

1 rRNA Operon Improves Species-Level Classification of Bacteria and Microbial Community Analysis
2 Compared to 16S rRNA

3

4 Sohyoung Won^{a,b}, Seoae Cho^b, Heebal Kim^{a,b,c#}

5

6 ^aInterdisciplinary Program in Bioinformatics, Seoul National University, Seoul, Republic of Korea

7 ^beGnome, Inc, Seoul, Republic of Korea

8 ^cDepartment of Agricultural Biotechnology and Research Institute of Agriculture and Life Sciences,
9 Seoul National University, Seoul, Republic of Korea

10

11

12 Running Head: rRNA operon and 16S rRNA in Bacteria Species Analysis

13

14 # Address correspondence to Heebal Kim, heebal@snu.ac.kr.

15

16

17 **ABSTRACT**

18 Precise identification of species is fundamental in microbial genomics, crucial for understanding the
19 microbial communities. While the 16S rRNA gene, particularly its V3-V4 regions, has been
20 extensively employed for microbial identification, however has limitations in achieving species-level
21 resolution. Advancements in long-read sequencing technologies have highlighted the rRNA operon as
22 a more accurate marker for microbial classification and analysis than the 16S rRNA gene. This study
23 aims to compare the accuracy of species classification and microbial community analysis using the
24 rRNA operon versus the 16S rRNA gene. We evaluated the species classification accuracy of the
25 rRNA operon, 16S rRNA gene, and 16S rRNA V3-V4 region using a BLAST based method and a *k*-
26 mer matching based method with public data available from NCBI. We further performed simulations
27 to model microbial community analysis. We assessed the performance using each marker in
28 community composition estimation and differential abundance analysis. Our findings demonstrate that
29 the rRNA operon offers an advantage over the 16S rRNA gene and its V3-V4 region for species-level
30 classification within genus. When applied to microbial community analysis, the rRNA operon enables
31 a more accurate determination of composition. Using the rRNA operon yielded more reliable results
32 in differential abundance analysis as well.

33 **IMPORTANCE**

34 We quantitatively demonstrated that the rRNA operon outperformed the 16S rRNA and its V3-V4
35 regions in accuracy, for both individual species identification and species-level microbial community
36 analysis. Our findings can provide guidelines for selecting appropriate markers in the field of
37 microbial research.

38 INTRODUCTION

39 Accurate taxonomic classification is crucial for reliable outcomes in microbial genomics research. As
40 analysis increasingly shifts towards species-level identification beyond the genus level, enhancing the
41 resolution of microbial identification becomes critical for discerning specific species (1). This plays a
42 significant role in discovering novel microbial species and fostering a comprehensive understanding
43 of microbial communities (1).

44 Next-generation sequencing (NGS) technologies have revolutionized microbial genomics, enabling
45 rapid and cost-effective sequencing of whole genomes and amplicons (2). Second-generation
46 sequencing platforms, like Illumina's HiSeq and MiSeq, generate millions of short reads (100-300 bp)
47 (3), while third-generation technologies, like PacBio's SMRT and Oxford Nanopore's MinION,
48 produce significantly longer reads (up to 100 kb or more) (4, 5).

49 16S rRNA gene sequencing is a widely used method for microbial identification and community
50 profiling (6, 7). It targets the highly conserved 16S rRNA gene, containing variable regions among
51 species. Some second-generation sequencing approaches using only specific variable regions (e.g., V3
52 and V4) as markers can be cost-effective but have limitations in taxonomic resolution (8). Even
53 utilizing the entire 16S rRNA gene, accurate species-level classification remains challenging,
54 potentially underestimating diversity and hindering accurate characterization of microbial
55 communities (9, 10).

56 The emergence of third-generation sequencing technologies has enabled the analysis of larger
57 genomic regions, paving the way for whole rRNA operon sequencing as a prominent approach (11).
58 Encompassing the 16S, 23S, and 5S rRNA genes, along with the Internal Transcribed Spacer (ITS)
59 regions, the rRNA operon provides a comprehensive framework for microbial identification and
60 phylogenetic studies (12). Compared to 16S rRNA sequencing, rRNA operon sequencing offers richer
61 information content, promising higher-resolution taxonomic classification, reaching the species level
62 and more accurate microbiome community analysis (13). However, further quantitative research is

63 required to fully validate these expectations.

64 This study utilizes public data to compare the accuracy of species classification within the same genus
65 using the entire rRNA operon sequence, the 16S rRNA sequence, and the V3 and V4 regions of the
66 16S rRNA. Additionally, we create simulated microbiome community data to compare how accurately
67 each region determines the proportion of each species. The aim is to provide guidelines for selecting
68 marker regions for bacterial species classification and species-level microbiome studies.

69

70 **RESULTS**

71 **Species classification accuracy within genus.**

72 Both BLAST and k-mer matching methods demonstrated significantly higher accuracy when utilizing
73 the entire rRNA operon compared to the 16S rRNA alone (Fig. 1). The average accuracy for BLAST-
74 based classification using the rRNA operon reached 0.999, with a standard deviation of 0.005. This
75 accuracy dropped to 0.936 with a standard deviation of 0.108 when using the 16S rRNA, and further
76 decreased to 0.689 with a standard deviation of 0.300 with the 16S rRNA V3-V4 regions. This trend
77 reflects that analyzing broader genomic regions leads to improved accuracy and reduced variability.

78 k-mer matching yielded comparable results. The average accuracy using the rRNA operon was 0.999,
79 exceeding the 0.918 observed for the 16S rRNA and 0.693 for the V3-V4 regions. The rRNA operon
80 also displayed the lowest standard deviation (0.006), compared to 0.123 for the 16S rRNA and 0.297
81 for the V3-V4 regions.

82 Across both methods, the *Haemophilus* genus exhibited the lowest accuracy with the rRNA operon,
83 which was 0.960. For the 16S rRNA, the lowest accuracy was observed in the *Serratia* genus, with
84 BLAST and k-mer matching methods reporting 0.402 and 0.496, respectively. Notably, employing the
85 rRNA operon compared to the 16S rRNA consistently achieved higher accuracy for all genera when
86 using the k-mer matching method. The BLAST method presented a single minor exception in

87 *Chlamydia*, where 16S rRNA yielded marginally higher accuracy with a difference of only 0.0002.

88 On average, the rRNA operon achieved a classification accuracy 0.084 higher with BLAST and 0.109
89 higher with k-mer matching compared to the 16S rRNA. The largest observed difference in a single
90 genus reached a substantial gap of 0.503 (BLAST) and 0.598 (k-mer matching). Using the rRNA
91 operon, the BLAST method achieved perfect species classification accuracy (1.0) in 89.6% (43) of
92 genera, and the k-mer match method did so in 83.3% (40) of genera. In contrast, with the 16S rRNA,
93 the BLAST method had less than 0.9 accuracy in 31.3% (15) of genera, and the k-mer matching
94 method in 37.5% (18) of genera. This indicates that using the rRNA operon enables more precise
95 species classification than the 16S rRNA.

96 The standard deviation of accuracy with the 16S rRNA was a significant 19.8 times higher (BLAST)
97 and 20.7 times higher (k-mer matching) compared to the rRNA operon. Additionally, the minimum
98 accuracy observed with the rRNA operon consistently exceeded 0.95, whereas the 16S rRNA dipped
99 below 0.5 in some cases.

100 **Microbial community composition prediction.**

101 We conducted simulations to evaluate the effectiveness of different regions for predicting the species
102 compositions (Fig. 2). These simulations assumed species existed in random proportions following a
103 Dirichlet distribution. The figure depicts the predicted proportions of the top 10 species for each
104 method (rRNA operon, 16S rRNA, and 16S rRNA V3-V4 regions). Predictions using the rRNA
105 operon closely matched the actual compositions. The 16S rRNA predictions displayed a similar trend
106 to the actual compositions, but with some discrepancies in the ratios between species. Predictions
107 based on the 16S rRNA V3-V4 regions deviated significantly from the actual compositions.

108 To numerically verify these observed trends, we calculated the Pearson correlation coefficient
109 between the actual and predicted proportions (Table 1). Across six simulations, the correlation
110 between actual and predicted proportions using the rRNA operon remained consistently high, with an

111 average of 0.999 and a standard deviation of 0.0001. This held true regardless of the number of reads
112 used in the simulation. The 16S rRNA exhibited a lower average correlation (0.849) with a higher
113 standard deviation (0.014), indicating a poorer match to the actual proportions and greater variability
114 between simulations compared to the rRNA operon. The 16S rRNA V3-V4 regions yielded the lowest
115 average correlation (0.288) with a standard deviation of 0.022. The correlation between actual and
116 predicted proportions increased with the number of simulated reads for both the 16S rRNA and its
117 V3-V4 regions, reaching a plateau beyond 500,000 reads.

118 To quantify the similarity between the actual and predicted microbial community compositions, we
119 employed the Bray-Curtis distance metric. A smaller Bray-Curtis distance signifies greater similarity
120 between the two datasets. The predicted composition using the rRNA operon yielded a remarkably
121 close distance to the actual composition, averaging only 0.001 across six simulations. Conversely, the
122 average distances observed when utilizing the 16S rRNA and 16S rRNA V3-V4 regions were higher,
123 at 0.088 and 0.367 respectively. Predictions based on the rRNA operon exhibited the closest match to
124 the actual community composition, with distances 71.2 times smaller than those obtained using the
125 16S rRNA.

126 To statistically validate these observations, we conducted an Analysis of Similarities (ANOSIM) test
127 using Bray-Curtis distance. This non-parametric method evaluates the probability of observed
128 differences in similarity between groups arising by chance. The results confirmed these findings.
129 Predictions made with the rRNA operon yielded a p-value of 0.272, indicating no statistically
130 significant difference from the actual community composition. Conversely, predictions utilizing the
131 16S rRNA and 16S rRNA V3-V4 regions produced p-values of 0.006 and 0.004, respectively. These
132 significant p-values ($p < 0.01$) demonstrate that these methods yielded compositions statistically
133 distinct from the actual community.

134 **Microbial community composition and differential abundance in human gut microbiome data.**

135 We evaluated the performance using the rRNA operon, 16S rRNA, and 16S rRNA V3-V4 regions for

136 microbial community composition prediction using real human gut microbiome data (Fig. 3). The
137 analysis included 558 overlapping species from samples of 14 healthy donors and 14 patients. The
138 correlation between the reference species proportions and those predicted using the rRNA operon was
139 remarkably high, averaging 1.00 with a minimal standard deviation of 0.000003. Predictions based on
140 the 16S rRNA and the 16S V3-V4 regions exhibited lower correlations with the reference, with
141 averages of 0.931 and 0.660, and standard deviations of 0.104 and 0.323, respectively. The rRNA
142 operon consistently achieved high correlations between predicted proportions and reference
143 proportions, with the lowest value still exceeding 0.999. Conversely, the 16S rRNA exhibited a
144 significantly lower correlation, dropping as low as 0.620. These results align with our observations
145 from the randomly generated data.

146 We further assessed the methods by conducting differential abundance analyses based on both the
147 reference compositions and those predicted by each classification method. We compared species
148 identified as significantly different in each case (Fig. 4). The reference data identified 132
149 significantly differentially abundant species, which were accurately reflected by the predictions made
150 with the rRNA operon. The 16S rRNA identified 151 significant species, with 127 overlapping with
151 the reference findings. It missed 5 significant species (false negatives) and identified 24 species as
152 significant that were not truly so (false positives). This translates to a false negative rate of 3.79% and
153 a false positive rate of 18.2%, highlighting a higher prevalence of false positives with 16S rRNA. The
154 16S rRNA V3-V4 regions performed even worse, with even greater false negative (22.0%) and false
155 positive (25.8%) rates.

156 Fig. 5 depicts the coefficients of species that were identified as false negatives or false positives when
157 using the 16S rRNA, as well as those whose coefficients differed in direction compared to using the
158 reference. Here, the coefficient represents the relative abundance of a species in patients compared to
159 donors. Some species identified as significantly different by 16S rRNA that were not detected by the
160 reference or the rRNA operon predictions. Additionally, one species exhibited an opposing abundance

161 trend between donors and patients when comparing the reference and 16S rRNA data. Furthermore,
162 the magnitudes of the coefficients often differed between the reference and 16S rRNA findings.
163 Conversely, the results using the rRNA operon displayed a high degree of agreement with the
164 reference data, with consistent signs and magnitudes of coefficients for most species. Only two
165 species showed minor discrepancies.

166

167 **DISCUSSION**

168 The rRNA operon demonstrably outperformed the 16S rRNA gene in terms of species classification
169 accuracy. Statistical tests confirmed this observation. Paired Wilcoxon rank sum tests revealed highly
170 significant differences ($p < 0.0001$ for BLAST and k-mer matching) in accuracy favoring the rRNA
171 operon. Furthermore, the rRNA operon exhibited considerably lower variability in accuracy across
172 genera. This signifies that the rRNA operon offers consistently high and stable classification accuracy
173 across various genera, while the 16S rRNA can yield unreliable results due to substantial variations in
174 accuracy depending on the genus.

175 Simulations revealed a clear advantage for the rRNA operon in predicting species compositions within
176 microbial communities. The correlation between actual and predicted proportions using the rRNA
177 operon consistently outperformed both the 16S rRNA and the 16S rRNA V3-V4 regions. Notably, the
178 correlation with the V3-V4 regions was significantly lower, rendering it unreliable for capturing
179 meaningful relationships with the actual data. When using the rRNA operon versus the 16S rRNA, the
180 difference in correlation between actual and predicted compositions in the data was greater than the
181 difference in accuracy of individual species predictions. Similarly, the Bray-Curtis distance metric
182 further supported the superiority of the rRNA operon. This suggests a far more accurate reflection of
183 the true community structure when employing the rRNA operon compared to the 16S rRNA.

184 Simulations replicating the composition of actual human gut microbiomes yielded consistent results.

185 Notably, this lower accuracy in predicting microbial compositions using 16S rRNA was particularly
186 problematic for patient groups. In the patient data, the average correlation between predicted and
187 reference compositions for the rRNA operon remained at 1.00, while the 16S rRNA only achieved an
188 average of 0.870.

189 The observed difference in accuracy for community composition predictions also impacted the results
190 of differential abundance analyses. When utilizing proportions derived from the rRNA operon, the
191 analysis identified the same 132 significant species as those identified using the reference proportions,
192 indicating perfect agreement. In contrast, the analysis based on 16S rRNA data and 16S rRNA V3-V4
193 region data yielded discrepancies, further solidifying the limitations of these methods for accurate
194 prediction.

195 Using the rRNA operon as a marker provided higher accuracy in individual species classification than
196 using the 16S rRNA or its V3-V4 regions, leading to more accurate community composition
197 predictions and more reliable results in differential abundance analyses. However, sequencing costs
198 may increase with the breadth of the region being read (24). Consequently, the choice of method
199 should consider the required resolution and available budget. The 16S rRNA can be a suitable option
200 when less precision is acceptable or species-level analysis is not necessary and genus-level
201 identification suffices. On the other hand, for research requiring precise species-level analysis, such as
202 discovering biomarkers, utilizing the microbiome for treatments, or other studies necessitating
203 accurate species identification, the rRNA operon is preferable. This is especially true for disease-
204 related microbial community studies, as the accuracy difference in community composition
205 predictions between methods was more pronounced in patient groups, highlighting the importance of
206 using the rRNA operon for more precise species differentiation in such contexts.

207 The accuracy of microbial community composition prediction using the 16S rRNA or its V3-V4
208 regions improves when using more number of reads, which mean sequencing more data. However,
209 this also raises data production costs and should be carefully weighed. Given the same budget,

210 producing less data with the rRNA operon may be more efficient than generating more data with the
211 16S rRNA.

212 Employing the rRNA operon as a marker demonstrably enhances individual species classification
213 accuracy compared to the 16S rRNA gene. This translates to more precise predictions of microbial
214 community compositions and more reliable differential abundance analysis results. The 16S rRNA
215 V3-V4 region exhibited even lower accuracy across all scenarios compared to the full 16S rRNA,
216 highlighting a significant decline in precision. Therefore, for research requiring accurate species
217 classification, employing the rRNA operon as a marker appears to be the most appropriate choice. In
218 microbial community studies aiming for precise species-level analysis, utilizing the rRNA operon is
219 advisable as using the 16S rRNA has its limitations, and relying solely on its V3-V4 regions may
220 make it challenging to achieve meaningful results.

221

222 **MATERIALS AND METHODS**

223 **Data collection.**

224 We collected complete bacterial genomes available in the NCBI database as of November 29, 2023
225 (14). To ensure robust comparisons, we only considered genera containing more than 50 complete
226 genomes. Our final dataset comprised 72 genera, 2,026 species, and 20,314 genome sequences
227 (Supplementary Table 1).

228 **rRNA operon and 16S rRNA sequence extraction.**

229 We utilized riboSeed with its default settings to extract rRNA operon sequences (15). Following
230 identification of rRNA gene regions using the riboscan command, we employed the riboselect
231 command to locate rRNA operon regions containing 16S, 23S, and 5S rRNA. The corresponding
232 rRNA operon sequences were then extracted. Implementing quality control, only sequences within the
233 4,000 to 6,000 base pair range were retained. 16S rRNA sequences were extracted from regions

234 identified as 16S rRNA by the riboscan results.

235 We employed EMBOSS's primer search tool to identify the V3-V4 regions within the previously
236 extracted 16S rRNA sequences (16). The primer sequences used were 'CCTACGGGNGGCWGCAG'
237 for the forward primer and 'GACTACHVGGGTATCTAATCC' for the reverse primer (17). A
238 mismatch percentage of 10% was allowed during the search. Following quality control, only
239 sequences with lengths between 430 bp and 550 bp were retained.

240 **Introduction of sequencing errors.**

241 To simulate real-world applications of using 16S rRNA and rRNA operon sequencing for species
242 classification, we introduced random sequencing errors into the extracted sequences. Error rates were
243 determined by referencing a 2022 study comparing the accuracy of Illumina and ONT technologies
244 (18). 1D ONT MinION read error rates were applied to the rRNA operon sequences, while the
245 average error rate of Illumina's read1 and read2 was used for the 16S rRNA sequences (Table 2).
246 Errors were introduced through random positional mismatches, insertions, and deletions.

247 For each position, a random number between 0 and 1 was generated. If this number was lower than
248 the error rate, an error was introduced. Mismatches involved replacing the original nucleotide with a
249 random one. Implementation was carried out using BioPython SeqIO (19).

250 **Species classification within genus.**

251 Two distinct methods were employed to classify species within the same genus: BLAST alignment
252 score (20) and k-mer matching.

253 The BLAST-based approach used the sequences of the rRNA operon or 16S rRNA (including random
254 errors) as the query, while the original sequences of the extracted regions served as the reference.
255 Nucleotide BLAST was run with default options. Each sequence was classified into the species with
256 the highest bitscore. In cases of ties, one species was randomly chosen for classification.

257 The k-mer matching method benchmarked the approach commonly used in microbiome data
258 classification by Kraken (21). This method involves finding the number of exact matches of 31-mers
259 and classifying the sequence to the species with the most 31-mer matches. Similar to the BLAST
260 approach, ties were resolved by randomly choosing one species for classification.

261 In both methods, when multiple copies of the rRNA operon or 16S rRNA were present, we classified
262 based on the copy with the highest score or the greatest number of matches. We assessed the accuracy
263 of species classification per genus, by calculating the proportion of samples within each genus that
264 were correctly assigned to their respective species.

265 **Simulation in microbial community data.**

266 To evaluate the accuracy of species classification in community data, a simulation was run on
267 community composition data. First, we used a Dirichlet distribution to randomly set proportions for
268 the species in our study and made a mock proportion data. Reads were initially distributed to match
269 these true species proportions. Based on the likelihood of their classification through k-mer matching
270 from the previous analysis, reads were then assigned to species. For example, if species A was
271 correctly classified 90% of the time and misclassified as species B 10% of the time, a read intended
272 for species A would have a 90% chance of being assigned to A and a 10% chance to B.

273 This process was applied to all reads. We conducted simulations for library sizes of 5,000, 10,000,
274 50,000, 100,000, 500,000, 1,000,000 reads.

275 **Microbial community analysis and differential abundance analysis in real world data.**

276 To further validate our findings using real-world gut microbiome data, we performed simulations with
277 publicly available metagenomic proportions. We leveraged gut microbiome data from both donors and
278 recipients of fecal microbiota transplantation (FMT) described in (22). Assuming the reported
279 proportions reflect reality, we predicted species proportions based on the classification accuracy of the
280 rRNA operon, 16S rRNA, and 16S rRNA V3-V4 regions. This procedure mirrored our previous

281 community composition simulation, again using a library size of 100,000 reads. Subsequently, we
282 employed the R package 'Maaslin2' to conduct differential abundance analysis comparing successful
283 donors and pre-FMT patients (23). The analysis utilized AST transformation, TSS normalization, and
284 a linear model. We considered findings with an adjusted false discovery rate (FDR) of less than 0.01
285 to be statistically significant.

286

287 **ACKNOWLEDGMENT**

288 We would like to express our gratitude to eGnome Inc. for their support throughout the course of this
289 research.

290

291

292

293 **REFERENCES**

- 294 1. Duvallet C, Gibbons SM, Gurry T, Irizarry RA, Alm EJ. Meta-analysis of gut microbiome
295 studies identifies disease-specific and shared responses. *Nature communications*. 2017;8(1):1784.
- 296 2. Loman NJ, Pallen MJ. Twenty years of bacterial genome sequencing. *Nature Reviews*
297 *Microbiology*. 2015;13(12):787-94.
- 298 3. Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, et al. Ultra-high-
299 throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. *The ISME*
300 *journal*. 2012;6(8):1621-4.
- 301 4. Eid J, Fehr A, Gray J, Luong K, Lyle J, Otto G, et al. Real-time DNA sequencing from single
302 polymerase molecules. *Science*. 2009;323(5910):133-8.
- 303 5. Deamer D, Akeson M, Branton D. Three decades of nanopore sequencing. *Nature*
304 *biotechnology*. 2016;34(5):518-24.
- 305 6. Langille MG, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA, et al. Predictive
306 functional profiling of microbial communities using 16S rRNA marker gene sequences. *Nature*
307 *biotechnology*. 2013;31(9):814-21.
- 308 7. Johnson JS, Spakowicz DJ, Hong B-Y, Petersen LM, Demkowicz P, Chen L, et al. Evaluation
309 of 16S rRNA gene sequencing for species and strain-level microbiome analysis. *Nature*
310 *communications*. 2019;10(1):5029.
- 311 8. Jeong J, Yun K, Mun S, Chung W-H, Choi S-Y, Nam Y-d, et al. The effect of taxonomic
312 classification by full-length 16S rRNA sequencing with a synthetic long-read technology. *Scientific*
313 *reports*. 2021;11(1):1727.
- 314 9. Poretsky R, Rodriguez-R LM, Luo C, Tsementzi D, Konstantinidis KT. Strengths and
315 limitations of 16S rRNA gene amplicon sequencing in revealing temporal microbial community
316 dynamics. *PloS one*. 2014;9(4):e93827.
- 317 10. Caudill MT, Brayton KA. The use and limitations of the 16S rRNA sequence for species
318 classification of *Anaplasma* samples. *Microorganisms*. 2022;10(3):605.
- 319 11. Kinoshita Y, Niwa H, Uchida-Fujii E, Nukada T. Establishment and assessment of an
320 amplicon sequencing method targeting the 16S-ITS-23S rRNA operon for analysis of the equine
321 gut microbiome. *Scientific reports*. 2021;11(1):11884.
- 322 12. Espejo RT, Plaza N. Multiple ribosomal RNA operons in bacteria; their concerted evolution
323 and potential consequences on the rate of evolution of their 16S rRNA. *Frontiers in microbiology*.
324 2018;9:338498.
- 325 13. Cuscó A, Catozzi C, Viñes J, Sanchez A, Francino O. Microbiota profiling with long
326 amplicons using Nanopore sequencing: full-length 16S rRNA gene and the 16S-ITS-23S of the *rrn*
327 operon. *F1000Research*. 2018;7.
- 328 14. Kitts PA, Church DM, Thibaud-Nissen F, Choi J, Hem V, Sapojnikov V, et al. Assembly: a
329 resource for assembled genomes at NCBI. *Nucleic acids research*. 2016;44(D1):D73-D80.

- 330 15. Waters NR, Abram F, Brennan F, Holmes A, Pritchard L. riboSeed: leveraging prokaryotic
331 genomic architecture to assemble across ribosomal regions. *Nucleic Acids Research*.
332 2018;46(11):e68-e.
- 333 16. Staden R, Judge DP, Bonfield JK. Analyzing sequences using the Staden package and
334 EMBOSS. *Introduction to bioinformatics: a theoretical and practical approach*: Springer; 2003. p.
335 393-410.
- 336 17. Klindworth A, Pruesse E, Schweer T, Peplies J, Quast C, Horn M, et al. Evaluation of
337 general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-based
338 diversity studies. *Nucleic acids research*. 2013;41(1):e1-e.
- 339 18. Zee A, Deng DZ, Adams M, Schimke KD, Corbett-Detig R, Russell SL, et al. Sequencing
340 Illumina libraries at high accuracy on the ONT MinION using R2C2. *Genome research*. 2022;32(11-
341 12):2092-106.
- 342 19. Cock PJ, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely
343 available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*.
344 2009;25(11):1422.
- 345 20. Ye J, McGinnis S, Madden TL. BLAST: improvements for better sequence analysis. *Nucleic*
346 *acids research*. 2006;34(suppl_2):W6-W9.
- 347 21. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact
348 alignments. *Genome biology*. 2014;15:1-12.
- 349 22. Kazemian N, Ramezankhani M, Sehgal A, Khalid FM, Kalkhoran AHZ, Narayan A, et al. The
350 trans-kingdom battle between donor and recipient gut microbiome influences fecal microbiota
351 transplantation outcome. *Scientific reports*. 2020;10(1):18349.
- 352 23. Mallick H, Rahnavard A, McIver LJ, Ma S, Zhang Y, Nguyen LH, et al. Multivariable
353 association discovery in population-scale meta-omics studies. *PLoS computational biology*.
354 2021;17(11):e1009442.
- 355

FIG 1. A boxplot of species classification accuracy across genera using the rRNA operon versus the 16S rRNA gene and its V3-V4 regions: (A) demonstrates the results from the BLAST method, showing a median accuracy of 0.999 (SD: 0.005) for the rRNA operon, 0.936 (SD: 0.108) for the 16S rRNA, and 0.689 (SD: 0.300) for the V3-V4 regions; (B) presents outcomes from the k-mer matching method, with a median accuracy of 0.999 (SD: 0.006) for the rRNA operon, 0.918 (SD: 0.123) for the 16S rRNA, and 0.693 (SD: 0.297) for the V3-V4 regions.

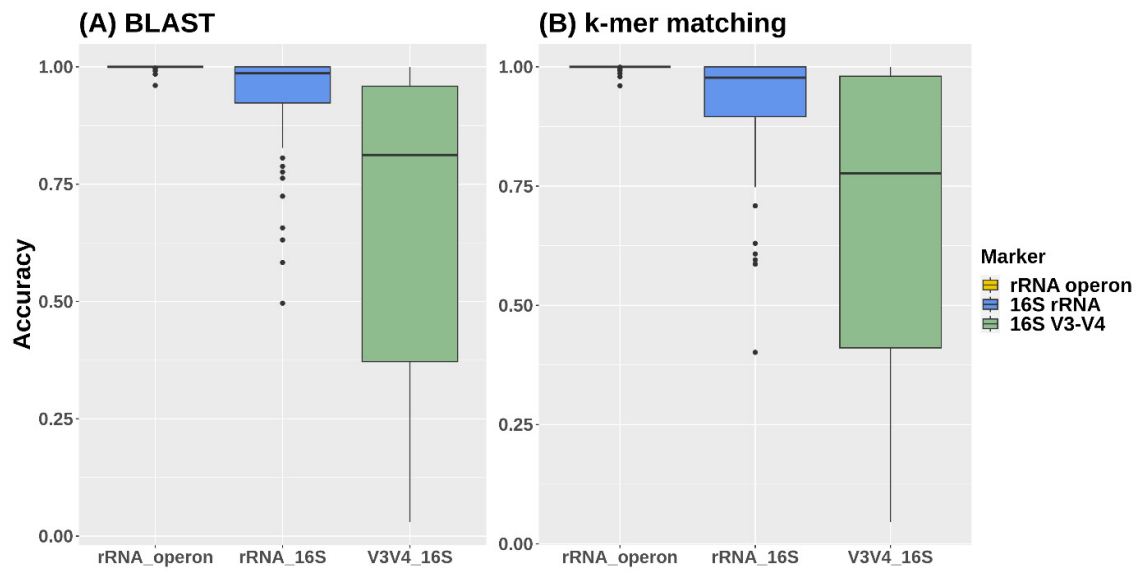


FIG 2. The relative abundance of the top 10 abundant species in the “True” data, where “True” represents the actual proportions. "rRNA operon," "16S rRNA," and "16S V3-V4" show the proportions of species predicted based on the accuracy of species classification within those genomic regions. Each color represents the same species across different predictions, with the x-axis indicating the number of reads used in the simulation and the y-axis showing the proportion of each species.

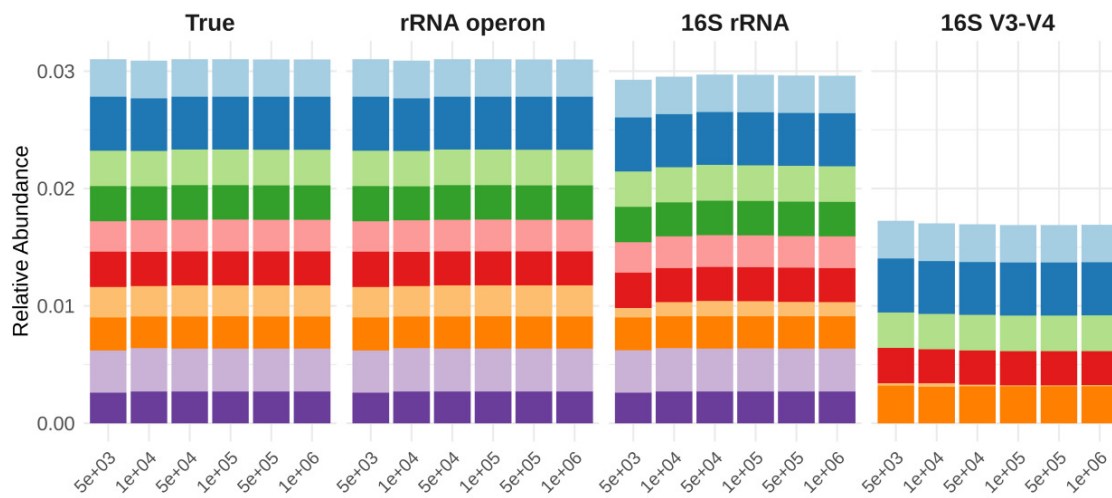


FIG 3. The relative abundance of the top 10 gut microbiome species in (A) FMT donors and (B) patients (B). "Reference" indicates actual proportions. "rRNA operon," "16S rRNA," and "16S V3-V4" show predicted proportions based on species classification accuracy.

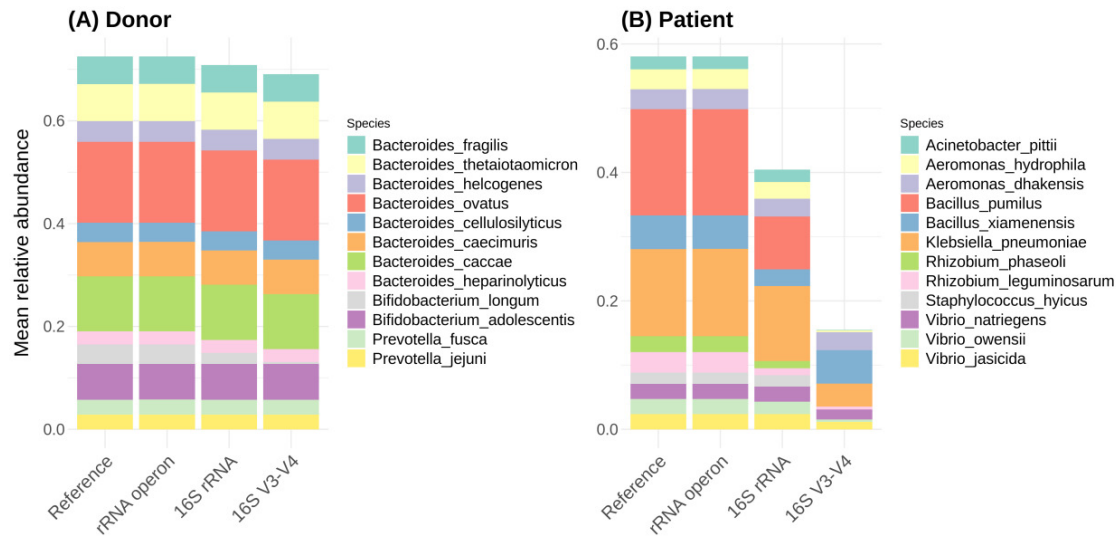


FIG 4. Venn diagrams of the significant species identified through differential abundance analysis based on proportions derived from the reference and rRNA operon, reference and 16S rRNA, and reference and 16S rRNA V3-V4 regions. Overlapping sections of the diagram represent the number of species significantly identified across both methods. The area exclusive to the “Reference” (left side) shows the number of species that were false negatives. Areas unique to each method (right side) indicate false positives. Complete overlap between the “Reference” and a method implies identical species significance findings.

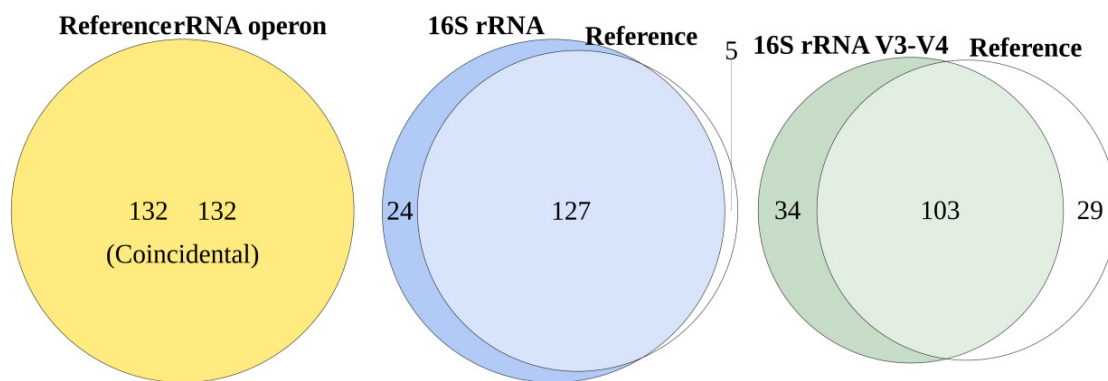


FIG 5. The coefficients from differential abundance analyses using the proportions obtained from reference, rRNA operon, and 16S rRNA. We only showed species that have discrepancy in the reference and 16S rRNA results; species identified as significant in one analysis but not the other and Species with differing direction of coefficient. A positive coefficient (depicted in pink) indicates a species is more abundant in patients, while a negative coefficient (shown in sky blue) suggests it is less abundant. The magnitude of the coefficient signifies the degree of abundance difference.

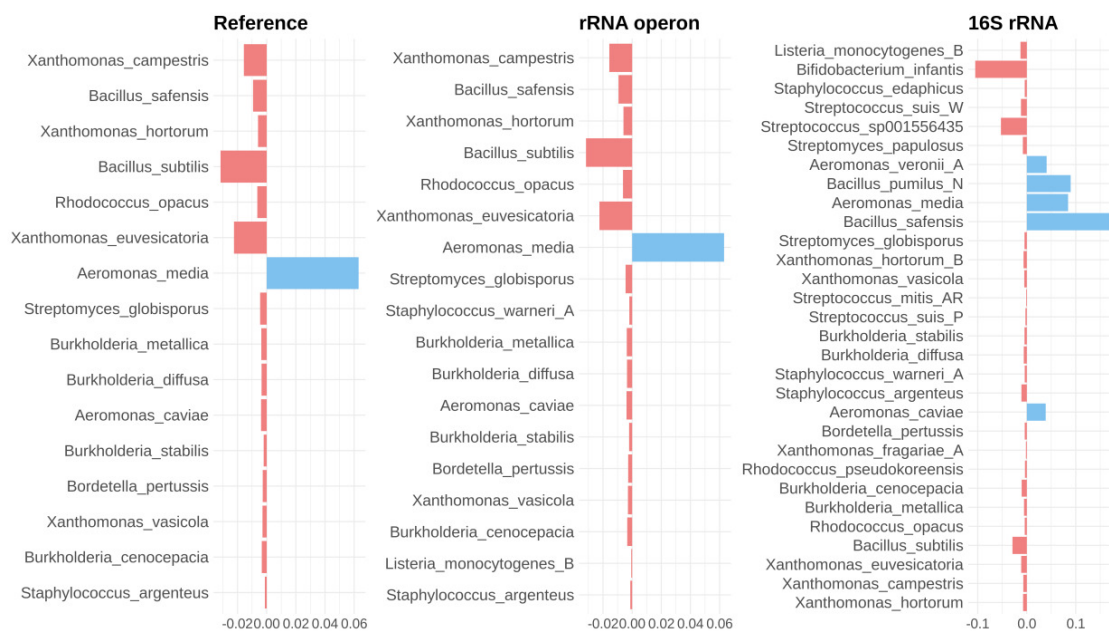


TABLE 1. Pearson correlations between actual species proportions and predicted species proportions using the rRNA operon, 16S rRNA, and 16S rRNA V3-V4 region, by the number of reads used for the simulation.

	5,000	10,000	50,000	100,000	500,000	1,000,000
rRNA operon	0.998	0.998	0.998	0.998	0.998	0.998
16S rRNA	0.804	0.821	0.835	0.837	0.838	0.838
16S V3-V4	0.295	0.321	0.344	0.346	0.348	0.348

TABLE 2. Error rates applied to simulate sequencing inaccuracies in rRNA operon and 16S rRNA sequences for species classification simulations. The table outlines the mismatch, insertion, and deletion rates for the rRNA operon sequenced with Nanopore technology and the 16S rRNA (including its V3-V4 regions) sequenced with Illumina.

Marker region	Sequencing tool	Mismatch rate	Insertion rate	Deletion rate
rRNA operon	Nanopore	0.0116	0.0081	0.0144
16S rRNA	Illumina	0.0089	0.00045	0.00045
16S rRNA V3-V4	Illumina	0.0089	0.00045	0.00045