

1 Removal of false positives in metagenomic-based 2 taxonomy profiling via targeting Type IIB restriction sites

3 *Supplementary Information*

4 **Supplemental Methods**

5 **Supplemental Figures**

6 **Figure S1.** Distribution of theoretically existent and unique 2b tags for all the 48,475 species
7 in the combined database of GTDB and Ensemble Fungi.

8 **Figure S2.** Sample-specific unique 2b database can largely expand the number of unique 2b
9 tags compared to pre-constructed unique 2b database.

10 **Figure S3.** Comparing MAP2B with other profilers using simulation data generated by shared
11 genome sources of different profilers.

12 **Figure S4.** Benchmarking MAP2B using simulation data generated by genome sources
13 independent of GTDB or with a high mutation rate.

14 **Figure S5.** Assessing MAP2B's microbial profiles in disease discrimination (validation cohort)
15 and prediction of metabolomic profiles (using mNODE).

16 **Figure S6.** Comparison of the algorithm implemented in MAP2B and conventional
17 metagenomic profiling tools.

18 **Figure S7.** Definition of the terminology used in the species identification procedures of
19 MAP2B.

20 **Figure S8.** Top-10 feature combinations in each machine learning classifier when selecting
21 the false positive recognition model.

22 **Figure S9.** Performance evaluations of the machine learning model on false positive
23 recognition based on 50 times five-fold cross-validation.

24 **Reference**

25

26 Supplemental Methods

27 Feature engineering and selection for the false positive recognition model

28 We sought to resolve the false-positive identification issue by leveraging a classifier that relies
29 on meaningful biological features (**Fig. S7**). Consider the *in silico* digestion of reference
30 genomes and a given WMS dataset. For species-*i* in the integrated database of GTDB¹ and
31 Ensembl Fungi², we denote its total number of 2b tags generated by *in silico* digestion of its
32 genome as H_i , representing the number of its theoretically existent 2b tags. Among the H_i
33 tags, there are E_i unique 2b tags, which are single-copy within species-*i*'s genome and are
34 unique to species-*i* w.r.t all other species in the database.

35

36 Given an input WMS dataset with total R reads, we *in silico* extract 2b tags, map them to the
37 species-specific 2b tags, and denote the number of tags unique to species-*i* as Q_i
38 (sequenced unique 2b tags). Among the Q_i tags unique to species-*i*, there are U_i
39 nonredundant (or distinct) sequenced unique 2b tags.

40

41 The genome coverage of species-*i* can be calculated as $C_i = U_i/E_i$. Usually, we have $C_i < 1$
42 for WMS data of complex microbial communities. But we can infer the actual number of
43 sequenced unique 2b tags of species-*i*, denoted as \tilde{Q}_i , by its genome coverage correction,
44 i.e., $\tilde{Q}_i = Q_i/C_i = (Q_i E_i)/U_i$.

45

46 Now we propose seven features to help us distinguish false positives from true positives.
47 Those seven features can be classified into three distinct categories:

48 (1) Related to the genome coverage:

$$49 \text{ Feature 1} = C_i = \frac{U_i}{E_i}, (\text{genome coverage}) \quad (7)$$

50 (2) Related to the taxonomic and sequence abundance:

$$51 \text{ Feature 2} = N_i = \frac{Q_i}{U_i}, (\text{taxonomic counts}) \quad (8)$$

$$52 \text{ Feature 3} = \frac{Q_i}{H_i}, \quad (9)$$

$$53 \text{ Feature 4} = \frac{Q_i}{R}, \quad (10)$$

$$54 \text{ Feature 5} = R_i = \frac{\tilde{Q}_i}{R}, (\text{sequence counts}) \quad (11)$$

55 (3) Related to both genome coverage and abundances:

$$56 \text{ Feature 6} = \sqrt{\frac{U_i}{E_i} * \frac{\tilde{Q}_i}{R}}, \quad (12)$$

$$57 \text{ Feature 7} = G_i = \sqrt{Q_i * U_i}, (\text{G - score}) \quad (13)$$

58

59 We then employed min-max scaling and log transformation separately for any combination of
60 the above seven features before passing to six classifiers in the scikit-learn 1.1.0 (Logistic
61 Regression, Support-Vector Machines, naive Bayes, K-neighbors, AdaBoost, as well as

62 Gradient Boosting Classifier) to select the best feature combination and the best classifier.
63 Specifically, for each feature combination using each normalization and classifier, we
64 performed five-fold cross-validation five times using the CAMI2 simulation datasets. As a result,
65 we found that the Random Forest classifier using log transformed features 1,2,5, and 7 (i.e.,
66 genome coverage C_i , taxonomic count N_i , sequence count R_i , and G-score G_i) has the best
67 performance (**Fig. S8**).
68

69 **Generating a GTDB version ground truth for CAMI 2 simulation datasets**

70 CAMI 2 simulation datasets were synthesized using novel assembled microbial genomes as
71 source genomes with only Refseq annotations³. By merging CAMI 2's source genomes with
72 GTDB and annotating them with GTDB-TK, we were able to generate a GTDB version ground
73 truth for MAP2B as well as improve the machine learning model's training accuracy (please
74 see our GitHub repository for ground truth, see
75 Manuscript/Figures/FigureS9/CAMI_50_abundance_change_in_ground_truth_100W.zip,
76 "GT_Abd").
77

78 **Validation of the false positive and true positive classifier**

79 As the most important component of MAP2B, the classification of false positives and true
80 positives largely determines the accuracy of the final profiling result. To best utilize the four
81 key features (genome coverage, sequence count, taxonomic count, and G-score) for
82 distinguishing true positives from false positives, we trained a Random Forest model based
83 on these features using all the simulated metagenomes from CAMI2. Specifically, we randomly
84 selected 80% of the samples in each of the three CAMI2 datasets (marine, plant-associated,
85 and strain madness datasets) to train a Random Forest model with default parameters, then
86 tested its performance using the remaining 20% of samples. We repeated the whole process
87 of five-fold cross-validation 50 times by randomly assigning samples in either train or test folds.
88 When evaluating the performance of the model, the low-abundance species in the ground
89 truth were gradually filtered out according to varying abundance thresholds from 10^{-6} to
90 10^{-4} . This is because the sensitivity for species identification can be limited due to low
91 sequencing depth (e.g., ~2GB/sample for the strain madness dataset). Some state-of-the-art
92 metagenomic profilers^{4,5} actually set the default abundance threshold as 10^{-4} . Note that the
93 minimum abundance in the CAMI2 datasets is 2×10^{-6} . Therefore, using threshold 10^{-6} is
94 equivalent to not setting any threshold.

95 We separately evaluated the performance of the Random Forest classification model on
96 each of the three simulation datasets. For the marine dataset, when using threshold 10^{-6} (or
97 equivalently, without setting any threshold), we achieved Accuracy~0.988, AUROC~0.999,
98 AUPRC~0.999, Precision~0.991, Recall~0.990, and F1 score~0.990, respectively (**Fig. S9a-**
99 **f**). We then found that filtering out the species in the ground truth with abundance less than
100 10^{-5} will further increase the performance of the model. For example, in the marine dataset,
101 the average of Accuracy, AUROC, AUPRC, Precision, Recall, and F1 increased to 0.989,
102 0.999, 0.999, 0.991, 0.990, and 0.991, respectively. Furthermore, filtering out species with

103 relative abundance less than 10^{-4} in the ground truth can maximize the performance of the
104 model for the marine datasets (Accuracy, AUROC, AUPRC, Precision, Recall, and F1 are
105 0.992, 1, 1, 0.991, 0.995, and 0.993, respectively). A similar trend was observed for the plant-
106 associated (green lines and dots in **Fig. S9**) and strain madness datasets (yellow lines and
107 dots in **Fig. S9**). The average Accuracy, AUROC, AUPRC, Precision, Recall, and F1 for the
108 three datasets (when using 10^{-4} as the threshold) are 0.993, 1, 0.997, 0.958, 0.975, and
109 0.966, respectively. Finally, the well-trained classifier with the best performance among 50
110 repeats will be used in MAP2B.

111

112 Usage of metagenomic profilers

113 For WMS data, we compared MAP2B with five state-of-the-arts metagenomic profilers:
114 MetaPhlAn4⁶, mOTUs3⁷, Bracken⁵, Kraken2⁴, and KrakenUniq⁸. The detailed procedures are
115 listed below.

116 (1) MetaPhlAn4 (v4.0.1) is a marker-gene alignment approach that relies on a precomputed
117 databases containing clade-specific marker genes. Query reads are aligned via bowtie2
118 to the marker genes for microbial identification and abundance estimation. The database
119 version used is mpa_vJan21_CHOCOPhlAnSGB_202103. The following MetaPhlAn4
120 command was used.

```
121 "metaphlan input_1.fastq.gz,input_2.fastq.gz --input_type fastq --bowtie2out  
122 output.bz2 --tax_level s --nproc 32 -o output.txt"
```

123 (2) mOTUs3 (v3.0.3, database version v3.0.3) is a marker-based method that compiles a
124 large variety of phylogenetic marker genes from multiple biomes. Query reads are aligned
125 using bwa mem and further processed to generate an abundance profile. The following
126 mOTUs3 command was used.

```
127 "motus profile -f input_1.fastq.gz -r input_2.fastq.gz -n sample_name -u -p -k mOTU -  
128 o output.txt -t 32"
```

129 (3) Kraken2 (v2.1.1) is a k-mers based taxonomic classification method. It searches for 35bp
130 k-mers from the query sequence in a precomputed database that matches k-mers to the
131 lowest common ancestor (LCA) taxon of all genomes that contain that taxon. The database
132 was constructed using complete bacterial, archaeal, human and viral genomes from NCBI
133 RefSeq (2020 Dec). A filtering abundance threshold of 0.01 (default) was selected. The
134 following Kraken2 command was used.

```
135 "kraken2 --threads 32 --fastq-input --gzip-compressed --paired input_1.fastq.gz  
136 input_2.fastq.gz --output output.reads --report output.report"
```

137 (4) Bracken (v2.5) utilizes the read classification output from standard Kraken for a Bayesian
138 re-estimation of taxonomic abundances, which significantly improves the false-positive
139 issue of standard Kraken and implicitly normalizes for genome length. The kraken-filter
140 was used to filter raw classifications at the 0.01 threshold. The below Bracken command
141 was used.

```
142 "est_abundance.py -i output_kraken2.report -k db -o output"
```

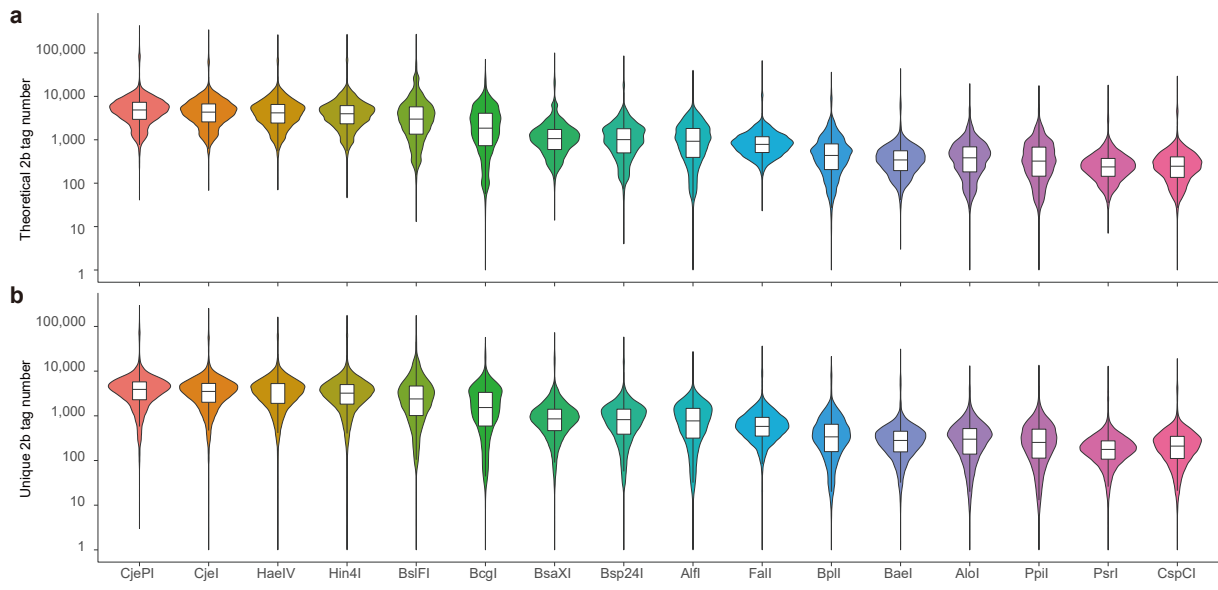
143 (5) KrakenUniq (v0.5.8) is a metagenomics classifier that combines the fast k-mer-based
144 classification of Kraken with an efficient algorithm for assessing the coverage of unique k-
145 mers found in each species in a dataset. The database was constructed using complete

146 bacterial, archaeal and viral genomes from NCBI RefSeq (2022 Jan). A filtering abundance
147 threshold of 0.01 (same with kraken2) was selected. The KrakenUniq command below
148 was used.

149 `krakenuniq --db db --threads 32 --report-file output.report --gzip-compressed`
150 `input_1.fastq.gz input_2.fastq.gz --fastq-input`

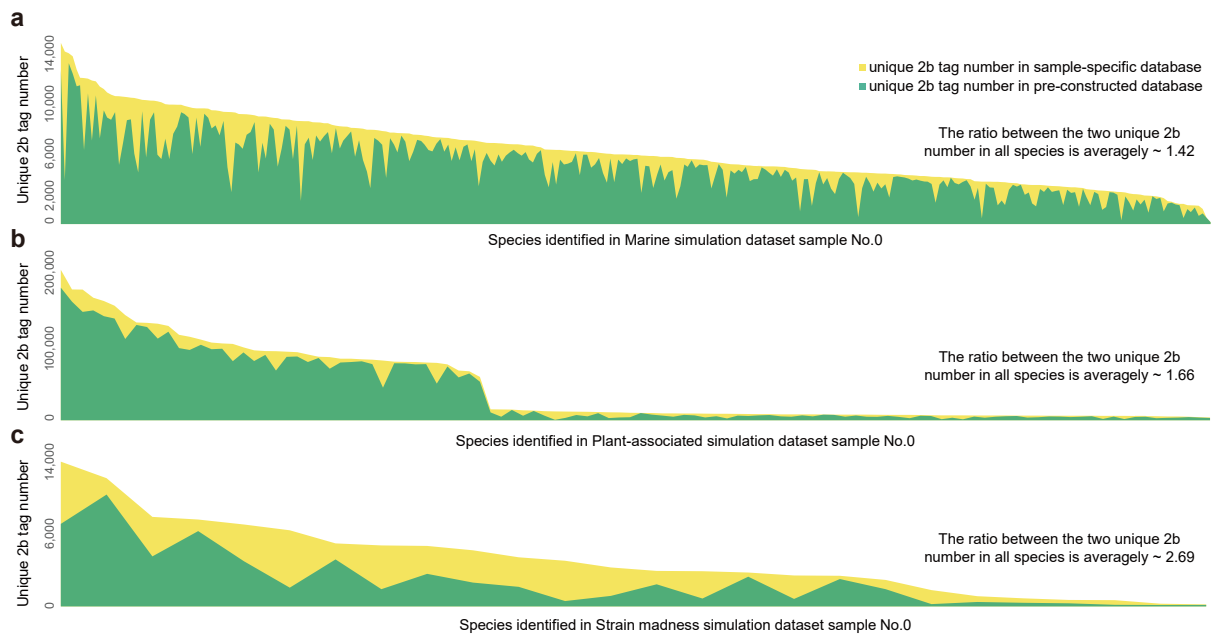
151 (6) As for MAP2B (v1), the default MAP2B command of `perl MAP2B.pl -t 2 -l data.list -d`
152 `MAP2B_DB_GTDB -o output` was used.

153 **Supplemental Figures**



154

155 **Figure S1. Distribution of theoretically existent and unique 2b tags for all the 48,475**
 156 **species in the combined database of GTDB and Ensembl Fungi. (a)** The theoretically
 157 existent 2b tags (H_i) were generated by *in silico* digestion (using 16 different Type IIB enzymes)
 158 for 48,475 species' 258,406 microbial genomes downloaded from GTDB and Ensembl Fungi.
 159 For species- i , its theoretically existent 2b tags in the integrated database of GTDB and
 160 Ensembl Fungi is denoted as H_i . **(b)** We then selected those 2b tags that are not
 161 duplicated/overlapped between any two species and named them as unique 2b tags. For
 162 species- i , its unique 2b tags in the integrated database of GTDB and Ensembl Fungi is
 163 denoted as E_i . The Type IIB restriction enzymes in X-axis are sorted by the median H_i (or
 164 E_i) in descending order. In this paper, we used *CjePI* as the Type IIB enzyme for *in silico*
 165 digestion since it has the highest median H_i and E_i . Using multiple IIB enzymes has limited
 166 improvement in the accuracy of species identification and abundance estimation.



167

168

169

170

171

172

173

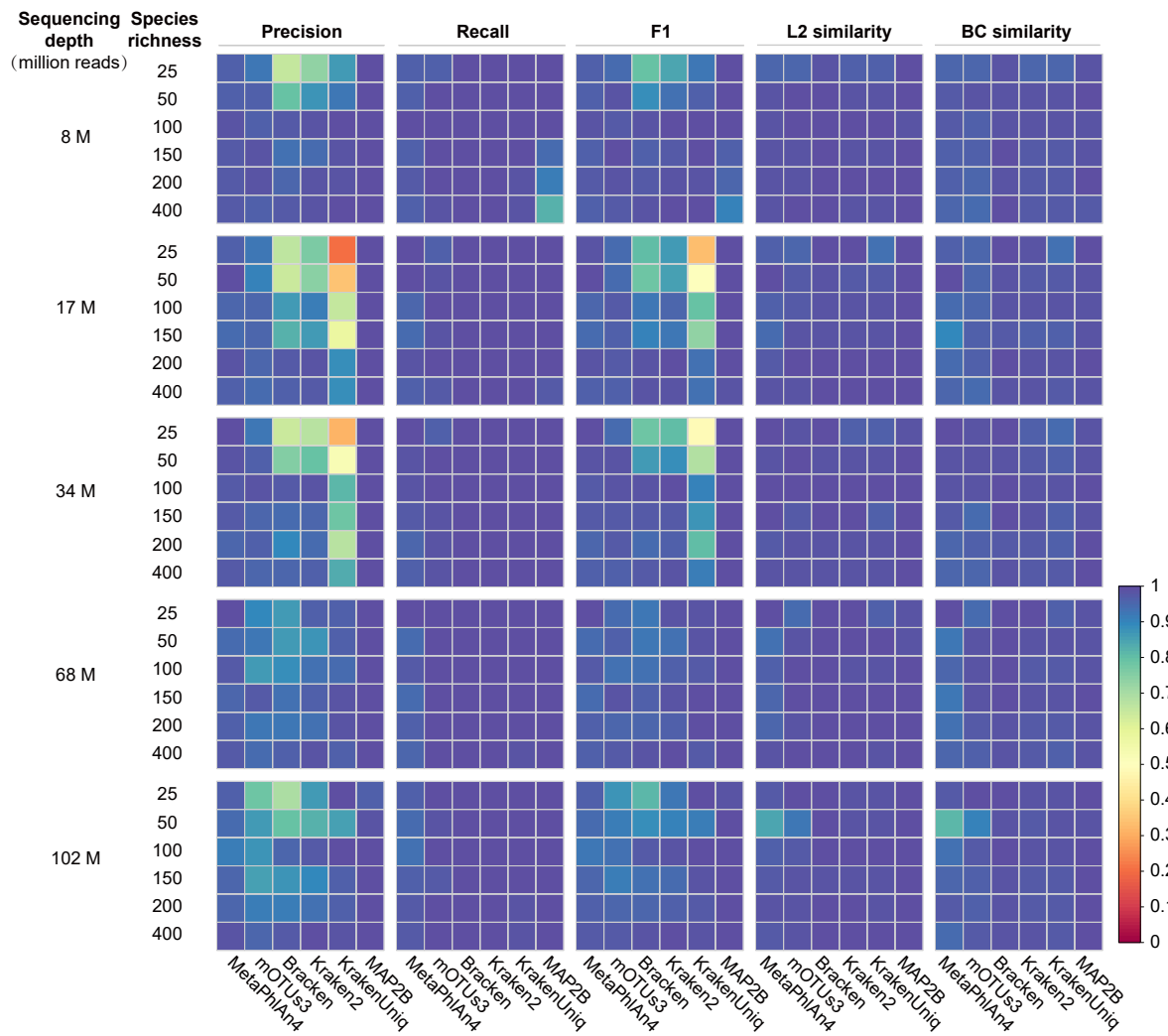
174

175

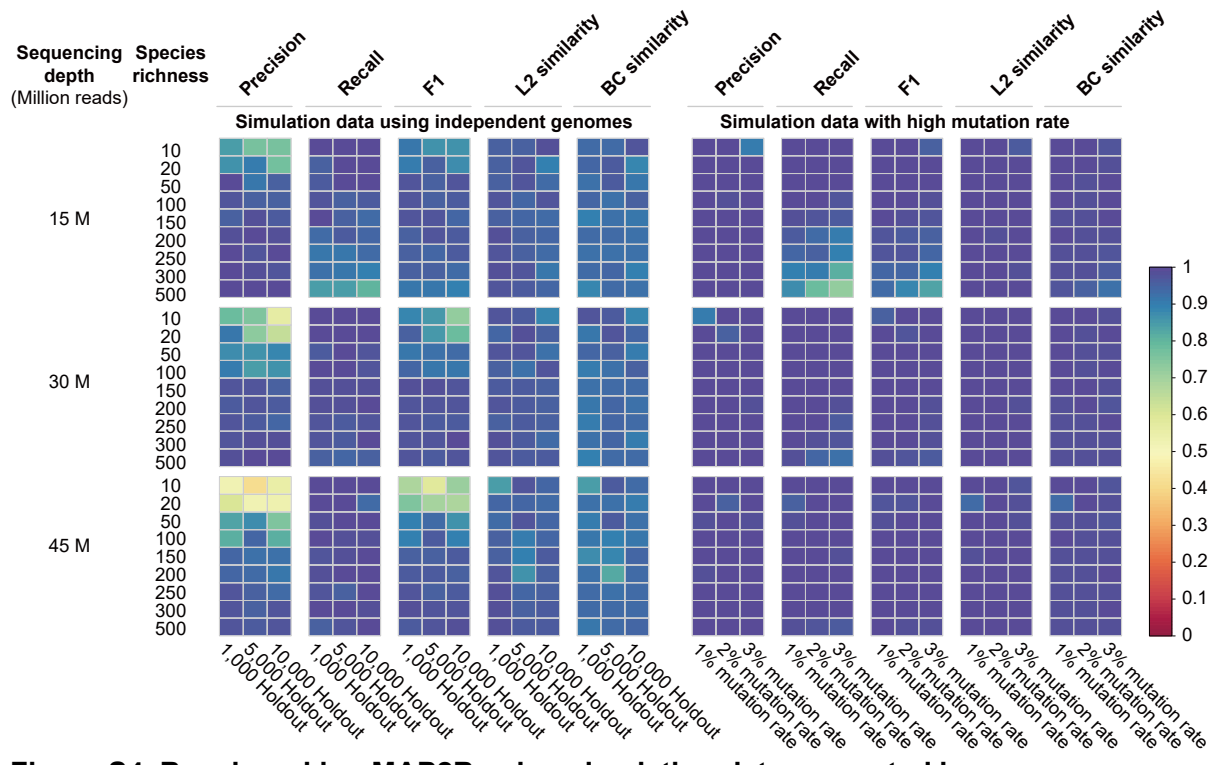
176

177

Figure S2. Sample-specific unique 2b database can largely expand the number of unique 2b tags compared to pre-constructed unique 2b database. To demonstrate the advantage of using a sample-specific unique 2b tag database for the second-round reads alignment, we consider the three CAMI2 simulation datasets of (a) marine, (b) plant-associated, and (c) strain madness as examples. We compare the number of unique 2b tags for all identified species between the preconstructed unique 2b tag database and sample-specific database and found that the average fold change for all the identified species is 1.42, 1.66, and 2.69 in the three datasets, respectively. The former contains unique 2b tags generated by comparing theoretically existent 2b tags among 48,475 species. By contrast, the latter usually contains twice unique 2b tags selected from a few hundreds of species.



178
 179 **Figure S3. Comparing MAP2B with other profilers using simulation data generated by**
 180 **shared genome sources of different profilers.** To minimize the influence of different
 181 reference databases on the evaluation, we selected the shared microbial genomes between
 182 different metagenomic profilers (e.g., mOTU2, MetaPhlaAn, and Kraken) as genomes to
 183 simulate WMS data. From left to right, the profiling results generated by different metagenomic
 184 profilers were compared with ground truth and illustrated by the Precision, Recall, F1 score,
 185 L2 similarity, and BC similarity. From top to bottom, the simulated sequencing depth increases
 186 from 8M to 102M, and the species richness increases from 25 to 400 under each sequencing
 187 depth. Since selecting the intersection of different metagenomic profilers' reference genomes
 188 dramatically decreased the number of source genomes for simulation, we slightly adjusted the
 189 species number and sequencing depth compared to Fig. 3.



190

191

192

193

194

195

196

197

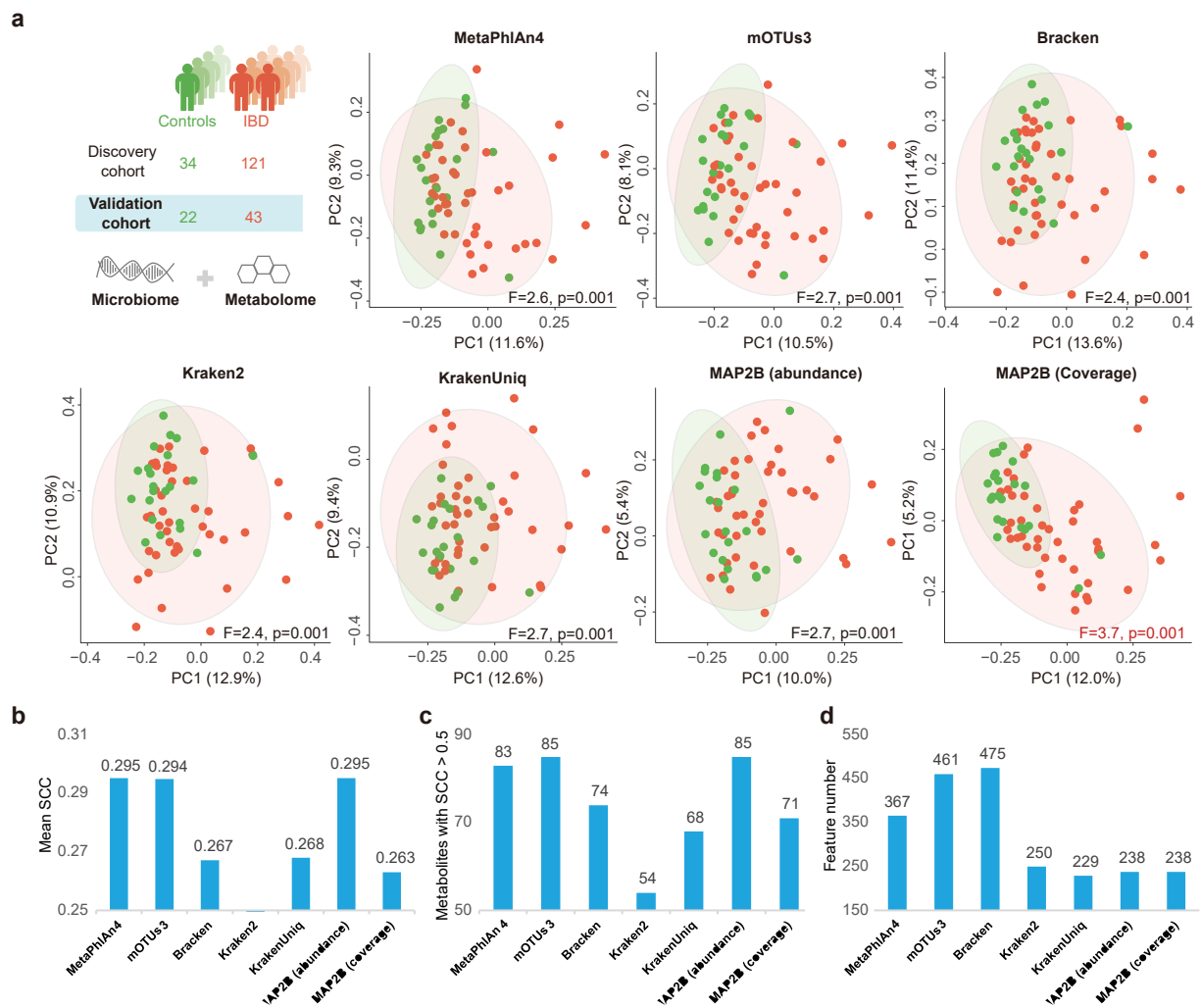
198

199

200

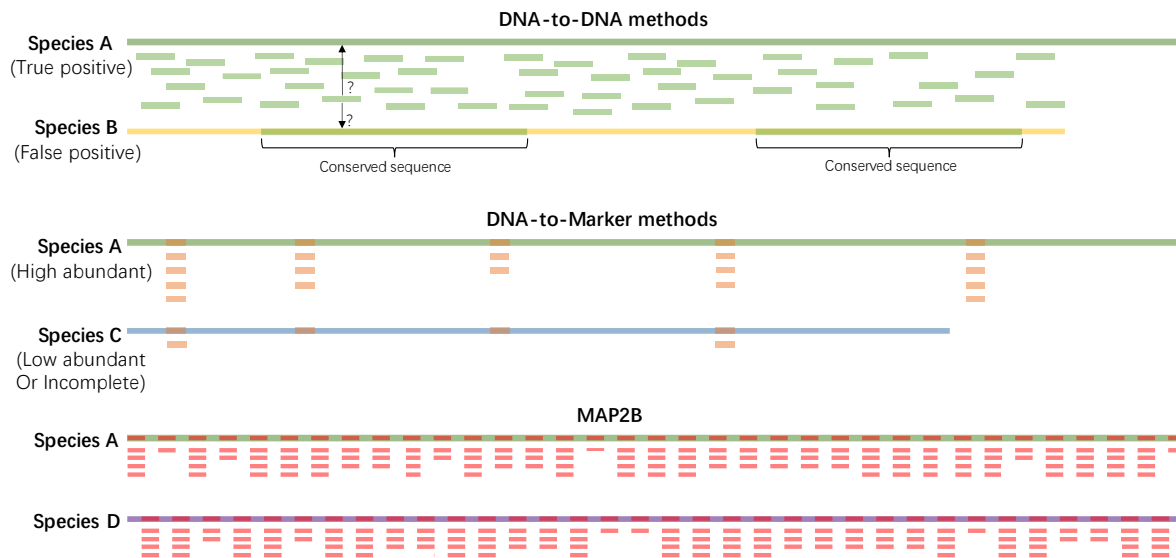
201

Figure S4. Benchmarking MAP2B using simulation data generated by genome sources independent of GTDB or with high a mutation rate. To overcome the challenge of accurately estimating species abundance in the absence of microbial genomes in the reference database, we implemented a systematic partitioning approach for the GTDB database. We utilized held-out genomes to simulate whole metagenome sequencing (WMS) data and systematically varied the mutation rate in the simulated data. The resulting profiles were compared to ground truth using Precision, Recall, F1 score, L2 similarity, and BC similarity metrics. We tested different numbers of holdout genomes and mutation rates, and the results were illustrated from left to right. In this evaluation, we also increased the simulated sequencing depth from 15M to 45M and increased species richness from 10 to 500 for each sequencing depth.



202
 203
 204
 205
 206
 207
 208
 209
 210
 211
 212
 213
 214

Figure S5. Assessing MAP2B's microbial profiles in disease discrimination (validation cohort) and prediction of metabolomic profiles. (a) PCoA plots for the validation cohort (n=65) based on the taxonomic profiles generated by different profilers. The ellipses with 95% CI are drawn to illustrate the difference between IBD (red dots) and non-IBD (green dots) in PCoA. F values of the PERMANOVA are also marked on the bottom of each plot to quantify the difference in disease status. (b) Comparison of prediction results by using different taxonomic profiling via mean SCC of the metabolite between its true values and predicted values across all individuals in the validation cohort. (c) Comparison of the number of metabolites with SCCs larger than 0.5 among different taxonomic profiling results. (d) Comparison of the number of taxonomic features used by different metagenomic profilers in the prediction for metabolomic profiles. The prediction results in (b) - (c) were generated by the mNODE which is in line with MiMeNet (Fig. 5b-c).



215

216

Figure S6. Comparison of the algorithm implemented in MAP2B and conventional

217

metagenomic profiling tools. Using the whole genome as a reference, DNA-to-DNA

218

methods (such as Bracken, Kraken, CLARK, Centrifuge, and PathSeq) may be confused by

219

multi-alignments in conserved sequences, leading to a high rate of false positives. Although

220

DNA-to-Marker methods (such as MetaPhlan and mOTUs) can naturally avoid this issue, they

221

may be limited by the availability of universal markers, such as missing markers from

222

incomplete microbial genomes during database construction, high marker similarity among

223

conspecific taxa during database construction (and sequencing), and undetectable markers

224

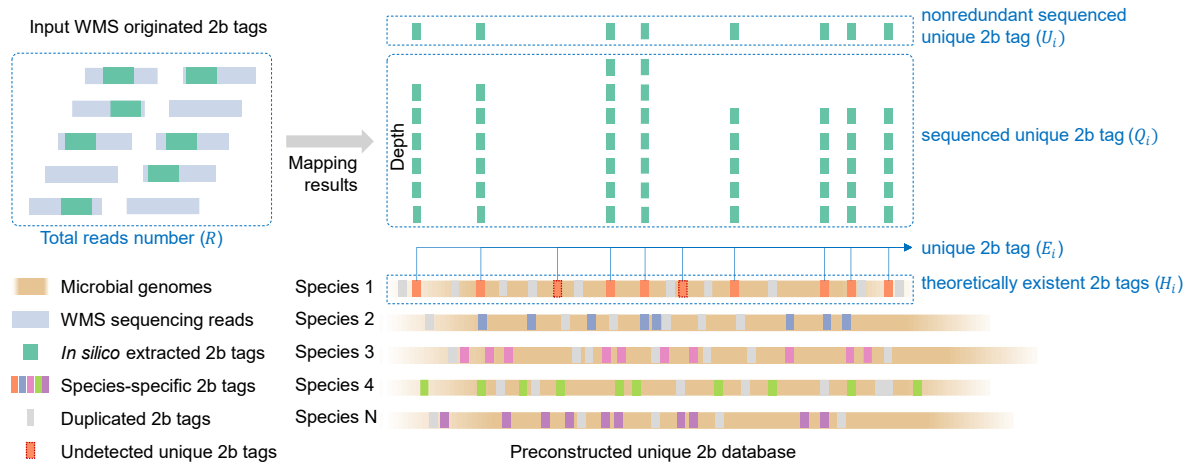
in low abundant taxa during sequencing. MAP2B is not relying on the whole genome or

225

universal marker genes as references. Using species-specific 2b tags can also avoid the multi-

226

alignment issue while providing ample small markers for species identification.



227

228

229

230

231

232

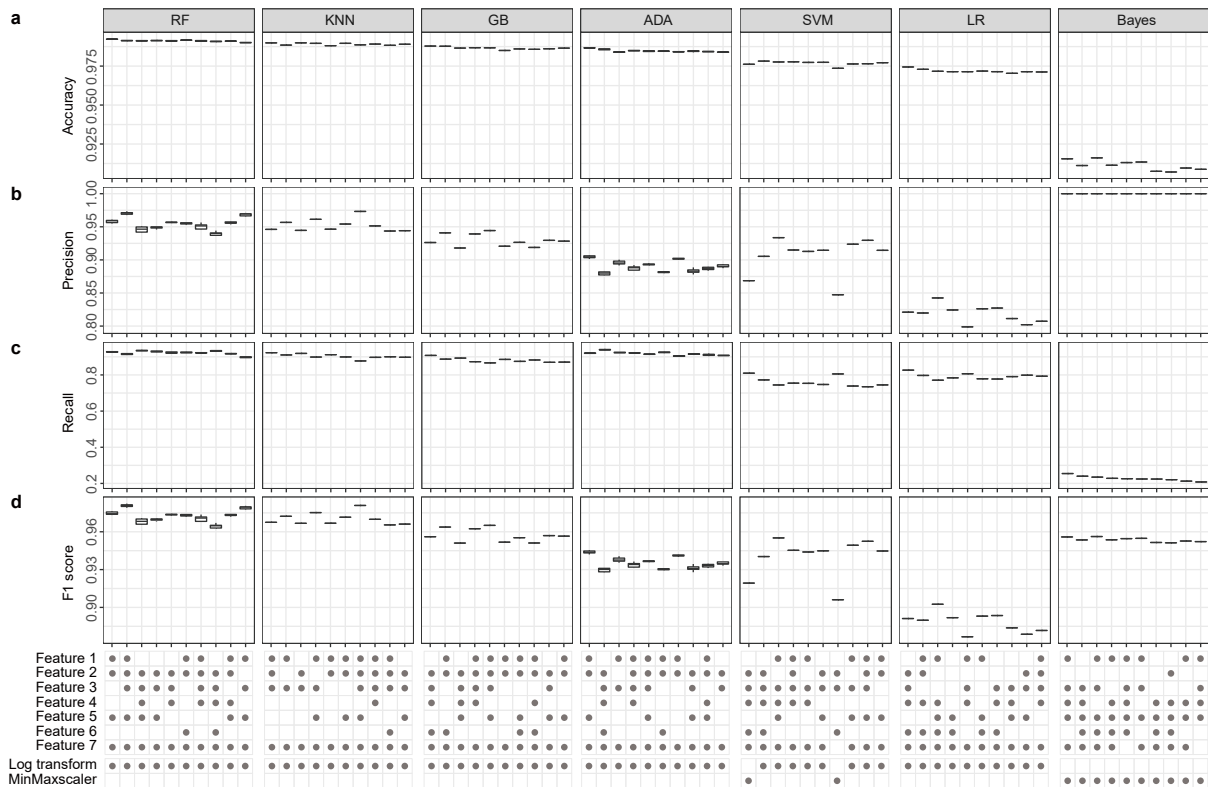
233

234

235

236

Figure S7. Definition of the terminology used in the species identification procedures of MAP2B. For species- i , we denote its total number of 2b tags generated by *in silico* digestion of its genome as H_i . Among the H_i tags, there are E_i tags that are single-copy within the genome of species- i , and are unique to species- i w.r.t all other species in the database of microbial genomes. Given an input WMS sequencing dataset, we *in silico* extract 2b tags, map them to the species-specific 2b tags, and denote the number of tags unique to species- i as Q_i . Among the Q_i tags unique to species- i , there are U_i distinct or nonredundant ones. R is the total number of reads in the WMS sequencing data, which can vary a lot across different samples.



237

238

239

240

241

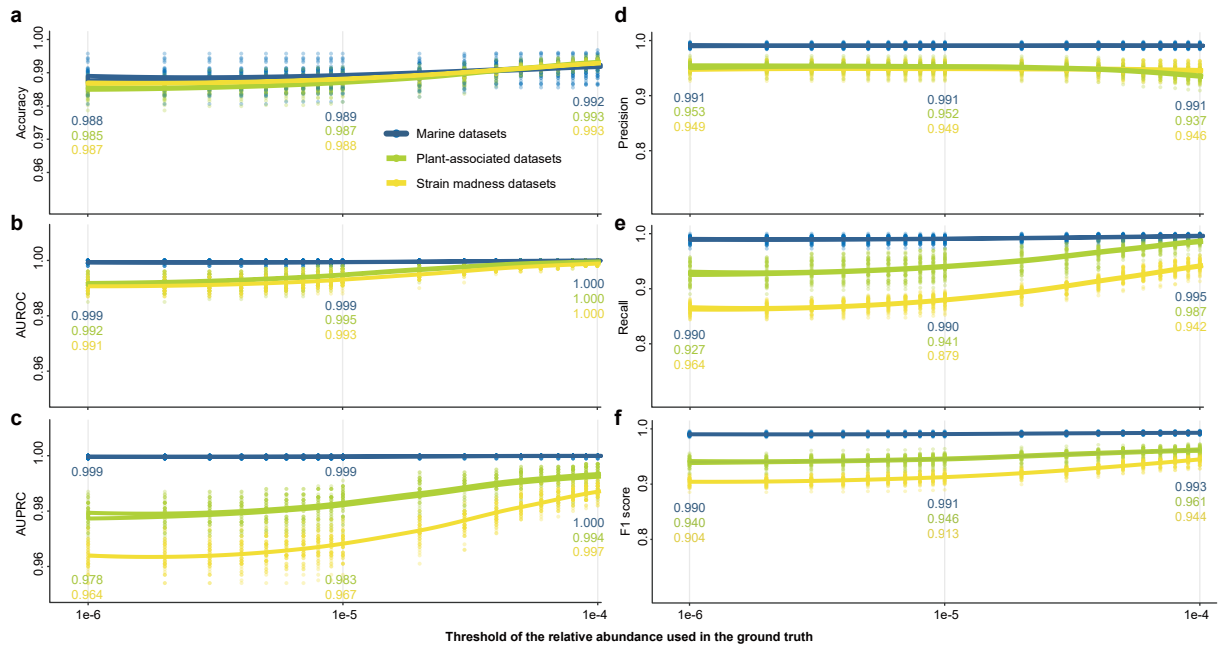
242

243

244

245

Figure S8. TOP-10 feature combinations for each classifier in discriminating false positives from true positives. To select the best feature combination, and the best classifier, both min-max scaling and log transformation were used separately for any combination of the seven features before passing to seven classifiers: Random Forest (RF), K-neighbors (KNN), Gradient Boosting (GB), AdaBoost (ADA), Support-Vector Machines (SVM), Logistic Regression (LR), and naive Bayes (Bayes). For each feature combination and each classifier, we performed five-fold cross-validation five times using the CAMI2 simulation datasets and compared their performance via (a) Accuracy, (b) Precision, (c) Recall, and (d) F1 score.



246

247

248

249

250

251

252

253

254

Figure S9. Performance evaluations of the machine learning model on false positive recognition based on 50 times five-fold cross-validation. When evaluating the performance of the model, the low-abundance species in the ground truth were filtered out according to different abundance thresholds. We gradually discard the true species with relative abundance from 10^{-6} to 10^{-4} and illustrate the performance of the model in determining false positives using metrics such as (a) the Area Under the Receiver Operating Characteristic (AUROC) curve; (b) the Area Under the Precision-Recall Curve (AUPRC); (c) Accuracy; (d) Precision; (e) Recall; and (f) F1-score.

255 **Reference**

- 256 1. Parks, D.H. et al. A complete domain-to-species taxonomy for Bacteria and Archaea.
257 *Nat Biotechnol* **38**, 1079-1086 (2020).
- 258 2. Yates, A.D. et al. Ensembl Genomes 2022: an expanding genome resource for non-
259 vertebrates. *Nucleic Acids Res* **50**, D996-D1003 (2022).
- 260 3. Meyer, F. et al. Critical Assessment of Metagenome Interpretation: the second round
261 of challenges. *Nature methods* **19**, 429-+ (2022).
- 262 4. Wood, D.E., Lu, J. & Langmead, B. Improved metagenomic analysis with Kraken 2.
263 *Genome biology* **20**, 257 (2019).
- 264 5. Lu, J., Breitwieser, F.P., Thielen, P. & Salzberg, S.L. Bracken: estimating species
265 abundance in metagenomics data. *PeerJ Computer Science* **3**, e104 (2017).
- 266 6. Blanco-Miguez, A. et al. Extending and improving metagenomic taxonomic profiling
267 with uncharacterized species using MetaPhlAn 4. *Nat Biotechnol* (2023).
- 268 7. Ruscheweyh, H.J. et al. Cultivation-independent genomes greatly expand taxonomic-
269 profiling capabilities of mOTUs across various environments. *Microbiome* **10**, 212
270 (2022).
- 271 8. Breitwieser, F.P., Baker, D.N. & Salzberg, S.L. KrakenUniq: confident and fast
272 metagenomics classification using unique k-mer counts. *Genome biology* **19**, 198
273 (2018).