

Gene copy normalization of the 16S rRNA gene cannot outweigh the methodological biases of sequencing

Robert Starke, Daniel Morais

1 Laboratory of Environmental Microbiology, Institute of Microbiology of the CAS, Vídeňská 1083,
2 14220 Praha 4, Czech Republic

3 *The 16S rRNA gene is the golden standard target of sequencing to uncover the composition of*
4 *bacterial communities but the presence of multiple copies of the gene makes gene copy*
5 *normalization (GCN) inevitable. Even though GCN resulted in abundances closer to the*
6 *metagenome, it should be validated by communities with known composition as both amplicon*
7 *and shotgun sequencing are prone to methodological biases. Here we compared the*
8 *composition of three mock communities to the composition derived from 16S sequencing*
9 *without and with GCN. In all of them, the 16S composition was different from the mock*
10 *community and GCN improved the picture only in the community with the lowest Shannon*
11 *diversity. Albeit with low abundance, half of the identified genera were not present in the mock*
12 *communities. Our approach provides empirical evidence to the methodological biases*
13 *introduced by sequencing that was only counteracted by GCN in the case of low α -diversity,*
14 *potentially due to the small number of bacterial taxa with known gene copy numbers. We thus*
15 *cannot recommend the use of GCN moving forward and it is questionable whether a complete*
16 *catalogue of 16S rRNA copy numbers can outweigh the methodological biases of sequencing.*

17 Amplicon sequencing of 16S rRNA is the golden standard to describe the composition of
18 bacterial communities due to (i) cost, (ii) availability, (iii) presence of extraction and preparation
19 kits, (iv) taxonomic resolution as deep as the level of genera and (v) previous research.
20 Unsurprisingly it outcompeted (46,473 papers as of February 2019) other possible techniques
21 to describe community structures such as metagenomics (7,699), metaproteomics (367) or
22 metatranscriptomics (439). The general practice as shown by the myriads of publications does
23 not comprise the correction of the obtained raw counts by 16S rRNA gene copy numbers per
24 bacterial genome even though it is known that bacteria can have multiple copy numbers of the
25 gene ¹. Logically, two bacteria with similar raw counts but different gene copy numbers cannot

26 be equally abundant which is why GCN seems necessary. The recent recommendation against
27 GCN based on the systematic evaluation of the predictability of 16S GCNs in bacteria ²
28 contradict the previous suggestion in favor of GCN based on the comparison of 16S and
29 metagenomics ¹. However, sequencing techniques are prone to similar methodological biases
30 introduced by extraction, PCR, sequencing and bioinformatics, and could thus similarly diverge
31 from the real picture as recently demonstrated ³. We therefore believe that the use of mock
32 communities as standard is inevitable to prove the viability of GCN. Here we compared three
33 randomly chosen taxonomically different mock communities (Mock-2, Mock-20 and Mock-21)
34 from *mockrobiata* provided elsewhere ⁴ that derived from the combination of extracted
35 genomic DNA from bacterial strains and the subsequent 16S rRNA gene amplicon sequencing to
36 estimate the impact of GCN on the bacterial community composition.

37 Operational taxonomic units (OTU) were annotated by *blastn* ⁵ as best hit down to the
38 genus level with an average similarity match of 97.7±1.7% for Mock-2, 97.4±1.7% for Mock-20
39 and 97.6±1.6% for Mock-21. In total 3,973 from 34,154 OTU counts (13.2%) in Mock-2 could not
40 be assigned to a bacterial genus compared to 378 from 173,460 (0.22%) for Mock-20 and 328
41 from 180,542 (0.18%) for Mock-21. Mock-2 comprised of 23 bacterial genera of which only 14
42 were identified by 16S sequencing opposed to 17 in Mock-20 and 16 in Mock-21 of which all
43 have been identified (**Figure 1**). These findings illustrated missed identifications that seem to be
44 related to the sequencing depth. 30,000 OTU counts for Mock-2 were not sufficient to identify
45 all of the 23 genera in the community whereas 180,000 OTU counts for Mock-20 and Mock-21
46 resulted in the identification of all genera. However, the three mock communities are simple
47 compared to the billions of organisms belonging to thousands of different species in one gram
48 of soil ⁶. Particularly considering the prokaryotic density of 10,000,000 organisms per gram of
49 soil ⁷ that is at least one magnitude of order higher than per milliliter of water in the ocean ⁸,
50 we conclude that 10,000 but at least 2,000 OTU counts per taxonomic rank of interest are
51 necessary to fully cover the members of the community.

52 In total, 19 genera in Mock-2 together with 18 in Mock-20 and 77 in Mock-22 were
53 wrongly identified during sequencing due to their absence when the extracted genomic DNA
54 was combined. Admittedly, the majority among them were found with low abundance,

55 presumably as noise during sequencing. However, *Klebsiella* of the family *Enterobacteriaceae*
56 comprised high abundances in each community. Together with highly abundant extracted DNA
57 from *Escherichia* but low sequencing abundances, we conclude the misidentification of
58 *Escherichia*, also an *Enterobacteriaceae*, as *Klebsiella*. In compliance with our results,
59 phylogenetic trees based on the 16S rRNA gene are ambiguous in *Enterobacteriaceae* and differ
60 in the relative position of several genera^{9,10}. Processes of recombination and gene conversion
61^{11,12}, and different sequences of the 16S rRNA gene found within a single species¹³ were
62 previously hold accountable. Here we provide empirical evidence for the misidentification of
63 *Escherichia* as *Klebsiella*, which could prove to be detrimental for proper prophylactic medical
64 treatment since both are pathogens causing a different array of diseases^{14,15}. Even though one
65 advantage of targeting the 16S rRNA gene is taxonomic resolution, we report the failure of
66 correct identification of the genus within *Enterobacteriaceae*, which could be true for other
67 bacterial families as well.

68 Non-metric multidimensional scaling (NMDS) revealed that the mock community
69 composition derived from combining extracted DNA was indeed different to the composition
70 derived from 16S rRNA gene amplicon sequencing (**Figure 2**). GCN did not impact the
71 community composition and its relative distance within two dimension of the NMDS to Mock-2
72 and Mock-20, which was further supported by the residual sum of squares between each mock
73 community and 16S sequencing without and with GCN (**Table 1**). However, in Mock-21, GCN
74 resulted in a picture closer to the mock community, presumably due to the low complexity of
75 the community as its Shannon diversity on genus level was 40% lower (1.69) than in Mock-2
76 (2.71) and in Mock-21 (2.76). Logically we suggest that GCN in communities of low Shannon
77 diversity (<<2.7) could be beneficial. However, the Shannon diversity of bacterial communities
78 typically range between 3 and 6 in both terrestrial¹⁶⁻¹⁹ and aquatic ecosystems^{20,21}, which are
79 likely too diverse for GCN to have an impact as shown for Mock-2 and Mock-20.

80 Concluding, together with the issues to predict 16S GCNs in bacteria², we cannot
81 recommend the use of GCN based on the *in vitro* comparison of sequenced amplicons from
82 three randomly chosen mock communities.

83 **Methods**

84 *Data generation*

85 The community data was obtained from the *mockrobiota* database provided by Bokulich and
86 colleagues⁴. Three mock communities that contain the reverse reads of sequencing or a clear
87 summary of the known community composition were randomly chosen: Mock-2 that has been
88 described elsewhere^{22,23} together with Mock-20 and Mock-21 that was obtained through BEI
89 Resources, NIAID, NIH as part of the Human Microbiome Project: Genomic DNA from Microbial
90 Mock Community B (Even, High Concentration), v5.1H, for Whole Genome Shotgun Sequencing,
91 HM-276D. The mock composition generated by the combination of extracted genomic DNA
92 from bacterial strains was downloaded from [https://github.com/caporaso-](https://github.com/caporaso-lab/mockrobiota/tree/master/data)
93 [lab/mockrobiota/tree/master/data](https://github.com/caporaso-lab/mockrobiota/tree/master/data) (Mock-2, Mock-20 and Mock-21) as expected taxonomy
94 using *SILVA* at a 99% identity criterion to remove highly identical sequences²⁴. The raw
95 sequencing data including the forward and reverse reads was processed with *SEED 2*²⁵. Briefly,
96 quality filtering at a cut-off of 30 was followed by clustering of the representative sequences
97 from the clusters as consensus and most abundant, and identification of OTUs by *blastn*⁵
98 against a 16S database including chloroplasts and archaea from the ribosomal database project
99 (RDP) as of December 2017²⁶.

100 *Gene copy number normalization and statistical analysis*

101 Known 16S rRNA gene copy numbers from bacterial genomes were obtained from the
102 Ribosomal RNA Database (*rrnDB*) as of September 2018²⁷. Of the 152 genera identified by 16S
103 sequencing, 116 were annotated with a gene copy number on genus level ranging from one to
104 21 copies from one to 621 genomes per genus with an average gene copy number of 5.54 ± 0.99 .
105 For the remaining 36 genera, the next higher taxonomic rank with a gene copy number derived
106 from it was used. For each OTU, The raw counts were divided by the mean gene copy number
107 of the annotated genus to obtain the absolute normalized OTU content. The absolute OTU
108 counts were divided by the total OTU counts to give relative abundances. Non-metric
109 multidimensional scaling (NMDS) using Bray Curtis distances in two dimensions of the known
110 community structure (Mock) and derived from 16S sequencing without (16S) and with GCN

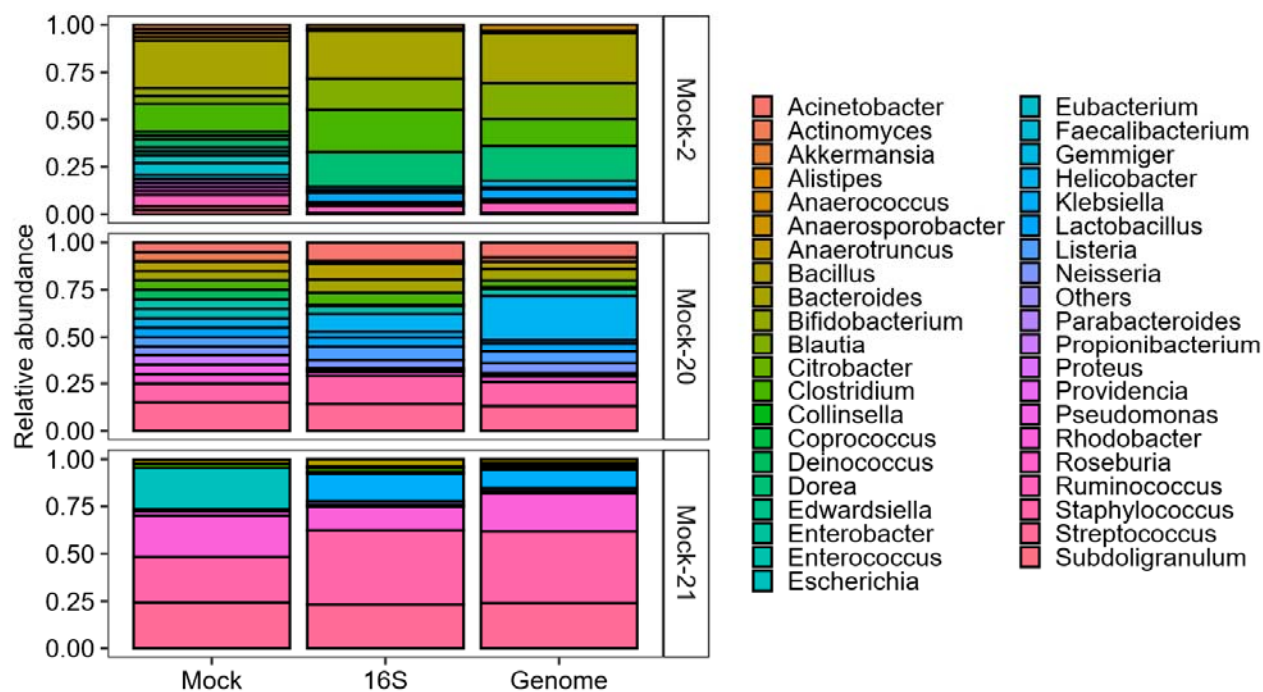
111 (Genome) was performed in *R* using the package *MASS* ²⁸. The distance (d) between the mock
112 community and the 16S sequencing without (16S) and with GCN (Genome) in two dimensions
113 was derived as straight line between two points (x_1, y_1) and (x_2, y_2) in a 2D-plane as given by the
114 Pythagorean Theorem (Equation 1). The residual sum of squares (RSS) was estimated from the
115 difference of the i^{th} value between the mock community as y_i and the 16S rRNA sequencing
116 without (16S) and with GCN (Genome) both as $f(x_i)$ given by Equation 2. The Shannon diversity
117 was calculated on the level of bacterial genera. Visualization was carried out in *R* using the
118 package *ggplot2* ²⁹.

119 (1)
$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

120 (2)
$$RSS = \sum_{i=1}^n (y_i - f(x_i))^2$$

121 **Figures**

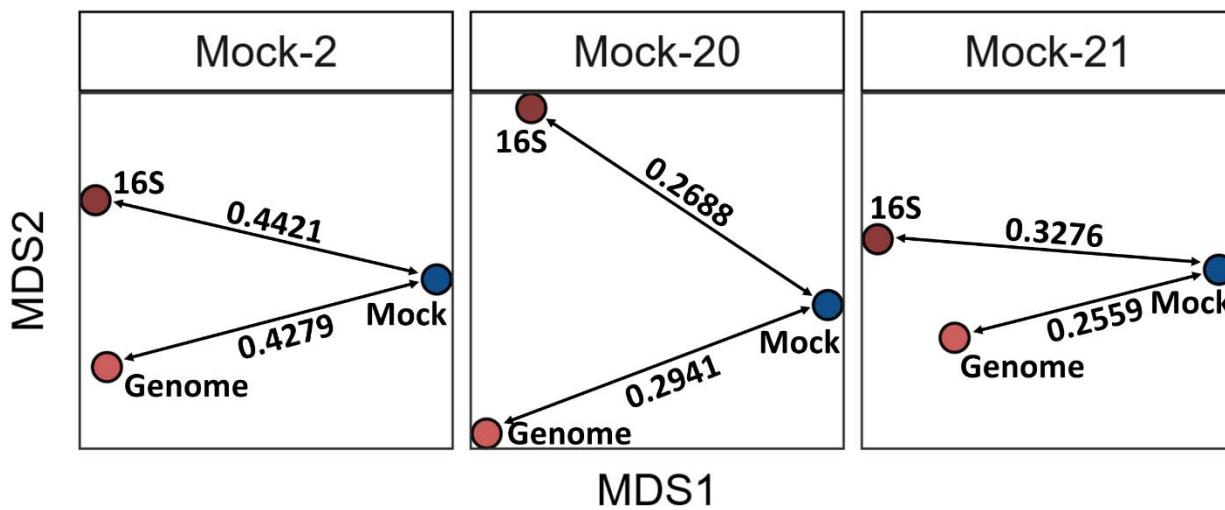
122 **Figure 1:** Community structure on the genus level of three mock communities (Mock) and
 123 estimated by 16S rRNA sequencing without (16S) and with GCN (Genome).



124

125

126 **Figure 2:** Non-metric multidimensional scaling (NMDS) using Bray Curtis distances in two
127 dimensions of the known community structure (in blue) and derived from 16S rRNA sequencing
128 without (16S) and with GCN (Genome).



129
130

131 **Tables**

132 **Table 1:** Residual sum of squares (RSS) as discrepancy between the known composition of the
133 mock community and the 16S rRNA sequencing without (16S) and with GCN (Genome) as well
134 as the α -diversity as Shannon diversity on genus level.

Community	16S	Genome	α-diversity
Mock-2	0.0576	0.0593	2.71
Mock-20	0.0204	0.0466	2.76
Mock-21	0.0968	0.0745	1.69

135

136 **References**

- 137 1. Větrovský, T. & Baldrian, P. The Variability of the 16S rRNA Gene in Bacterial Genomes
138 and Its Consequences for Bacterial Community Analyses. *PLoS One* (2013).
139 doi:10.1371/journal.pone.0057923
- 140 2. Louca, S., Doebeli, M. & Parfrey, L. W. Correcting for 16S rRNA gene copy numbers in
141 microbiome surveys remains an unsolved problem. *Microbiome* (2018).
142 doi:10.1186/s40168-018-0420-9
- 143 3. McLaren, M. R., Willis, A. D. & Callahan, B. J. Consistent and correctable bias in
144 metagenomic sequencing experiments. *bioRxiv* (2019).
145 doi:http://dx.doi.org/10.1101/559831doi
- 146 4. Bokulich, N. A. *et al.* mockrobiota: a Public Resource for Microbiome Bioinformatics
147 Benchmarking. *mSystems* (2016). doi:10.1128/mSystems.00062-16
- 148 5. BLAST. BLAST Basic Local Alignment Search Tool. *Blast Program Selection Guide* (2013).
149 doi:10.1006/jmbi.1990.9999
- 150 6. Torsvik, V. & Øvreås, L. Microbial diversity and function in soil: From genes to
151 ecosystems. *Curr. Opin. Microbiol.* (2002). doi:10.1016/S1369-5274(02)00324-7
- 152 7. Raynaud, X. & Nunan, N. Spatial ecology of bacteria at the microscale in soil. *PLoS One*
153 (2014). doi:10.1371/journal.pone.0087217
- 154 8. Whitman, W. B., Coleman, D. C. & Wiebe, W. J. Prokaryotes: The unseen majority. *Proc.*
155 *Natl. Acad. Sci.* (1998). doi:10.1073/pnas.95.12.6578
- 156 9. Granier, S. A. *et al.* Recognition of two genetic groups in the *Klebsiella oxytoca* taxon on
157 the basis of chromosomal β -lactamase and housekeeping gene sequences as well as
158 ERIC-1R PCR typing. *Int. J. Syst. Evol. Microbiol.* (2003). doi:10.1099/ijs.0.02408-0
- 159 10. Hedegaard, J., Sperling-Petersen, H. U., Nørskov-Lauritsen, N., Steffensen, S. A. d. A. &
160 Mortensen, K. K. Identification of Enterobacteriaceae by partial sequencing of the gene
161 encoding translation initiation factor 2. *Int. J. Syst. Bacteriol.* (2009).
162 doi:10.1099/00207713-49-4-1531

- 163 11. Martinez-Murcia, A. J., Anton, A. I. & Rodriguez-Valera, F. Patterns of sequence variation
164 in two regions of the 16S rRNA multigene family of *Escherichia coli*. *Int. J. Syst. Bacteriol.*
165 (2009). doi:10.1099/00207713-49-2-601
- 166 12. Hashimoto, J. G., Stevenson, B. S. & Schmidt, T. M. Rates and consequences of
167 recombination between rRNA operons. *J. Bacteriol.* (2003). doi:10.1128/JB.185.3.966-
168 972.2003
- 169 13. Ueda, K., Seki, T., Kudo, T., Yoshida, T. & Kataoka, M. Two distinct mechanisms cause
170 heterogeneity of 16S rRNA. *J. Bacteriol.* (1999).
- 171 14. Chaudhury, A., Nath, G., Tikoo, A. & Sanyal, S. C. Enteropathogenicity and antimicrobial
172 susceptibility of new *Escherichia* spp. *J. Diarrhoeal Dis. Res.* (1999).
- 173 15. Podschun, R. & Ullmann, U. *Klebsiella* spp. as nosocomial pathogens: Epidemiology,
174 taxonomy, typing methods, and pathogenicity factors. *Clinical Microbiology Reviews*
175 (1998).
- 176 16. Bastida, F. *et al.* Differential sensitivity of total and active soil microbial communities to
177 drought and forest management. *Glob. Chang. Biol.* **23**, (2017).
- 178 17. Fierer, N. & Jackson, R. B. The diversity and biogeography of soil bacterial communities.
179 *Proc. Natl. Acad. Sci.* (2006). doi:10.1073/pnas.0507535103
- 180 18. Peng, M., Zi, X. & Wang, Q. Bacterial community diversity of oil-contaminated soils
181 assessed by high throughput sequencing of 16s rRNA genes. *Int. J. Environ. Res. Public*
182 *Health* (2015). doi:10.3390/ijerph121012002
- 183 19. Kaiser, K. *et al.* Driving forces of soil bacterial community structure, diversity, and
184 function in temperate grasslands and forests. *Sci. Rep.* (2016). doi:10.1038/srep33696
- 185 20. Zhang, H. H. *et al.* Vertical distribution of bacterial community diversity and water quality
186 during the reservoir thermal stratification. *Int. J. Environ. Res. Public Health* (2015).
187 doi:10.3390/ijerph120606933
- 188 21. Liu, K. *et al.* Bacterial community changes in a glacial-fed Tibetan lake are correlated with
189 glacial melting. *Sci. Total Environ.* (2019). doi:10.1016/j.scitotenv.2018.10.104

- 190 22. Bokulich, N. A. *et al.* Quality-filtering vastly improves diversity estimates from Illumina
191 amplicon sequencing. *Nat. Methods* (2013). doi:10.1038/nmeth.2276
- 192 23. Bokulich. A standardized, extensible framework for optimizing classification improves
193 marker-gene taxonomic assignments. *PeerJ Prepr.* (2015).
194 doi:10.7287/peerj.preprints.49v1
- 195 24. Pruesse, E. *et al.* SILVA: A comprehensive online resource for quality checked and aligned
196 ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Res.* (2007).
197 doi:10.1093/nar/gkm864
- 198 25. Větrovský, T., Baldrian, P. & Morais, D. SEED 2: A user-friendly platform for amplicon
199 high-throughput sequencing data analyses. in *Bioinformatics* (2018).
200 doi:10.1093/bioinformatics/bty071
- 201 26. Cole, J. R. *et al.* Ribosomal Database Project: Data and tools for high throughput rRNA
202 analysis. *Nucleic Acids Res.* (2014). doi:10.1093/nar/gkt1244
- 203 27. Stoddard, S. F., Smith, B. J., Hein, R., Roller, B. R. K. & Schmidt, T. M. rrnDB: Improved
204 tools for interpreting rRNA gene abundance in bacteria and archaea and a new
205 foundation for future development. *Nucleic Acids Res.* (2015). doi:10.1093/nar/gku1201
- 206 28. Venables, W. N. & Ripley, B. D. *Modern Applied Statistics with S Fourth edition.* World
207 (2002). doi:10.2307/2685660
- 208 29. Wickham, H. ggplot2: Elegant Graphics for Data Analysis. *Journeal Stat. Softw.* (2017).
209 doi:10.1007/978-0-387-98141-3
- 210

211 **Acknowledgements**

212 RS thanks the Czech Science Foundation for the project 18-25706S.

213 **Author contributions**

214 RS and DM designed the experiment. RS and DM performed data analysis. The paper was
215 written by RS and DM. All authors approved the final version of the manuscript.

216 **Conflict of Interest**

217 The authors declare no competing financial interests.

218 **Corresponding author**

219 Correspondence and requests for materials should be addressed to
220 robert.starke@biomed.cas.cz.