

1 **Amplicon sequence variants artificially split bacterial genomes into**  
2 **separate clusters**

3 **Running title:** ASVs artificially split bacterial genomes

4 Patrick D. Schloss<sup>†</sup>

5 <sup>†</sup> To whom correspondence should be addressed:

6 pschloss@umich.edu

7 Department of Microbiology & Immunology

8 University of Michigan

9 Ann Arbor, MI 48109

10 **Observation Format**

## 11 **Abstract**

12 Amplicon sequencing variants (ASVs) have been proposed as an alternative to operational taxonomic units  
13 (OTUs) for analyzing microbial communities. ASVs have grown in popularity, in part, because of a desire to  
14 reflect a more refined level of taxonomy since they do not cluster sequences based on a distance-based  
15 threshold. However, ASVs and the use of overly narrow thresholds to identify OTUs increase the risk of  
16 splitting a single genome into separate clusters. To assess this risk, I analyzed the intragenomic variation of  
17 16S rRNA genes from the bacterial genomes represented in a *rrn* copy number database, which contained  
18 20,427 genomes from 5,972 species. As the number of copies of the 16S rRNA gene increased in a  
19 genome, the number of ASVs also increased. There was an average of 0.58 ASVs per copy of the 16S  
20 rRNA gene for full length 16S rRNA genes. It was necessary to use a distance threshold of 5.25% to cluster  
21 full length ASVs from the same genome into a single OTU with 95% confidence for genomes with 7 copies of  
22 the 16S rRNA, such as *E. coli*. This research highlights the risk of splitting a single bacterial genome into  
23 separate clusters when ASVs are used to analyze 16S rRNA gene sequence data. Although there is also a  
24 risk of clustering ASVs from different species into the same OTU when using broad distance thresholds,  
25 those risks are of less concern than artificially splitting a genome into separate ASVs and OTUs.

## 26 **Importance**

27 16S rRNA gene sequencing has engendered significant interest in studying microbial communities. There  
28 has been a tension between trying to classify 16S rRNA gene sequences to increasingly lower taxonomic  
29 levels and the reality that those levels were defined using more sequence and physiological information than  
30 is available from a fragment of the 16S rRNA gene. Furthermore, naming of bacterial taxa reflects the biases  
31 of those who name them. One motivation for the recent push to adopt ASVs in place of OTUs in microbial  
32 community analyses is to allow researchers to perform their analyses at the finest possible level that reflects  
33 species-level taxonomy. The current research is significant because it quantifies the risk of artificially splitting  
34 bacterial genomes into separate clusters. Far from providing a better representation of bacterial taxonomy  
35 and biology, the ASV approach can lead to conflicting inferences about the ecology of different ASVs from  
36 the same genome.

37 16S rRNA gene sequencing is a powerful technique for describing and comparing microbial communities (1).  
38 Efforts to link 16S rRNA gene sequences to taxonomic levels based on distance thresholds date to at least  
39 the 1990s. The distance-based threshold that was developed and is now widely used was based on  
40 DNA-DNA hybridization approaches that are not as precise as genome sequencing (2, 3). Instead, genome  
41 sequencing technologies have suggested that the widely used 3% distance threshold to operationally define  
42 bacterial taxa is too coarse (4–6). As an alternative to operational taxonomic units (OTUs), amplicon  
43 sequencing variants (ASVs) have been proposed as a way to adopt the thresholds suggested by genome  
44 sequencing to microbial community analysis using 16S rRNA gene sequences (7–10). Approaches for  
45 identifying ASVs do not cluster sequences based on a distance-based threshold (11). Proponents of ASVs  
46 are largely dismissive of concerns that most bacterial genomes have more than one copy of the *rrn* operon  
47 and that those copies are not identical (12, 13). Yet, ASVs and using too fine a threshold to identify OTUs  
48 could split a single genome into multiple clusters. Conversely, using too broad of a threshold to define OTUs  
49 could cluster together multiple bacterial species into the same OTU. An example of both is seen in the  
50 comparison of *Staphylococcus aureus* (NCTC 8325) and *S. epidermidis* (ATCC 12228) where each genome  
51 has 5 copies of the 16S rRNA gene. Each of the 10 copies of the 16S rRNA gene in these two genomes is  
52 distinct and represent 10 ASVs. Conversely, if the copies were clustered using a 3% distance threshold, then  
53 all 10 ASVs would cluster into the same OTU. The goal of this study was to quantify the tradeoff of splitting a  
54 single genome into multiple clusters and the risk of clustering different bacterial species into the same  
55 cluster when using ASVs and various OTU definitions.

56 To investigate the variation in the number of copies of the 16S rRNA gene per genome and the intragenomic  
57 variation among copies of the 16S rRNA gene, I obtained 16S rRNA sequences from the *rrn* copy number  
58 database (*rrnDB*)(14). Among the 5,972 species represented in the *rrnDB* there were 20,427 genomes. The  
59 median *rrn* copy number per species ranged between 1 (e.g., *Mycobacterium tuberculosis*) and 19  
60 (*Metabacillus litoralis*). As the *rrn* copy number for a genome increased, the number of variants of the 16S  
61 rRNA gene in each genome also increased. On average, there were 0.58 variants per copy of the full length  
62 16S rRNA gene and an average of 0.32, 0.25, and 0.27 variants when considering the V3-V4, V4, and  
63 V4-V5 regions of the gene, respectively. Although a species tended to have a consistent number of 16S  
64 rRNA gene copies per genome, the number of total variants increased with the number of genomes that  
65 were sampled (Figure S1). For example, the 271 genome accessions of *Mycobacterium tuberculosis* in the  
66 *rrnDB* each had 1 copy of the gene per genome. However, across those accessions, there were 17 versions  
67 of the gene. An *E. coli* genome typically had 7 copies of the 16S rRNA gene with a median of 5 distinct full  
68 length ASVs per genome (intraquartile range between 3 and 6). Across the 1,390 *E. coli* genomes in the

69 *rrnDB*, there were 1,402 versions of the gene. These observations highlight the risk of selecting a threshold  
70 for defining clusters that is too narrow because it is possible to split a single genome into multiple clusters.

71 A method to avoid splitting a single genome into multiple clusters is to cluster 16S rRNA gene sequences  
72 together based on their distances between each other. Therefore, I assessed the impact of the distance  
73 threshold used to define clusters of 16S rRNA genes on the propensity to split a genome into separate  
74 clusters. To control for uneven representation of genomes across species, I randomly selected one genome  
75 from each species and repeated each randomization 100 times. I observed that as the *rrn* copy number  
76 increased, the distance threshold required to reduce the ASVs in each genome to a single OTU increased  
77 (Figure 1). Among species with 7 copies of the *rrn* operon (e.g., *E. coli*), a distance threshold of 5.25% was  
78 required to reduce full length ASVs into a single OTU for 95% of the species. Similarly, thresholds of 5.25,  
79 2.50, and 3.75% were required for the V3-V4, V4, and V4-V5 regions, respectively. But, if a 3% distance  
80 threshold was used, then ASVs from genomes containing fewer than 6, 6, 8, and 6 copies of the *rrn* operon  
81 would reliably be clustered into a single OTU for ASVs from the V1-V9, V3-V4, V4, and V4-V5 regions,  
82 respectively. Consequently, these results demonstrate that broad thresholds must be used to avoid splitting  
83 different operons from the same genome into separate clusters.

84 At broad thresholds, 16S rRNA gene sequences from multiple species could be clustered into the same ASV  
85 or OTU. I again randomly selected one genome from each species to control for uneven representation of  
86 genomes across species and for this analysis I measured the percentage of ASVs and OTUs that contained  
87 16S rRNA gene sequences from multiple species (Figure 2). Without using distance-based thresholds, 4.1%  
88 of the ASVs contained sequences from multiple species when considering full length sequences and 10.9,  
89 16.2, and 13.1% when considering the V3-V4, V4, and V4-V5 regions, respectively. At the commonly used  
90 3% threshold for defining OTUs, 27.4% of the OTUs contained 16S rRNA gene sequences from multiple  
91 species when considering full length sequences and 31.7, 34.3, and 34.8% when considering the V3-V4, V4,  
92 and V4-V5 regions, respectively. Considering that species designations are inconsistently applied and reflect  
93 multiple human-imposed biases, the risk of splitting a genome into multiple OTUs is more problematic than  
94 clustering species together. Therefore, larger thresholds are advisable.

95 The results of this analysis demonstrate that there is a significant risk of splitting a single genome into  
96 multiple clusters if using ASVs or too fine of a threshold to define OTUs. An ongoing problem for  
97 amplicon-based studies is defining a meaningful taxonomic unit (11, 15, 16). Since there is no consensus for  
98 a biological definition of a bacterial species (17–19), microbiologists must accept that how bacterial species  
99 are named is biased and that taxonomic rules are not applied in a consistent manner (e.g., (19, 20)). This  
100 makes it impossible to fit a distance threshold to define an OTU definition that matches a set of species

101 names (21). Furthermore, the 16S rRNA gene does not evolve at the same rate across all bacterial lineages  
102 (15), which limits the biological interpretation of a common OTU definition. A distance-based definition of a  
103 taxonomic unit based on 16S rRNA gene or full genome sequences is, at best, operational and not grounded  
104 in biological theory (15, 22–24). There is general agreement in bacterial systematics that to classify an  
105 organism to a bacterial species, phenotypic and genome sequence data are needed (17–20). A short  
106 sequence from a bacterial genome simply cannot differentiate between species. Moreover, it is difficult to  
107 defend a clustering threshold that would split a single genome into multiple taxonomic units. It is not  
108 biologically plausible to entertain the possibility that different *rrn* operons from the same genome would have  
109 different ecologies. Although there are multiple reasons that proponents favor ASVs, the significant risk of  
110 artificially splitting genomes into separate clusters is too high to warrant their use.

111 **Materials and Methods. (i) Data availability.** The 16S rRNA gene sequences used in this study were  
112 obtained from the *rrnDB* (<https://rrndb.umms.med.umich.edu>; version 5.7, released January 18, 2021) (14).  
113 At the time of submission, this was the most current version of the database. The *rrnDB* obtained the  
114 curated 16S rRNA gene sequences from the KEGG database, which ultimately obtained them from NCBI's  
115 non-redundant RefSeq database. The *rrnDB* provided downloadable versions of the sequences with their  
116 taxonomy as determined using the naive Bayesian classifier trained on the RDP reference taxonomy. For  
117 some genomes this resulted in multiple classifications since a genome's 16S rRNA gene sequences were  
118 not identical. Instead, I mapped the RefSeq accession number for each genome in the database to obtain a  
119 single taxonomy for each genome. Because strain names were not consistently given to genomes across  
120 bacterial species, I disregarded the strain level designations.

121 **(ii) Definition of regions within the 16S rRNA gene.** The full length 16S rRNA gene sequences were  
122 aligned to a SILVA reference alignment of the 16S rRNA gene (v. 138) using the mothur software package (v.  
123 1.44.2) (25, 26). Regions of the 16S rRNA gene were selected because of their use in the microbial ecology  
124 literature. Full length sequences corresponded to *E. coli* str. K-12 substr. MG1655 (NC\_000913) positions  
125 28 through 1491, V4 to positions 534 through 786, V3-V4 to positions 358 through 786, and V4-V5 to  
126 positions 534 through 908. The positions between these coordinates reflect the fragments that would be  
127 amplified using commonly used PCR primers.

128 **(iii) Clustering sequences into OTUs.** Pairwise distances between sequences were calculated using the  
129 `dist.seqs` command from mothur. The OptiClust algorithm, as implemented in mothur, was used to assign  
130 16S rRNA gene sequences to OTUs (27). Distance thresholds between 0.25 and 10.00% in 0.25 percentage

131 point increments were used to assign sequences to OTUs.

132 **(iv) Controlling for uneven sampling of genomes by species.** Because of the uneven distribution of  
133 genome sequences across species I randomly selected one genome from each species for the analysis of  
134 splitting genomes and clustering ASVs from different species (Figures 1 and 2). The random selection was  
135 repeated 100 times. Analyses based on this randomization reported the median of the 100 randomizations.  
136 The intraquartile range between randomizations was less than 0.0024. Because the range was so small, the  
137 confidence intervals were more narrow than the thickness of the lines in Figures 1 and 2 and were not  
138 included.

139 **(v) Reproducible data analysis.** The code to perform the analysis in this manuscript and its history are  
140 available as a git-based version control repository on GitHub  
141 ([https://github.com/SchlossLab/Schloss\\_rrnAnalysis\\_mSphere\\_2021](https://github.com/SchlossLab/Schloss_rrnAnalysis_mSphere_2021)). The analysis can be regenerated  
142 using a GNU Make-based workflow that made use of built-in bash tools (v. 3.2.57), mothur (v. 1.44.2), and R  
143 (v. 4.0.4). Within R, I used the tidyverse (v. 1.3.0), data.table (v. 1.13.2), Rcpp (v. 1.0.5), furrr (v. 0.2.1), here  
144 (v. 1.0.1) and rmarkdown (v. 2.5) packages. The conception and development of this analysis is available as  
145 a playlist on the Riffomonas YouTube channel  
146 ([https://youtube.com/playlist?list=PLmNrK\\_nkqBpL7m\\_tyWdQgdyurerttCsPY](https://youtube.com/playlist?list=PLmNrK_nkqBpL7m_tyWdQgdyurerttCsPY)).

147 **(vi) Note on usage of ASV, OTU, and cluster.** I used “ASV” to denote the cluster of true 16S rRNA gene  
148 sequences that were identical to each other and “OTU” to denote the product of distance-based clustering of  
149 sequences. Although ASVs do represent a type of operational definition of a taxonomic unit and can be  
150 thought of as an OTU formed using a distance of zero, proponents of the ASV approach prefer to avoid the  
151 term OTU given the long history of OTUs being formed by distance-based clustering  
152 (<https://github.com/benjjneb/dada2/issues/62>; accessed 2021-02-26). For this reason, when an ASV split a  
153 genome into different units, those units were called clusters rather than OTUs.

154 **Acknowledgements.** I am grateful to Robert Hein and Thomas Schmidt, who maintain the *rrnDB*, for their  
155 help in understanding the curation of the database and for making the 16S rRNA gene sequences and  
156 related metadata publicly available. I am also grateful to community members who watched the serialized  
157 version of this analysis on YouTube and provided suggestions and questions over the course of the  
158 development of this project. This work was supported, in part, through grants from the NIH (P30DK034933,  
159 U01AI124255, and R01CA215574).

## 160 References

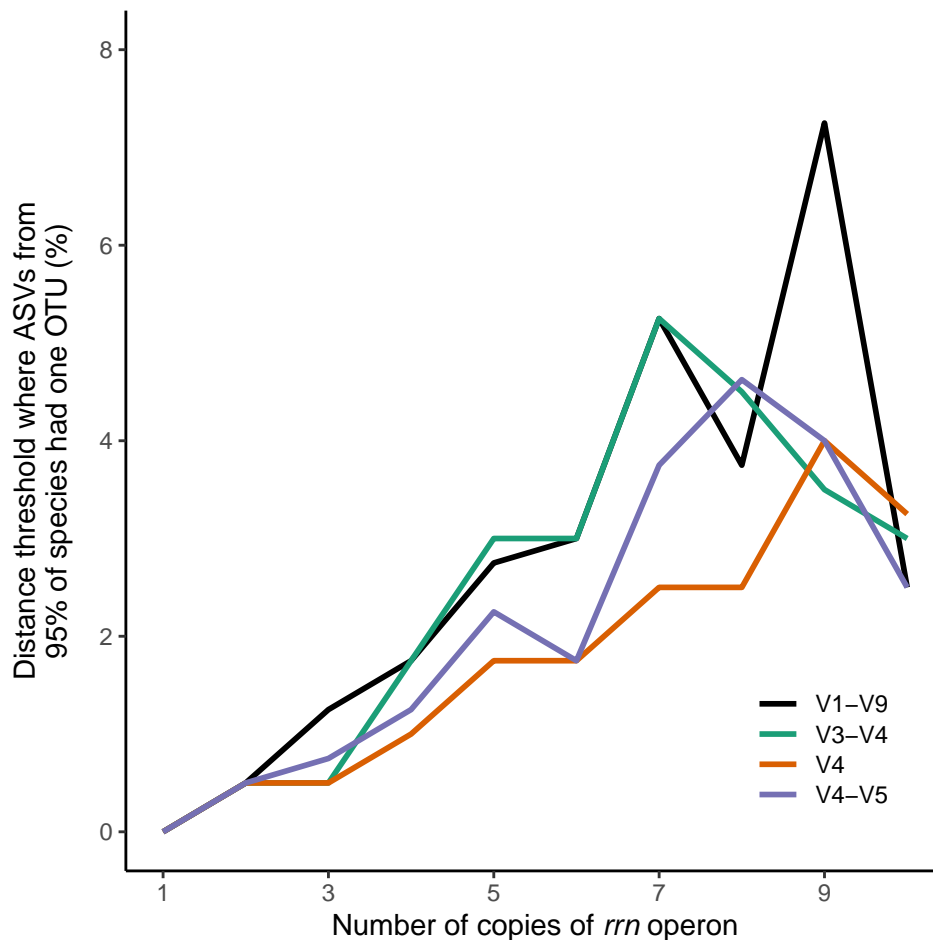
- 161 1. Lane DJ, Pace B, Olsen GJ, Stahl DA, Sogin ML, Pace NR. 1985. Rapid determination of 16S ribosomal  
162 RNA sequences for phylogenetic analyses. *Proceedings of the National Academy of Sciences*  
163 82:6955–6959. doi: 10.1073/pnas.82.20.6955.
- 164 2. Stackebrandt E, Goebel BM. 1994. Taxonomic note: A place for DNA-DNA reassociation and 16S rRNA  
165 sequence analysis in the present species definition in bacteriology. *International Journal of Systematic  
166 and Evolutionary Microbiology* 44:846–849. doi: 10.1099/00207713-44-4-846.
- 167 3. Goris J, Konstantinidis KT, Klappenbach JA, Coenye T, Vandamme P, Tiedje JM. 2007. DNA-DNA  
168 hybridization values and their relationship to whole-genome sequence similarities. *International Journal  
169 of Systematic and Evolutionary Microbiology* 57:81–91. doi: 10.1099/ijs.0.64483-0.
- 170 4. Rodriguez-R LM, Castro JC, Kyrpides NC, Cole JR, Tiedje JM, Konstantinidis KT. 2018. How much do  
171 rRNA gene surveys underestimate extant bacterial diversity? *Applied and Environmental Microbiology*  
172 84:e00014–18. doi: 10.1128/aem.00014-18.
- 173 5. Stackebrandt E, Ebers J. 2006. Taxonomic parameters revisited: Tarnished gold standards. *Microbiol  
174 Today* 33:152–155.
- 175 6. Edgar RC. 2018. Updating the 97% identity threshold for 16S ribosomal RNA OTUs. *Bioinformatics*  
176 34:2371–2375. doi: 10.1093/bioinformatics/bty113.
- 177 7. Edgar RC. 2016. UNOISE2: Improved error-correction for illumina 16S and its amplicon sequencing.  
178 *bioRxiv*. doi: 10.1101/081257.
- 179 8. Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Zech Xu Z, Kightley EP, Thompson LR,  
180 Hyde ER, Gonzalez A, Knight R. 2017. Deblur rapidly resolves single-nucleotide community sequence  
181 patterns. *mSystems* 2:e00191–16. doi: 10.1128/mSystems.00191-16.
- 182 9. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. 2016. DADA2: High-resolution  
183 sample inference from illumina amplicon data. *Nature Methods* 13:581–583. doi: 10.1038/nmeth.3869.
- 184 10. Eren AM, Morrison HG, Lescault PJ, Reveillaud J, Vineis JH, Sogin ML. 2014. Minimum entropy  
185 decomposition: Unsupervised oligotyping for sensitive partitioning of high-throughput marker gene  
186 sequences. *The ISME Journal* 9:968–979. doi: 10.1038/ismej.2014.195.
- 187 11. Callahan BJ, McMurdie PJ, Holmes SP. 2017. Exact sequence variants should replace operational  
188 taxonomic units in marker-gene data analysis. *The ISME Journal* 11:2639–2643. doi:



- 189 10.1038/ismej.2017.119.
- 190 12. Pei AY, Oberdorf WE, Nossa CW, Agarwal A, Chokshi P, Gerz EA, Jin Z, Lee P, Yang L, Poles M, Brown  
191 SM, Sotero S, DeSantis T, Brodie E, Nelson K, Pei Z. 2010. Diversity of 16S rRNA genes within  
192 individual prokaryotic genomes. *Applied and Environmental Microbiology* 76:3886–3897. doi:  
193 10.1128/aem.02953-09.
- 194 13. Sun D-L, Jiang X, Wu QL, Zhou N-Y. 2013. Intragenomic heterogeneity of 16S rRNA genes causes  
195 overestimation of prokaryotic diversity. *Applied and Environmental Microbiology* 79:5962–5969. doi:  
196 10.1128/aem.01282-13.
- 197 14. Stoddard SF, Smith BJ, Hein R, Roller BRK, Schmidt TM. 2014. *rrnDB*: Improved tools for interpreting  
198 rRNA gene abundance in bacteria and archaea and a new foundation for future development. *Nucleic  
199 Acids Research* 43:D593–D598. doi: 10.1093/nar/gku1201.
- 200 15. Schloss PD, Westcott SL. 2011. Assessing and improving methods used in operational taxonomic  
201 unit-based approaches for 16S rRNA gene sequence analysis. *Applied and Environmental Microbiology*  
202 77:3219–3226. doi: 10.1128/aem.02810-10.
- 203 16. Johnson JS, Spakowicz DJ, Hong B-Y, Petersen LM, Demkowicz P, Chen L, Leopold SR, Hanson BM,  
204 Agresta HO, Gerstein M, Sodergren E, Weinstock GM. 2019. Evaluation of 16S rRNA gene sequencing  
205 for species and strain-level microbiome analysis. *Nature Communications* 10:5029. doi:  
206 10.1038/s41467-019-13036-1.
- 207 17. Staley JT. 2006. The bacterial species dilemma and the genomicphylogenetic species concept.  
208 *Philosophical Transactions of the Royal Society B: Biological Sciences* 361:1899–1909. doi:  
209 10.1098/rstb.2006.1914.
- 210 18. Oren A, Garrity GM. 2013. Then and now: A systematic review of the systematics of prokaryotes in the  
211 last 80 years. *Antonie van Leeuwenhoek* 106:43–56. doi: 10.1007/s10482-013-0084-1.
- 212 19. Sanford RA, Lloyd KG, Konstantinidis KT, Löffler FE. 2021. Microbial taxonomy run amok. *Trends in  
213 Microbiology*. doi: 10.1016/j.tim.2020.12.010.
- 214 20. Baltrus DA, McCann HC, Guttman DS. 2016. Evolution, genomics and epidemiology of *Pseudomonas  
215 syringae*. *Molecular Plant Pathology* 18:152–168. doi: 10.1111/mpp.12506.
- 216 21. Konstantinidis KT, Tiedje JM. 2005. Towards a genome-based taxonomy for prokaryotes. *Journal of  
217 Bacteriology* 187:6258–6264. doi: 10.1128/jb.187.18.6258-6264.2005.

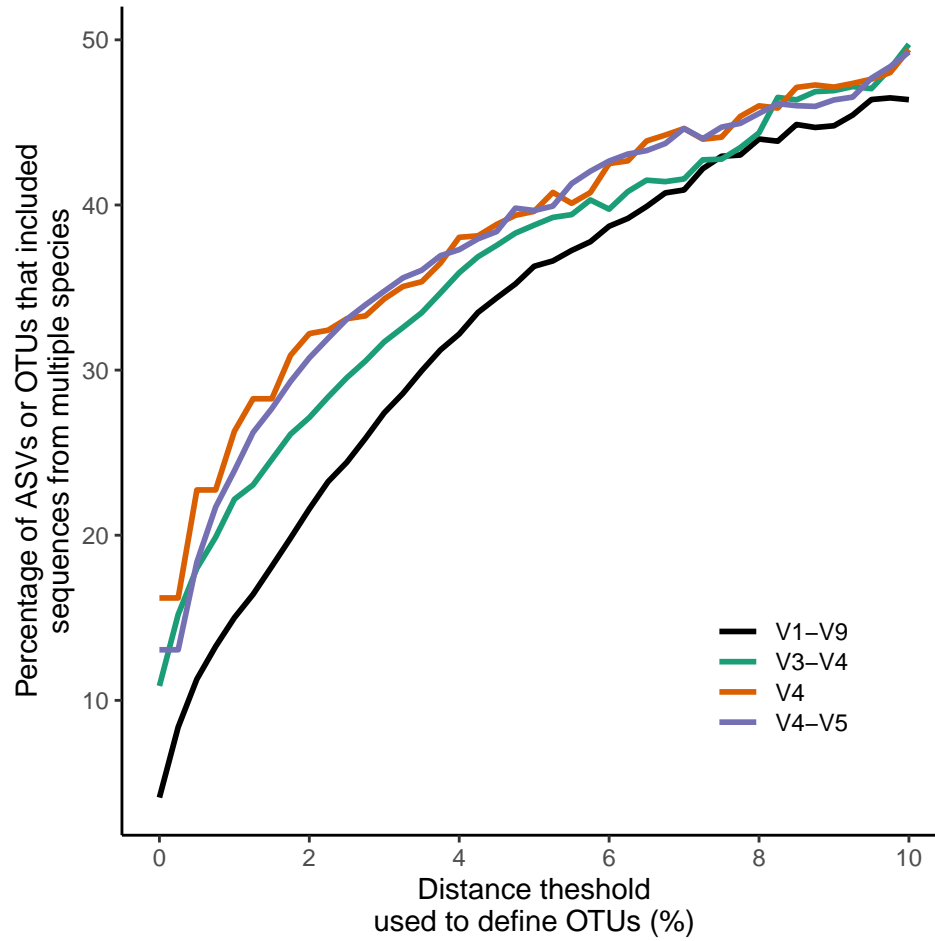


- 218 22. Barco RA, Garrity GM, Scott JJ, Amend JP, Neelson KH, Emerson D. 2020. A genus definition for  
219 bacteria and archaea based on a standard genome relatedness index. *mBio* 11:02475–19. doi:  
220 10.1128/mbio.02475-19.
- 221 23. Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, Hugenholtz P. 2020. A complete  
222 domain-to-species taxonomy for bacteria and archaea. *Nature Biotechnology* 38:1079–1086. doi:  
223 10.1038/s41587-020-0501-8.
- 224 24. Yarza P, Yilmaz P, Pruesse E, Glöckner FO, Ludwig W, Schleifer K-H, Whitman WB, Euzéby J, Amann R,  
225 Rosselló-Móra R. 2014. Uniting the classification of cultured and uncultured bacteria and archaea using  
226 16S rRNA gene sequences. *Nature Reviews Microbiology* 12:635–645. doi: 10.1038/nrmicro3330.
- 227 25. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB,  
228 Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Horn DJV, Weber CF. 2009. Introducing  
229 mothur: Open-source, platform-independent, community-supported software for describing and  
230 comparing microbial communities. *Applied and Environmental Microbiology* 75:7537–7541. doi:  
231 10.1128/aem.01541-09.
- 232 26. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO. 2012. The SILVA  
233 ribosomal RNA gene database project: Improved data processing and web-based tools. *Nucleic Acids  
234 Research* 41:D590–D596. doi: 10.1093/nar/gks1219.
- 235 27. Westcott SL, Schloss PD. 2017. OptiClust, an improved method for assigning amplicon-based sequence  
236 data to operational taxonomic units. *mSphere* 2:00073–17. doi: 10.1128/mspheredirect.00073-17.



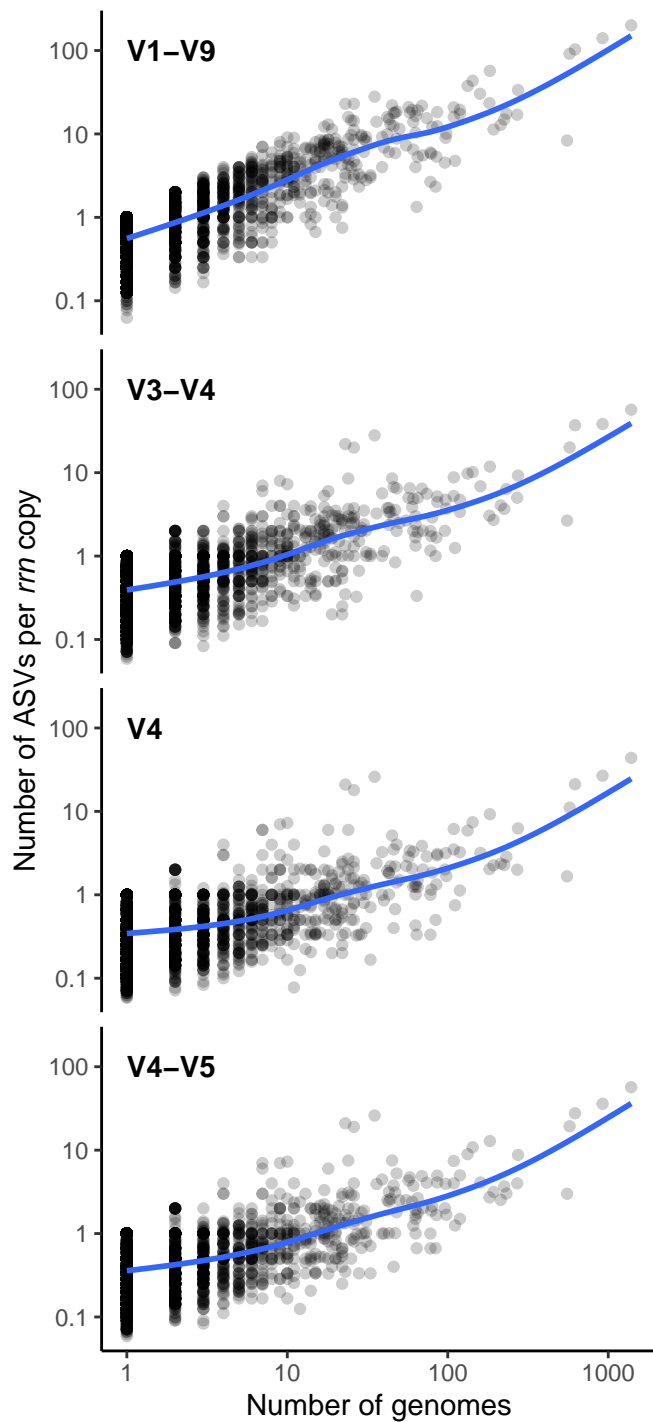
237

238 **Figure 1. The distance threshold required to prevent the splitting of genomes into multiple OTUs**  
239 **increased as the number of *rrn* operons in the genome increased.** Each line represents the median  
240 distance threshold for each region of the 16S rRNA gene that is required for 95% of the genomes with the  
241 indicated number of *rrn* operons to cluster their ASVs to a single OTU. The median distance threshold was  
242 calculated across 100 randomizations in which one genome was sampled from each species. Only those  
243 number of *rrn* operons that were found in more than 100 species are included.



244

245 **Figure 2. As the distance threshold used to define an OTU increased, the percentage of ASVs and**  
246 **OTUs representing multiple species increased.** These data represent the median fractions for both  
247 measurements across 100 randomizations. In each randomization, one genome was sampled from each  
248 species.



249

250 **Figure S1. The ratio of number of distinct ASVs per copy of the *rrn* operon increased for a species as**  
251 **the number of genomes in the *rrn*DB for that species increased.** Each point represents a different  
252 species and was shaded to be 80% transparent so that when points overlap they become darker. The blue  
253 line represents a smoothed fit through the data. Both axes use a logarithmic scale (base 10).