data sets (Figure 2A). Fungal sequences exhibited similar performance, with
195 the exception that mean BLAST+ and vsearch performance was markedly lower at all
196 taxonomic levels, indicating high sensitivity to parameter configurations, and species-level
197 F-measures were in general much lower (median minimum: BLAST+ F=0.17, maximum:
198 UCLUST F=0.45) than those of bacterial sequence classifications (Figure 2A).

199 Species-level classifications of 16S rRNA gene simulated sequences were best with
200 optimized UCLUST and SortMeRNA configurations for V4 domain, and Naive Bayes and

201 RDP for V1-3 domain and full-length 16S rRNA gene sequences (Figure 2B). UCLUST
202 achieved the highest F-measure for ITS classification ($F = 0.51$). However, all optimized
203 classifiers achieved similar F-measure ranges, with the exception of legacy BLAST for ITS
204 sequences (Figure 2B).

205 Species-level classification performance of 16S rRNA gene simulated reads was
206 significantly correlated between each sub-domain and the full-length gene sequences
207 (Figure 2C). In our tests, full-length sequences exhibited slightly lower accuracy than V1-3
208 and V4 sub-domains. The relative performance of full-length 16S rRNA genes versus
209 hypervariable sub-domain reads is variable in the literature [12, 16–21], and our results
210 add another data point to the ongoing discussion of this topic. Nevertheless, species-level
211 classifications yielded strong correlation between method configurations (Figure 2C) and
212 optimized method performance (Figure 2B), suggesting that primer choice impacts
213 classification accuracy uniformly across all methods. Hence, we focused on V4 sub-domain
214 reads for downstream analyses.

215

216 **Novel taxon classification evaluation**

217 Novel taxon classification offers a unique perspective on classifier behavior,
218 assessing how classifiers perform when challenged with a “novel” clade that is not
219 represented in the reference database [22–25]. An ideal classifier should identify the
220 nearest taxonomic lineage to which this taxon belongs, but no further. In this evaluation, a
221 reference database is subsampled k times to generate query and reference sequence sets,

10

222 as for cross-validated classification, but two important distinctions exist: 1) the reference
223 database used for classification excludes any sequence that matches the taxonomic
224 affiliation of the query sequences at taxonomic level L , the taxonomic rank at which
225 classification is being attempted; and 2) this is performed at each taxonomic level, in order
226 to assess classification performance when each method encounters a “novel” species,
227 genus, family, etc.

228 Due to these differences, interpretation of novel taxon classification results is
229 different from that of mock community and cross-validated classifications. For the latter,
230 classification accuracy may be assessed at each taxonomic level for each classification
231 result: mean classification accuracy at family level and species level evaluate the same
232 results but focus on different taxonomic levels of classification. For novel taxa, however,
233 different query and reference sequences are compiled for classification at each taxonomic
234 level and separate classifications are performed for each. Hence, classifications at family
235 and species level are independent events — one assesses how accurately each method
236 performs when it encounters a “novel” family that is not represented in the reference
237 database, the other when a “novel” species is encountered.

238 Novel taxon evaluations employ a suite of modified metrics, to provide more
239 information on what types of classification errors occur. Precision, recall, and F-measure
240 calculations at each taxonomic level L assess whether an accurate taxonomy classification
241 was made at level $L-1$: for example, a “novel” species should be assigned a genus, because
242 the correct species class is not represented within the reference database. Any species-
243 level classification in this scenario is an *overclassification* (affecting both recall and

244 precision) [25]. Overclassification is one of the key metrics for novel taxa evaluation,
245 indicating the degree to which novel sequences will be interpreted as known organisms.
246 This overclassification is often highly undesirable because it leads, for example, to the
247 incorrect classification of unknown but harmless environmental sequences as known
248 pathogens. Novel sequences that are classified within the correct clade, but to a less specific
249 level than L , are *underclassified* (affecting recall but not precision) [25]. Sequences that are
250 classified into a completely different clade are *misclassified* (affecting both recall and
251 precision) [25].

252 Precision, recall, and F-measure all gradually increase from average scores near 0.0
253 at class level, reaching peak scores at genus level for bacteria and species level for fungi
254 (Figure 3A-C). These trends are paired with gradual decreases in underclassification and
255 misclassification rates for all classification methods, indicating that all classifiers perform
256 poorly when they encounter sequences with no known match at the class, order, or family
257 levels (Figure 3D-F). At species level, UCLUST, BLAST+, and vsearch achieved significantly
258 better F-measures than all other methods for 16S rRNA gene classifications ($P < 0.05$)
259 (Figure 3G). UCLUST achieved significantly better F-measures than all other methods for
260 ITS classifications (Figure 3G). Over-, under-, and misclassification scores are less
261 informative for optimizing classifiers for real use cases, as most methods could be
262 optimized to yield near-zero scores for each of these metrics separately, but only through
263 extreme configurations, leading to F-measures that would be unacceptable under any
264 scenario. Note that all comparisons were made between methods optimized to maximize
265 (or minimize) a single metric, and hence the configurations that maximize precision are

266 frequently different from those that maximize recall or other metrics. This trade-off
267 between different metrics is discussed in more detail below.

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

278

279 **Multi-evaluation method optimization**

280 The mock community and cross-validation classification evaluations yielded similar trends
281 in configuration performance, but optimizing parameters choices for the novel taxa
282 generally lead to suboptimal choices for the mock community and cross-validation tests
283 (Figure 4). We sought to determine the relationship between method configuration
284 performance for each evaluation, and use this information to select configurations that
285 perform best across all evaluations. For 16S rRNA gene sequence species-level
286 classification, method configurations that achieve maximum F-measures for mock and

287 cross-validated sequences perform poorly for novel taxon classification (Figure 4B).
288 Optimization is more straightforward for genus-level classification of 16S rRNA gene
289 sequences (Figure 4A) and for fungal sequences (Figure 4C-D), for which configuration
290 performance (measured as mean F-measure) is maximized by similar configurations
291 among all three evaluations.

292 To identify optimal method configurations, we set accuracy score minimum
293 thresholds for each evaluation by identifying natural breaks in the range of quality scores,
294 selecting methods and parameter ranges that met these criteria. Table 2 lists method
295 configurations that maximize species-level classification accuracy scores for mock
296 community, cross-validated, and novel taxon evaluations under several common operating
297 conditions. “Balanced” configurations are recommended for general use, and are methods
298 that maximize F-measure scores. “Precision” and “Recall” configurations maximize
299 precision and recall scores, respectively, for mock, cross-validated, and novel-taxa
300 classifications (Table 2). “Novel” configurations optimize F-measure scores for novel taxon
301 classification, and secondarily for mock and cross-validated performance (Table 2). These
302 configurations are recommended for use with sample types that are expected to contain
303 large proportions of unidentified species, for which overclassification can be excessive.
304 However, these configurations may not perform optimally for classification of known
305 species (i.e., underclassification rates will be higher). For fungi, the same configurations
306 recommended for “Precision” perform well for novel taxon classification (Table 2). For 16S
307 rRNA gene sequences, BLAST+, UCLUST, and vsearch consensus classifiers perform best for
308 novel taxon classification (Table 2).

309

310 Computational runtime

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

322 We assessed computational runtime as a linear function of 1) the number of query
323 sequences and 2) the number of reference sequences. Linear dependence is empirically
324 evident in Figure 5. For both of these metrics, the slope is the most important measure of
325 performance. The intercept may include the amount of time taken to train the classifier,
326 preprocess the reference sequences, load preprocessed data, or other “setup” steps that
327 will diminish in significance as sequence counts grow, and hence is negligible.

328 UCLUST (0.000028 s/sequence), vsearch (0.000072 s/sequence), BLAST+
329 (0.000080 s/sequence), and legacy BLAST (0.000100 s/sequence) all exhibit shallow

330 slopes with increasing numbers of reference sequences. Naive Bayes (0.000483
331 s/sequence) and SortMeRNA (0.000543 s/sequence) yield moderately higher slopes, and
332 RDP (0.001696 s/sequence) demonstrates the steepest slope (Figure 5A). For runtime as a
333 function of query sequence count, UCLUST (0.002248 s/sequence), RDP (0.002920
334 s/sequence), and SortMeRNA (0.003819 s/sequence) have relatively shallow slopes
335 (Figure 5B). Naive Bayes (0.022984 s/sequence), BLAST+ (0.026222s/sequence) , and
336 vsearch (0.030190 s/sequence) exhibit greater slopes. Legacy BLAST (0.133292
337 s/sequence) yielded a slope magnitudes higher than other methods, rendering this method
338 impractical for large data sets.

339

340

341 **Discussion**

342 We have developed and validated several machine-learning and alignment-based
343 classifiers provided in q2-feature-classifier and benchmarked these classifiers, as well as
344 other common classification methods, to evaluate their strengths and weaknesses for
345 marker-gene amplicon sequence classification across a range of parameter settings for
346 each (Table 2).

347 Each classifier required some degree of optimization to define top-performing
348 parameter configurations, with the sole exception of QIIME 1's legacy BLAST wrapper,
349 which was unaffected by its only user-defined parameter, e-value, over a range of 10^{-10} to

350 1000. For all other methods, performance varied widely depending on parameter settings,
351 and a single method could achieve among the worst performance with one configuration
352 but among the best performance with another. Configurations greatly affected accuracy
353 with mock community, cross-validated, and novel taxon evaluations, indicating that
354 optimization is necessary under a variety of performance conditions, and optimization for
355 one condition may not necessarily translate to another. Mock community and cross-
356 validated evaluations exhibited similar results, but novel taxon evaluations selected
357 different optimal configurations for most methods (Figure 4), indicating that configurations
358 optimized to one condition, e.g., high-recall classification of known sequences, may be less
359 suited for other conditions, e.g., classification of novel sequences. Table 2 lists the top-
360 performing configuration for each method for several standard performance conditions.

361 Optimal configurations also varied among different evaluation metrics. Precision
362 and recall, in particular, exhibited some mutual opposition, such that methods increasing
363 precision reduced recall. For this reason, F-measure, the harmonic mean of precision and
364 recall, is a useful metric for choosing configurations that are well balanced for average
365 performance. “Balanced” method configurations — which maximize F-measure scores for
366 mock, cross-validated, and novel taxon evaluations (Table 2) — are best suited for a wide
367 range of user conditions. The naive Bayes classifier with k-mer lengths of 6 or 7 and
368 confidence = 0.7 (or confidence ≥ 0.9 if using bespoke class weights), RDP with confidence
369 = 0.6-0.7, and UCLUST (minimum consensus = 0.51, minimum similarity = 0.9, max accepts
370 = 3) perform best under these conditions (Table 2). Performance is dramatically improved
371 using bespoke class weights for 16S rRNA sequences (Figure 4A-B), though this approach is

372 developmental and only applicable when the expected composition of samples is known in
373 advance (a scenario that is becoming increasingly common with the increasing quantity of
374 public microbiome data, and which could be aided by microbiome data sharing resources
375 such as Qiita (<http://qiita.microbio.me>)). For ITS sequences, the naive Bayes classifier with
376 k-mer lengths of 6 or 7 and confidence ≥ 0.9 , or RDP with confidence = 0.7-0.9, perform
377 best, and the effects of bespoke class weights are less pronounced (Figure 4C-D).

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

383 Other users may require high-recall classifiers when false-negatives and
384 underclassification hinder interpretation, but false positives (mostly overclassification to a
385 closely related species) are less damaging. For example, in environments with high
386 numbers of unidentified species, a high-precision classifier may yield large numbers of
387 unclassified sequences; in such cases, a second pass with a high-recall configuration (Table
388 2) may provide useful inference of what taxa are most similar to these unclassified
389 sequences. When recall is optimized, precision tends to suffer slightly (leading to similar F-
390 measure scores to “balanced” configurations) but novel taxon classification accuracy is
391 minimized, as these configurations tend to overclassify (Table 2). Any user prioritizing
392 recall ought to be aware of and acknowledge these risks, e.g., when sharing or publishing
393 their results, and understand that many of the species-level classifications may be wrong,

394 particularly if the samples are expected to contain many uncharacterized species. For 16S
395 rRNA gene sequences, naive Bayes bespoke classifiers with k-mer lengths between 12-32
396 and confidence = 0.5 yield maximal recall scores, but RDP (confidence = 0.5) and naive
397 Bayes (uniform class weights, confidence = 0.5, k-mer length = 11, 12, or 18) also perform
398 well (Table 2). Fungal recall scores are maximized by the same configurations
399 recommended for “Balanced” classification, i.e., naive Bayes classifiers with k-mer lengths
400 between 6-7 and confidence between 0.92-0.98, or RDP with confidence between 0.7-0.9
401 (Table 2).

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

407 This study did not compare methods for classification of shotgun metagenome
408 sequencing data sets, which present a series of unique challenges that do not exist for
409 marker-gene amplicon sequence data. These include much higher unique sequence counts
410 (making runtime a greater priority), the use of fully sequenced genomes as reference
411 sequences, and different analysis and quality control protocols. Metagenome sequences
412 also exhibit heterogenous coverage and length, unlike marker-gene amplicon sequences,
413 which typically have uniform start sites and read lengths within a single sequencing run. A
414 recent benchmark of metagenome taxonomic profiling methods describes similar results to
415 our benchmark of marker-gene sequence classifiers: most profilers perform well from

416 phylum to family level but performance degrades at genus and species levels; different
417 methods display superior performance according to different performance metrics; and
418 parameter configuration dramatically impacts performance [26]. In the current study we
419 focused on benchmarking and optimizing classifiers for marker-gene amplicon sequence
420 data, in light of the distinct needs of metagenome and marker-gene sequence datasets.

421 **Conclusions**

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

429 We also present the results of a benchmark of several widely used taxonomy
430 classifiers for marker-gene amplicon sequences, and recommend the top-performing
431 methods and configurations for the most common user scenarios. Our recommendations
432 for “balanced” methods (Table 2) will be appropriate for most users who are classifying
433 16S rRNA gene or fungal ITS sequences, but other users may prioritize high-precision (low
434 false-positive) or high-recall (low false-negative) methods.

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

441

442 **Methods**

443 **Mock communities**

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

450 The inclusion of multiple mock community samples is important to avoid overfitting;
451 optimizing method performance to a small set of data could result in overfitting to the

452 specific community compositions or conditions under which those data were generated,
453 which reduces the robustness of the classifier.

454 **Cross-validated simulated reads**

455 The simulated reads used here were derived from the reference databases using the
456 “Cross-validated classification performance” notebooks in our project repository. The
457 reference databases were either Greengenes or UNITE (99% OTUs) that were cleaned
458 according to taxonomic label to remove sequences with ambiguous or null labels.
459 Reference sequences were trimmed to simulate amplification using standard PCR primers
460 and slice out the first 250 bases downstream (3’) of the forward primer. The bacterial
461 primers used were 27F/1492R [27] to simulate full-length 16S rRNA gene sequences,
462 515F/806R [28] to simulate 16S rRNA gene V4 domain sequences, and 27F/534R [29] to
463 simulate 16S rRNA gene V1-3 domain sequences; the fungal primers used were
464 BITSf/B58S3r [30] to simulate ITS1 internal transcribed spacer DNA sequences. The exact
465 sequences were used for cross validation, and were not altered to simulate any sequencing
466 error; thus, our benchmarks simulate denoised sequence data [4] and isolate classifier
467 performance from impacts from sequencing errors. Each database was stratified by
468 taxonomy and 10-fold randomised cross-validation data sets were generated using scikit-
469 learn’s library functions. Where a taxonomic label had less than 10 instances, taxonomies
470 were amalgamated to make sufficiently large strata. If, as a result, a taxonomy in any test
471 set was not present in the corresponding training set, the expected taxonomy label was
472 truncated to the nearest common taxonomic rank observed in the training set (e.g.,

473 *Lactobacillus casei* would become *Lactobacillus*). The notebook detailing simulated read
474 generation (for both cross-validated and novel taxon reads) prior to taxonomy
475 classification is available at [https://github.com/caporaso-lab/tax-credit-](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/novel-taxa/dataset-generation.ipynb)
476 [data/blob/0.1.0/ipynb/novel-taxa/dataset-generation.ipynb](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/novel-taxa/dataset-generation.ipynb).

477 Classification performance was also slightly modified from a standard machine-learning
478 scenario as the classifiers in this study are able to refuse classification if they are not
479 confident above a taxonomic level for a given sample. This also accommodates the
480 taxonomy truncation that we performed for this test. The methodology was consistent with
481 that used below for novel taxon evaluations, but we defer this description to the next
482 section.

483 **“Novel taxon” simulation analysis**

484 “Novel taxon” classification analysis was performed to test the performance of classifiers
485 when assigning taxonomy to sequences that are not represented in a reference database,
486 e.g., as a simulation of what occurs when a method encounters an undocumented species
487 [22–25]. In this analysis, simulated amplicons were filtered from those used for the cross-
488 validation analysis. For all sequences present in each test set, sequences sharing taxonomic
489 affiliation at a given taxonomic level L (e.g., to species level) in the corresponding training
490 set were removed. Taxa are stratified among query and test sets such that for each query
491 taxonomy at level L , no reference sequences match that taxonomy, but at least one
492 reference sequence will match the taxonomic lineage at level $L-1$ (e.g., same genus but
493 different species). An ideal classifier would assign taxonomy to the nearest common

494 taxonomic lineage (e.g., genus), but would not “overclassify” [25] to near neighbors (e.g.,
495 assign species-level taxonomy when species X is removed from the reference database).
496 For example, a “novel” sequence representing the species *Lactobacillus brevis* should be
497 classified as “*Lactobacillus*”, without species-level annotation, in order to be considered a
498 true positive in this analysis. As described above for cross-validated reads, these novel taxa
499 simulated communities were also tested in both bacterial (B) and fungal (F) databases on
500 simulated amplicons trimmed to simulate 250-nt sequencing reads.

501 Novel taxon classification performance is evaluated using precision, recall, F-
502 measure, overclassification rates, underclassification rates, and misclassification rates [25]
503 for each taxonomic level (phylum to species), computed with the following definitions (see
504 below, *Performance analyses using simulated reads*, for full description of precision, recall,
505 and F-measure calculations):

- 506 1) A true positive is considered the nearest correct lineage contained in the reference
507 database. For example, if *Lactobacillus brevis* is removed from the reference
508 database and used as a query sequence, the only correct taxonomy classification
509 would be “*Lactobacillus*”, without species-level classification.
- 510 2) A false positive would be either a classification to a different *Lactobacillus* species
511 (*Overclassification*), or any genus other than *Lactobacillus* (*Misclassification*).
- 512 3) A false negative occurs if an expected taxonomy classification (e.g., “*Lactobacillus*”)
513 is not observed in the results. Note that this will be the modified taxonomy expected
514 when using a naive reference database, and is not the same as the true taxonomic
515 affiliation of a query sequence in the novel taxa analysis. A false negative results

516 from misclassification, overclassification, or when the classification contains the
517 correct basal lineage, but does not assign a taxonomy label at level L
518 (*Underclassification*). E.g., classification as "*Lactobacillaceae*", but no genus-level
519 classification.

520 **Taxonomy classification**

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

524 1. q2-feature-classifier multinomial naive Bayes classifier. Varied k-mer length
525 in {4, 6, 7, 8, 9, 10, 11, 12, 14, 16, 18, 32} and confidence threshold in {0, 0.5, 0.7, 0.9,
526 0.92, 0.94, 0.96, 0.98, 1}.

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

531 3. vsearch [10] global sequence alignment, followed by consensus taxonomy
532 classification implemented in q2-feature-classifier. Varied max accepts from 1 to
533 100; percent identity from 0.80 to 0.99; and minimum consensus from 0.51 to 0.99.
534 See description below.

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

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

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

542 7. UCLUST [14] (QIIME1 wrapper) varying minimum consensus fraction from
543 0.51 to 0.99; similarity from 0.8 to 0.9; and max accepts from 1 to 10.

544

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

552 Classification of bacterial/archaeal 16S rRNA gene sequences was made using the
553 Greengenes (13_8 release) [5] reference sequence database preclustered at 99% ID, with
554 amplicons for the domain of interest extracted using primers 27F/1492R [27], 515F/806R
555 [28], or 27F/534R [29] with q2-feature-classifier's `extract_reads` method. Classification of
556 fungal ITS sequences was made using the UNITE database (version 7.1 QIIME developer

26

557 release) [31] preclustered at 99% ID. For the cross validation and novel taxon
558 classification tests we prefiltered to remove sequences with incomplete or ambiguous
559 taxonomies (containing the substrings 'unknown', 'unidentified', or '_sp' or terminating at
560 any level with '_').

561
562 The notebooks detailing taxonomy classification sweeps of mock communities are available
563 at <https://github.com/caporaso-lab/tax-credit-data/tree/0.1.0/ipynb/mock-community>.
564 Cross-validated read classification sweeps are available at <https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/cross-validated/taxonomy-assignment.ipynb>. Novel
565 taxon classification sweeps are available at <https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/novel-taxa/taxonomy-assignment.ipynb>.

568

569 **Runtime analyses**

570 The tax-credit framework employs two different runtime metrics: as a function of 1) the
571 number of query sequences or 2) the number of reference sequences. Taxonomy classifier
572 runtimes were logged while performing classifications of pseudorandom subsets of 1,
573 2,000, 4,000, 6,000, 8,000, and 10,000 sequences from the Greengenes 99% OTU database.
574 Each subset was drawn once then used for all of the tests as appropriate. All runtimes were
575 computed on the same Linux workstation (Ubuntu 16.04.2 LTS, Intel Xeon CPU E7-4850 v3
576 @ 2.20GHz, 1TB memory). The exact commands used for runtime analysis are presented in

577 the “Runtime analyses” notebook in the project repository ([https://github.com/caporaso-](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/runtime/analysis.ipynb)
578 [lab/tax-credit-data/blob/0.1.0/ipynb/runtime/analysis.ipynb](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/runtime/analysis.ipynb)).

579 **Performance analyses using simulated reads**

580 Cross-validated and novel taxa reads are evaluated using the classic precision, recall, and F-
581 measure metrics [5] (novel taxa use the standard calculations as described below, but
582 modified definitions for true positive (TP), false positive (FP), and false negative (FN), as
583 described above for novel taxon classification analysis).

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

- 585 ○ $Precision = TP/(TP+FP)$ or the fraction of sequences that were classified correctly at
586 level L.
- 587 ○ $Recall = TP/(TP+FN)$ or the fraction of expected taxonomic labels that were
588 predicted at level L.
- 589 ○ $F\text{-measure} = 2 \times Precision \times Recall / (Precision + Recall)$, or the harmonic mean of
590 precision and recall.

591 The Jupyter notebook detailing commands used for evaluation of cross-validated read
592 classifications is available at [https://github.com/caporaso-lab/tax-credit-](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/cross-validated/evaluate-classification.ipynb)
593 [data/blob/0.1.0/ipynb/cross-validated/evaluate-classification.ipynb](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/cross-validated/evaluate-classification.ipynb). The notebook for
594 evaluation of novel taxon classifications is available at [https://github.com/caporaso-](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/novel-taxa/evaluate-classification.ipynb)
595 [lab/tax-credit-data/blob/0.1.0/ipynb/novel-taxa/evaluate-classification.ipynb](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/novel-taxa/evaluate-classification.ipynb).

596 **Performance analyses using mock communities**

597 The Jupyter notebook detailing commands used for evaluation of mock communities,
598 including the three evaluation types described below, is available at
599 [https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/mock-](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/mock-community/evaluate-classification-accuracy.ipynb)
600 [community/evaluate-classification-accuracy.ipynb](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/mock-community/evaluate-classification-accuracy.ipynb).

601 **Precision and Recall**

602 Classic precision, recall, and F-measure are used to calculate mock community
603 classification accuracy, using the definitions given above for simulated reads. These metrics
604 require knowing the expected classification of each sequence, which we determine by
605 performing a gapless alignment between each representative sequence in the mock
606 community and the marker-gene sequences of each microbial strain added to the mock
607 community. These “expected sequences” are provided for the mock communities in
608 mockrobiota [11]. Representative sequences are assigned the taxonomy of the best
609 alignment, and any representative sequence with more than 3 mismatches to the expected
610 sequences are excluded from precision/recall calculations. If a representative sequence
611 aligns to more than one expected sequence equally well, all top hits are accepted as the
612 “correct” classification. This scenario is rare and typically only occurred when different
613 strains of the same species were added to the same mock community to intentionally
614 produce this challenge (e.g., for mock-12 as described by [4]). Precision, recall, and F-
615 measure are then calculated by comparing the “expected” classification for each mock

616 community sequence to the classifications predicted by each taxonomy classifier using the
617 full reference databases, as described above.

618 **Taxon accuracy rate and taxon detection rate**

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

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

- 628 ○ true positive (*TP*), if that taxon is both observed and expected.
- 629 ○ false positive (*FP*), if that taxon is observed but not expected.
- 630 ○ false negative (*FN*), if a taxon is expected but not observed.

631 These are used to calculate TAR and TDR as:

- 632 ○ $TAR = TP/(TP+FP)$ or the fraction of observed taxa that were expected at level L.
- 633 ○ $TDR = TP/(TP+FN)$ or the fraction of expected taxa that are observed at level L.

634

635 **Bray-Curtis Dissimilarity**

636 Bray-Curtis dissimilarity [32] is used to measure the degree of dissimilarity between two
637 samples as a function of the abundance of each species label present in each sample,
638 treating each species as equally related. This is a useful metric for evaluating classifier
639 performance by assessing the relative distance between each predicted mock community
640 composition (abundance of taxa in a sample based on results of a single classifier) and the
641 expected composition of that sample. For each classifier, Bray-Curtis distances between the
642 expected and observed taxonomic compositions are calculated for each sample in each
643 mock community dataset; this yields a single expected-observed distance for each
644 individual observation. The distance distributions for each method are then compared
645 statistically using paired or unpaired t-tests to assess whether one method (or
646 configuration) performs consistently better than another.

647 **New taxonomy classifiers**

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

653 **Machine learning taxonomy classifiers**

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

661

662 The plugin provides a default method which is to extract k-mer counts from reference
663 sequences and train the scikit-learn multinomial naive Bayes classifier, and it is this
664 method that we test extensively here. Specifically, the pipeline consists of a
665 `sklearn.feature_extraction.text.HashingVectorizer` feature extraction step followed by a
666 `sklearn.naive_bayes.MultinomialNB` classification step. The use of a hashing feature
667 extractor allows the use of significantly longer k-mers than the 8-mers that are used by
668 RDP Classifier, and we tested up to 32-mers. Like most scikit-learn classifiers, we are able
669 to set class weights when training the multinomial naive Bayes classifiers. In the naive
670 Bayes setting, setting class weights means that class priors are not derived from the
671 training data or set to be uniform, as they are for the RDP Classifier. For more detail on how
672 class weights enter the calculations please refer to the scikit-learn User Guide
673 (<http://scikit-learn.org>).

674

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

683

684 For this study, we performed two parameter sweeps on the mock communities: an initial
685 broad sweep to optimize feature extraction parameters and then a more focused sweep to
686 optimise k-mer length and confidence parameter settings. These sweeps included varying
687 the assumptions regarding class weights. The focussed sweeps were also performed for the
688 cross-validated and novel taxa evaluations, but only for the assumption of uniform class
689 priors. The results for the focussed sweeps across all data sets are those which are
690 compared against the other classifiers in this work.

691

692 The broad sweeps used a modified scikit-learn pipeline which consisted of the
693 `sklearn.feature_extraction.text.HashingVectorizer`, followed by the
694 `sklearn.feature_extraction.text.TfidfTransformer`, then the
695 `sklearn.naive_bayes.MultinomialNB`. We performed a full grid search over the parameters

33

696 shown in Table 3. The conclusion from the initial sweep was that the TfidfTransformer step
697 did not significantly improve classification, that `n_features` should be set to 8192, feature
698 vectors should be normalised using L2 normalisation and that the alpha parameter for the
699 naive Bayes classifier should be set to 0.001. Please see [https://github.com/caporaso-](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/mock-community/evaluate-classification-accuracy-nb-extra.ipynb)
700 [lab/tax-credit-data/blob/0.1.0/ipynb/mock-community/evaluate-classification-accuracy-](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/mock-community/evaluate-classification-accuracy-nb-extra.ipynb)
701 [nb-extra.ipynb](https://github.com/caporaso-lab/tax-credit-data/blob/0.1.0/ipynb/mock-community/evaluate-classification-accuracy-nb-extra.ipynb) for details.

702 **Consensus taxonomy alignment-based classifiers**

703
704 Two new classifiers implemented in `q2-feature-classifier` perform consensus taxonomy
705 classification based on alignment of a query sequence to a reference sequence. The
706 methods `classify_consensus_vsearch` and `classify_consensus_blast` use the global aligner
707 `vsearch` [10] or the local aligner `BLAST+` [9], respectively, to return up to `maxaccepts`
708 reference sequences that align to the query with at least `perc_identity` similarity. A
709 consensus taxonomy is then assigned to the query sequence by determining the taxonomic
710 lineage on which at least `min_consensus` of the aligned sequences agree. This consensus
711 taxonomy is truncated at the taxonomic level at which less than `min_consensus` of
712 taxonomies agree. For example, if a query sequence is classified with `maxaccepts=3`,
713 `min_consensus=0.51`, and the following top hits:

714
715 `k__Bacteria; p__Firmicutes; c__Bacilli; o__Lactobacillales; f__Lactobacillaceae;`
716 `g__Lactobacillus; s__brevis`

717 k__Bacteria; p__Firmicutes; c__Bacilli; o__Lactobacillales; f__Lactobacillaceae;

718 g__Lactobacillus; s__brevis

719 k__Bacteria; p__Firmicutes; c__Bacilli; o__Lactobacillales; f__Lactobacillaceae;

720 g__Lactobacillus; s__delbrueckii

721

722 The taxonomy label assigned will be k__Bacteria; p__Firmicutes; c__Bacilli;

723 o__Lactobacillales; f__Lactobacillaceae; g__Lactobacillus; s__brevis. However, if

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

725 c__Bacilli; o__Lactobacillales; f__Lactobacillaceae; g__Lactobacillus.

726

727

728 **Declarations**

729 **Ethics approval and consent to participate**

730 Not applicable

731 **Consent for publication**

732 Not applicable

733 Availability of data and materials

734 Mock community sequence data used in this study are publicly available in mockrobiota
735 [11] under the study identities listed in Table 1. All other data generated in this study, and
736 all new software, is available in our GitHub repositories under the BSD license. The tax-
737 credit repository can be found at: <https://github.com/caporaso-lab/tax-credit>, and static
738 versions of all analysis notebooks, which contain all code and analysis results, can be
739 viewed there. The q2-feature-classifier repository can be accessed at
740 <https://github.com/qiime2/q2-feature-classifier>; as a QIIME2 core plugin, it is
741 automatically installed any time QIIME2 (<https://qiime2.org/>) is installed.

742
743 **Project name:** q2-feature-classifier

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

745 **Operating system(s):** macOS, Linux

746 **Programming language:** Python

747 **Other requirements:** QIIME2

748 **License:** BSD-3-Clause

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

750

751 **Project name:** tax-credit

752 **Project home page:** <https://github.com/caporaso-lab/tax-credit>

753 **Operating system(s):** macOS, Linux

754 **Programming language:** Python

755 **Other requirements:** None (QIIME2 required for some optional functions)

756 **License:** BSD-3-Clause

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

758

759

760 Funding

761 This work was funded in part by National Science Foundation award 1565100 to JGC and

762 RK, awards from the Alfred P. Sloan Foundation to JGC and RK, awards from the

763 Partnership for Native American Cancer Prevention (NIH/NCI U54CA143924 and

764 U54CA143925) to JGC, and National Health and Medical Research Council of Australia

765 award APP1085372 to GAH, JGC and RK. These funding bodies had no role in the design of
766 the study, the collection, analysis, or interpretation of data, or in writing the manuscript.

767 **Acknowledgments**

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

770 **Authors' Contributions**

771 NAB, RK, and JGC conceived and designed tax-credit. NAB, BDK, JGC, and JRR contributed
772 to tax-credit. BDK, MD, JGC, and NAB contributed to q2-feature-classifier. BDK, JGC, MD,
773 JRR, and EB provided QIIME 2 integration with q2-feature-classifier. JGC and GAH provided
774 materials and support. NAB, BDK, JGC, and GAH wrote the manuscript with input from all
775 co-authors.

776 **Competing Interests**

777 The authors declare that they have no competing interests.

778

779 **Tables and Figures**

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

Study ID*	Target gene**	Platform	Species	Strains	Citation
mock-1	16S	HiSeq	46	48	[33]
mock-2	16S	MiSeq	46	48	[33]
mock-3	16S	MiSeq	21	21	[33]
mock-4	16S	MiSeq	21	21	[33]
mock-5	16S	MiSeq	21	21	[33]
mock-7	16S	HiSeq	67	67	[34]
mock-8	16S	HiSeq	67	67	[11]
mock-9	ITS	HiSeq	13	16	[11]
mock-10	ITS	HiSeq	13	16	[11]
mock-12	16S	MiSeq	26	27	[4]
mock-16	16S	MiSeq	56	59	[35]
mock-18	16S	MiSeq	15	15	[36]
mock-19	16S	MiSeq	15	27	[36]
mock-20	16S	MiSeq	20	20	[37]
mock-21	16S	MiSeq	20	20	[37]
mock-22	16S	MiSeq	20	20	[37]
mock-23	16S	MiSeq	20	20	[37]
mock-24	ITS	MiSeq	8	8	[38]
mock-26	ITS	FLX Titanium	11	11	[39]

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

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

783 **Abbreviations: 16S = 16S rRNA gene; HiSeq = Illumina HiSeq; MiSeq = Illumina MiSeq.

784

785 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 gene	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	

38

			[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	

786

787 ^aF = F-measure, P = precision, R = recall788 ^bNaive Bayes parameters: k-mer range, confidence789 ^cRDP parameters: confidence790 ^dBLAST+/vsearch parameters: max accepts, minimum consensus, minimum percent
791 identity792 ^eUCLUST parameters: minimum consensus, similarity, max accepts793 ^fThreshold describes the score cutoffs used to define optimal method ranges, in the format:
794 [metric = (mock score, cross-validated score, novel-taxa score)]. If two cutoffs are given,
795 the second indicates a higher cutoff used to select parameters for the developmental NB-
796 bespoke method, and the configurations listed are the union of the two cutoffs: the second
797 cutoff for selecting NB-bespoke, the first for selecting all other methods.
798

799

800 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]

40

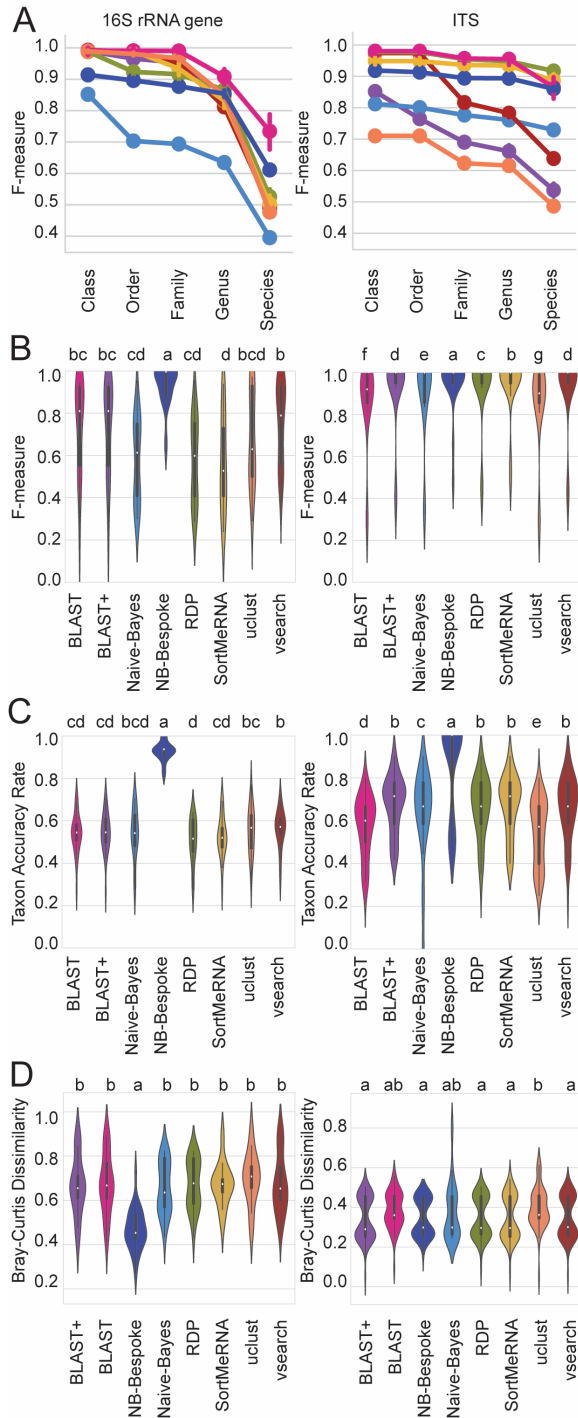
sklearn.feature_extraction.text.TfidfTransformer	norm	l1, l2, None
	use_idf	True, False
sklearn.naive_bayes.MultinomialNB	alpha	0.001, 0.01, 0.1
	class_prior	None, array of class weights
post processing	confidence	0, 0.2, 0.4, 0.6, 0.8

801

802

803

804



805

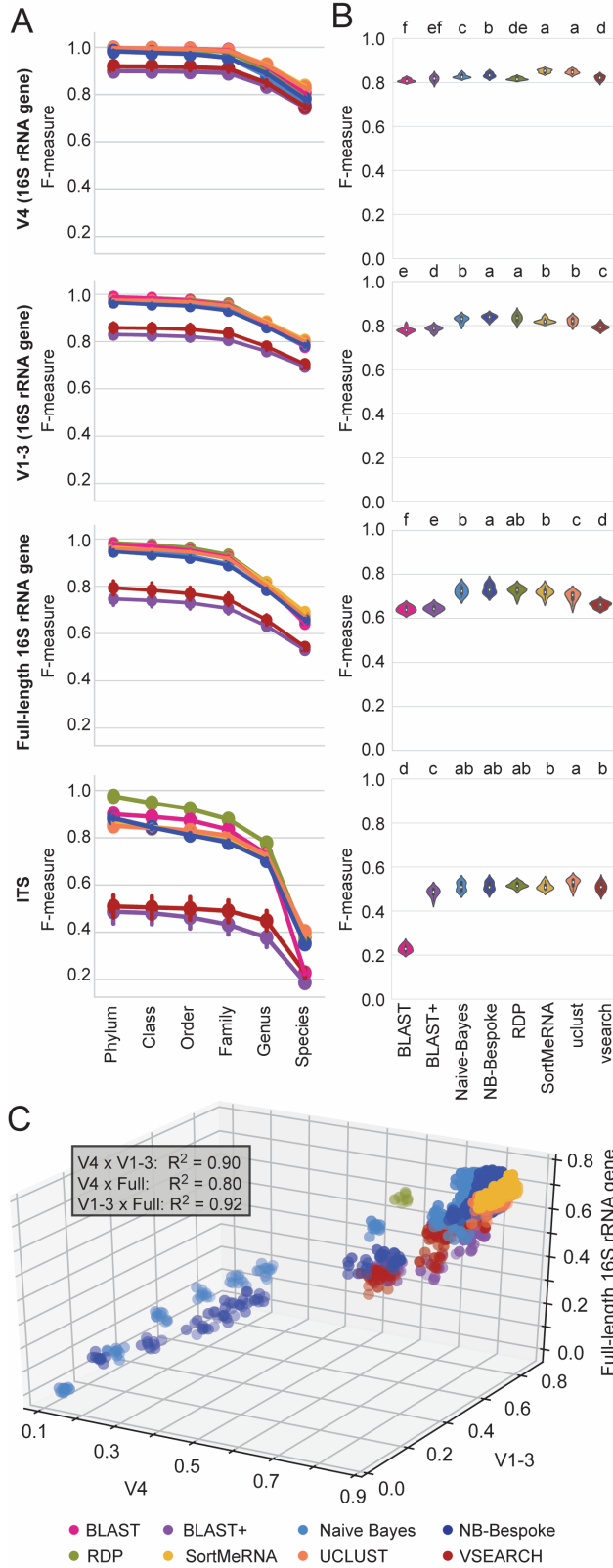
806 Figure 1. Classifier performance on mock community datasets for 16S rRNA gene

807 sequences (left column) and fungal ITS sequences (right column). A, Average F-measure for

42

808 each taxonomy classification method (averaged across all configurations and all mock
809 community datasets) from class to species level. Error bars = 95% confidence intervals. B,
810 Average F-measure for each optimized classifier (averaged across all mock communities) at
811 species level. C, Average taxon accuracy rate for each optimized classifier (averaged across
812 all mock communities) at species level. D, Average Bray-Curtis distance between the
813 expected mock community composition and its composition as predicted by each
814 optimized classifier (averaged across all mock communities) at species level. Violin plots
815 show median (white point), quartiles (black bars), and kernel density estimation (violin)
816 for each score distribution. Violins with different lower-case letters have significantly
817 different means (paired t-test false detection rate-corrected $P < 0.05$).

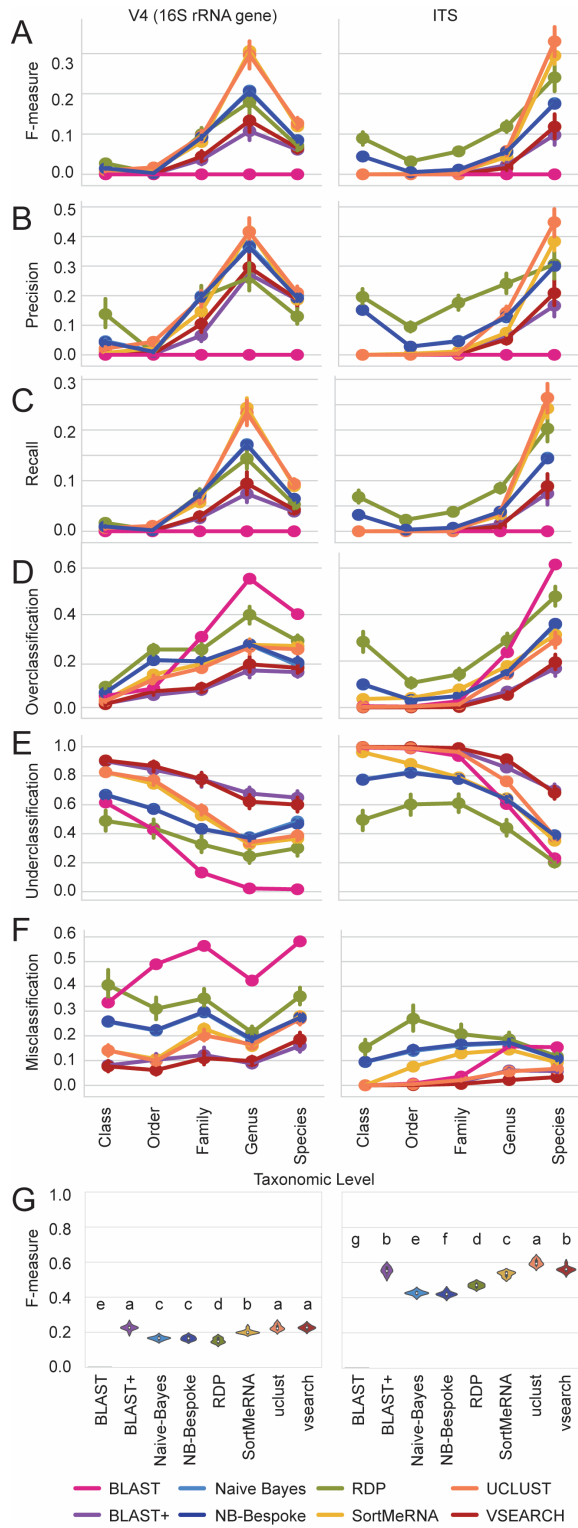
818



819

44

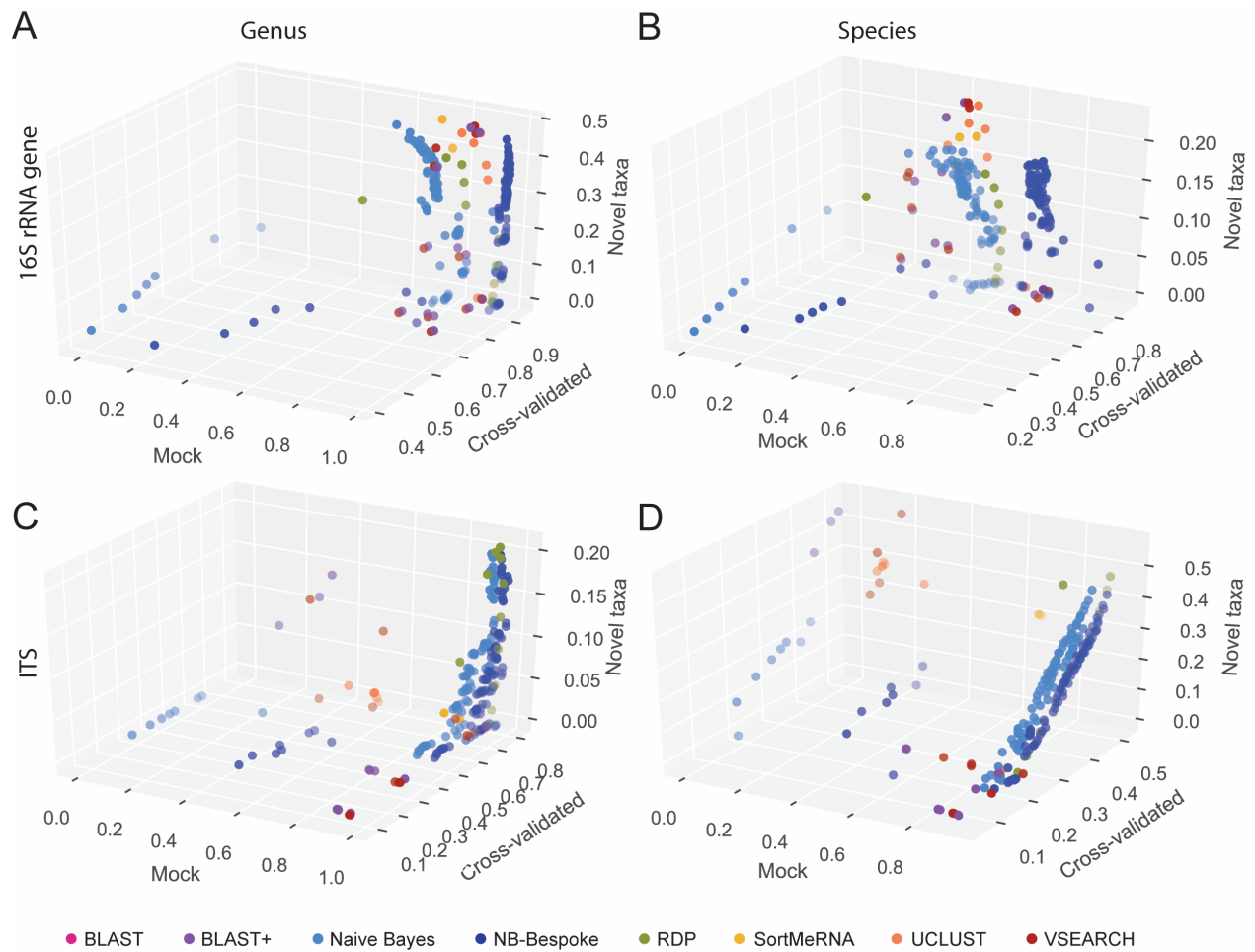
820 Figure 2. Classifier performance on cross-validated sequence datasets. Classification
821 accuracy of 16S rRNA gene V4 sub-domain (first row), V1-3 sub-domain (second row), full-
822 length 16S rRNA gene (third row), and fungal ITS sequences (fourth row). A, Average F-
823 measure for each taxonomy classification method (averaged across all configurations and
824 all cross-validated sequence datasets) from class to species level. Error bars = 95%
825 confidence intervals. B, Average F-measure for each optimized classifier (averaged across
826 all cross-validated sequence datasets) at species level. Violins with different lower-case
827 letters have significantly different means (paired t-test false detection rate-corrected $P <$
828 0.05). C, correlation between F-measure performance for each method/configuration
829 classification of V4 sub-domain (x-axis), V1-3 sub-domain (y-axis), and full-length 16S
830 rRNA gene sequences (z-axis). Inset lists the pearson R^2 value for each pairwise
831 correlation; each correlation is significant ($P < 0.001$).
832



833

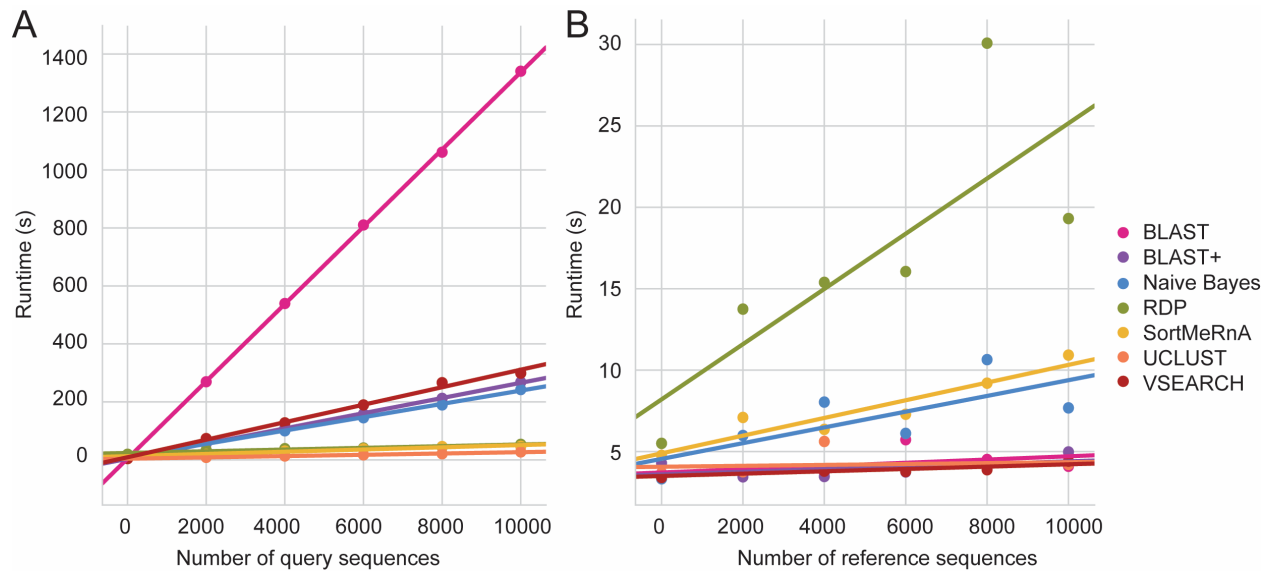
834 Figure 3. Classifier performance on novel-taxa simulated sequence datasets for 16S rRNA
835 gene sequences (left column) and fungal ITS sequences (right column). A-F, Average F-
836 measure (A), precision (B), recall (C), overclassification (D), underclassification (E), and
837 misclassification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

843



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

850



851

852 Figure 5. Runtime performance comparison of taxonomy classifiers. Runtime (s) for each
 853 taxonomy classifier either varying the number of query sequences and keeping a constant
 854 10000 reference sequences (A) or varying the number of reference sequences and keeping
 855 a constant 1 query sequence (B).

856

857

858 References

859 1. Human Microbiome Project Consortium. A framework for human microbiome research. *Nature*.
 860 2012;486:215–21.

861 2. Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, et al. A communal catalogue
 862 reveals Earth's multiscale microbial diversity. *Nature*. 2017;551:457–63.

863 3. Wang Q, Quensen JF 3rd, Fish JA, Lee TK, Sun Y, Tiedje JM, et al. Ecological patterns of nifH genes
 864 in four terrestrial climatic zones explored with targeted metagenomics using FrameBot, a new

49

- 865 informatics tool. *MBio*. 2013;4:e00592-13.
- 866 4. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-resolution
867 sample inference from Illumina amplicon data. *Nat Methods*. 2016;13:581-3.
- 868 5. McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, et al. An improved
869 Greengenes taxonomy with explicit ranks for ecological and evolutionary analyses of bacteria and
870 archaea. *ISME J*. 2012;6:610-8.
- 871 6. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al. QIIME allows
872 analysis of high-throughput community sequencing data. *Nat Methods*. 2010;7:335-6.
- 873 7. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B, Grisel, O., Blondel, M.,
874 Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M.,
875 Perrot, M., Duchesnay, E. Scikit-learn: Machine Learning in Python. *J Mach Learn Res*.
876 2011;12:2825-30.
- 877 8. Buitinck, L., Louppe, G., Blondel, M., Pedregosa, F., Mueller, A., Grisel, O., Niculae, V., Prettenhofer,
878 P., Gramfort, A., Grobler, J., Layton, R., VanderPlas, J., Joly, A., Holt, B., Varoquaux G. API design for
879 machine learning software: experiences from the scikit-learn project. In: *ECML PKDD Workshop:
880 Languages for Data Mining and Machine Learning*. 2013. p. 108-22.
- 881 9. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST : architecture
882 and applications. *BMC Bioinformatics*. 2009;10:421.
- 883 10. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool for
884 metagenomics. *PeerJ*. 2016;4:e2584.
- 885 11. Bokulich NA, Rideout JR, Mercurio WG, Shiffer A, Wolfe B, Maurice CF, et al. mockrobiota: a
886 Public Resource for Microbiome Bioinformatics Benchmarking. *mSystems*. 2016;1.
887 doi:10.1128/mSystems.00062-16.
- 888 12. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian classifier for rapid assignment of rRNA
889 sequences into the new bacterial taxonomy. *Appl Environ Microbiol*. 2007;73:5261-7.
- 890 13. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol*.
891 1990;215:403-10.
- 892 14. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*.
893 2010;26:2460-1.
- 894 15. Kopylova E, Noé L, Touzet H. SortMeRNA: fast and accurate filtering of ribosomal RNAs in
895 metatranscriptomic data. *Bioinformatics*. 2012;28:3211-7.
- 896 16. Soergel DAW, Dey N, Knight R, Brenner SE. Selection of primers for optimal taxonomic
897 classification of environmental 16S rRNA gene sequences. *ISME J*. 2012;6:1440-4.
- 898 17. Liu K-L, Wong T-T. Naïve Bayesian Classifiers with Multinomial Models for rRNA Taxonomic
899 Assignment. *IEEE/ACM Trans Comput Biol Bioinform*. 2013;10:1-1.

- 900 18. Chaudhary N, Sharma AK, Agarwal P, Gupta A, Sharma VK. 16S classifier: a tool for fast and
901 accurate taxonomic classification of 16S rRNA hypervariable regions in metagenomic datasets.
902 PLoS One. 2015;10:e0116106.
- 903 19. Claesson MJ, Wang Q, O'Sullivan O, Greene-Diniz R, Cole JR, Ross RP, et al. Comparison of two
904 next-generation sequencing technologies for resolving highly complex microbiota composition
905 using tandem variable 16S rRNA gene regions. Nucleic Acids Res. 2010;38:e200.
- 906 20. Liu Z, DeSantis TZ, Andersen GL, Knight R. Accurate taxonomy assignments from 16S rRNA
907 sequences produced by highly parallel pyrosequencers. Nucleic Acids Res. 2008;36:e120.
- 908 21. Liu Z, Lozupone C, Hamady M, Bushman FD, Knight R. Short pyrosequencing reads suffice for
909 accurate microbial community analysis. Nucleic Acids Res. 2007;35:e120.
- 910 22. Lanzén A, Jørgensen SL, Huson DH, Gorfer M, Grindhaug SH, Jonassen I, et al. CREST –
911 Classification Resources for Environmental Sequence Tags. PLoS One. 2012;7:e49334.
- 912 23. Lan Y, Wang Q, Cole JR, Rosen GL. Using the RDP classifier to predict taxonomic novelty and
913 reduce the search space for finding novel organisms. PLoS One. 2012;7:e32491.
- 914 24. Deshpande V, Wang Q, Greenfield P, Charleston M, Porrás-Alfaro A, Kuske CR, et al. Fungal
915 identification using a Bayesian classifier and the Warcup training set of internal transcribed spacer
916 sequences. Mycologia. 2016;108:1–5.
- 917 25. Edgar R. SINTAX: a simple non-Bayesian taxonomy classifier for 16S and ITS sequences. 2016.
918 doi:10.1101/074161.
- 919 26. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, et al. Critical Assessment of
920 Metagenome Interpretation—a benchmark of metagenomics software. Nat Methods. 2017;14:1063–
921 71.
- 922 27. Weisburg WG, Barns SM, Pelletier DA, Lane DJ. 16S ribosomal DNA amplification for
923 phylogenetic study. J Bacteriol. 1991;173:697–703.
- 924 28. Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, et al. Ultra-high-
925 throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. ISME J.
926 2012;6:1621–4.
- 927 29. Muyzer G, de Waal EC, Uitterlinden AG. Profiling of complex microbial populations by
928 denaturing gradient gel electrophoresis analysis of polymerase chain reaction-amplified genes
929 coding for 16S rRNA. Appl Environ Microbiol. 1993;59:695–700.
- 930 30. Bokulich NA, Mills DA. Improved Selection of Internal Transcribed Spacer-Specific Primers
931 Enables Quantitative, Ultra-High-Throughput Profiling of Fungal Communities. Appl Environ
932 Microbiol. 2013;79:2519–26.
- 933 31. Kõljalg U, Nilsson RH, Abarenkov K, Tedersoo L, Taylor AFS, Bahram M, et al. Towards a unified
934 paradigm for sequence-based identification of fungi. Mol Ecol. 2013;22:5271–7.

- 935 32. Bray JR, Curtis JT. An Ordination of the Upland Forest Communities of Southern Wisconsin. *Ecol*
936 *Monogr.* 1957;27:325–49.
- 937 33. Bokulich NA, Subramanian S, Faith JJ, Gevers D, Gordon JI, Knight R, et al. Quality-filtering vastly
938 improves diversity estimates from Illumina amplicon sequencing. *Nat Methods.* 2013;10:57–9.
- 939 34. Maurice CF, Haiser HJ, Turnbaugh PJ. Xenobiotics Shape the Physiology and Gene Expression of
940 the Active Human Gut Microbiome. *Cell.* 2013;152:39–50.
- 941 35. Schirmer M, Ijaz UZ, D'Amore R, Hall N, Sloan WT, Quince C. Insight into biases and sequencing
942 errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic Acids Res.* 2015;43:e37.
- 943 36. Tourlousse DM, Yoshiike S, Ohashi A, Matsukura S, Noda N, Sekiguchi Y. Synthetic spike-in
944 standards for high-throughput 16S rRNA gene amplicon sequencing. *Nucleic Acids Res.*
945 2016;;gkw984.
- 946 37. Gohl DM, Vangay P, Garbe J, MacLean A, Hauge A, Becker A, et al. Systematic improvement of
947 amplicon marker gene methods for increased accuracy in microbiome studies. *Nat Biotechnol.*
948 2016;34:942–9.
- 949 38. Taylor DL, Walters WA, Lennon NJ, Bochicchio J, Krohn A, Caporaso JG, et al. Accurate
950 Estimation of Fungal Diversity and Abundance through Improved Lineage-Specific Primers
951 Optimized for Illumina Amplicon Sequencing. *Appl Environ Microbiol.* 2016;82:7217–26.
- 952 39. Ihrmark K, Bödeker ITM, Cruz-Martinez K, Friberg H, Kubartova A, Schenck J, et al. New primers
953 to amplify the fungal ITS2 region--evaluation by 454-sequencing of artificial and natural
954 communities. *FEMS Microbiol Ecol.* 2012;82:666–77.
- 955